



Supplementary Information for

RNA Sequencing by Direct Tagmentation of RNA/DNA Hybrids

Lin Di,^{a,#} Yusi Fu,^{a,1,#} Yue Sun,^a Jie Li,^b Lu Liu,^b Jiacheng Yao,^b Guanbo Wang,^{c,d} Yalei Wu,^e Kaiqin Lao,^e Raymond W. Lee,^e Genhua Zheng,^e Jun Xu,^f Juntaek Oh,^f Dong Wang,^f X. Sunney Xie,^{a,d,*} Yanyi Huang,^{a,d,g,*} and Jianbin Wang^{b,*}

^a Beijing Advanced Innovation Center for Genomics (ICG), Biomedical Pioneering Innovation Center (BIOPIC), School of Life Sciences, and Peking-Tsinghua Center for Life Sciences, Peking University, Beijing 100871, China.

^b School of Life Sciences, and Tsinghua-Peking Center for Life Sciences, Tsinghua University, Beijing 100084, China.

^c School of Chemistry and Materials Science, Nanjing Normal University, Nanjing, Jiangsu Province, China

^d Institute for Cell Analysis, Shenzhen Bay Laboratory, Shenzhen, Guangdong Province, China

^e XGen US Co, South San Francisco, CA

^f Department of Cellular and Molecular Medicine, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San Diego, La Jolla, CA 92093

^g College of Engineering, Peking University, Beijing 100871, China

[#] These authors contributed equally to this work.

¹ Present address: Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030, USA

* Corresponding authors: X. Sunney Xie (sunneyxie@pku.edu.cn), Yanyi Huang (yanyi@pku.edu.cn), and Jianbin Wang (jianbinwang@tsinghua.edu.cn).

This PDF file includes:

Supplementary text
Figures S1 to S11
SI References

Materials and Methods

Purification of pTXB1 Tn5 and D188E mutation

The pTXB1 cloning vector, which introduced hyperactive E54K and L372P mutation into wildtype Tn5, was acquired from Addgene. The pTXB1 Tn5 and its mutant were expressed and purified mainly according to the protocol published by Picelli S et al. [1] To construct D188E mutation into Tn5, pTXB1 vector was firstly amplified into two parts by two sets of primers. Mutagenesis primers used for the first part (3771-7979) which contained site 188 were 5'-GGCAGCATGATGAGCAACGTGATTGCGGTGTGCGAACG TGAAGCGGATATTCATGC-3' and 5'-TATCAGCTCACTCAAAGG-3'. Amplified primers for the remaining part were 5'-GTATTACCGCCTTTGAGT-3' and 5'-CAATCACGTTGCTCATCA-3'. The purified PCR products were then assembled into intact plasmid using Gibson Assembly Master Mix (NEB, Cat.No. E2611). The newly assembled plasmid was transformed into *E. coli* Trans5α chemically competent cells (Transgene, CD201-01). After growing overnight on LB medium plate, single colony was picked and shaken in SOC liquid medium for at least 9 hours. Plasmid was extracted by PurePlasmid Mini Kit (CW BIO, Cat CW0500S) and confirmed carrying D188E mutation by Sanger sequencing. Then the plasmid was transformed into *E. coli* Transetta (DE3) chemically competent cells (Transgene, CD801-01) for further protein expression and purification.

Cell culture

HEK293T and HeLa cell lines are acquired from ATCC. Both of them were cultured in Dulbecco's Modified Eagle Medium (Gibco, Cat 11965092), supplemented with 10% fetal bovine serum (Gibco, Cat 1600044) and 1% penicillin-streptomycin (Gibco, Cat 15140122). The cell incubator (Thermo Scientific) was set at temperature of 37°C with 5% CO₂ injected. Adherent cells were washed twice by DPBS (Gibco, Cat 14190136) and detached by 0.05% Trypsin-EDTA (Gibco, Cat 25300062) at 37°C for 4min. Then double volume of culture media was added to terminate trypsinization. Cells were collected by centrifugation at 200g for 5min and resuspended for downstream experiment or passage cultivation.

Nucleic acids extraction and messenger RNA isolation

Genomic DNA was extracted using PureLink Genomic DNA Mini Kit (Invitrogen, Cat K182002), and total RNA was extracted using RNeasy Mini Kit (Qiagen, Cat.No.74104). The resulting total RNA was then reacted with 10ul DNase I (NEB, Cat.No.M0303) to remove remaining DNA thoroughly, and concentrated by RNA Clean & Concentrator-5 kit (Zymo Research, Cat R1015). The quality of extracted DNA and RNA was assessed by the Fragment Analyzer Automated CE System (AATI) and quantification was done by Qubit 2.0 (Invitrogen, Cat Q33230/ Q32852). We followed standard protocol of NEBNext Poly(A) mRNA Magnetic Isolation Module (NEB, Cat E7490) to isolate messenger RNA from the purified total RNA and stored them at -80°C.

Single cell preparation

We used pipette tips to form some drops made up of PBS (containing 1% BSA) (Thermo Scientific, Cat 37525) on a clean petri dish. The cell resuspension was pipetted up and down gently to disperse into single cells and we took ~5μl of them diluted in one of the drops. The mouth pipette with 50μm inside diameter was then used to pick one cell in the drop and release it in another clean drop. The picked cell was passed by at least three clean drops in order to wash away any debris and confirm that only one cell was in the last drop. We then aspirated the cell

with as little buffer as possible and blew it into 4µl lysis buffer [4 units of Recombinant RNase Inhibitor (Takara, Cat.No.2313), 2.5µM poly(T)30VN primer (Sangon), 2.5mM dNTP (NEB, Cat.No.N0447) and 0.48% Triton X-100 (Sigma, Cat.No.T9284)]. Successful transfer was confirmed by blowing mouth pipette again in a clean drop and no cell was to be seen in visual field. Reaction was carried out at 72°C for 3min after violent vortex.

mRNA/DNA hybrid formation

Total RNA was reverse transcribed into mRNA/DNA hybrid mainly referring to Smart-seq2 protocol [2], but with several modifications: 1) The ISPCR part in Oligo-dT primer was removed; 2) TSO was omitted, but in TSO-RT SHERRY, it should be kept; 3) The reaction was performed at 42°C for 1.5h without cycling. If input was purified RNA, Triton X-100 was omitted. When inputting more than 10ng total RNA, we would slightly upregulate amount of dNTP, poly(T)30VN primer and Superscript II (Invitrogen, Cat.No.18064014).

Tn5 transposome in vitro assembly and tagmentation

Functional mosaic-end (ME) oligonucleotides (5'-CTGTCTCTTATACACATCT-3', 100µM) was separately annealed with equal amounts of Adaptor A (5'-TCGTGGCAGCGTCAG ATGTGTATAAGAGACAG-3', 100µM) and Adaptor B (5'-GTCTCGTGGGCTCGGAGATG TGTATAAGAGACAG-3', 100µM). The concentration of purified Tn5 was quantified by Qubit Protein Assay Kit (Invitrogen, Cat.No.Q33212) and took around 100µg for transposome assembly. We mixed the Tn5 transposase with annealed ME-Adaptor A or B (20µM) in 45% glycerol (Sigma, Cat.No.G5516) thoroughly and incubated the mixture at 30°C for one hour. These two resulting transposomes (assembled with ME-Adaptor A/B) were then mixed together, ready for tagmentation or stored at -20°C. Specifically, to assemble rCrArG Tn5 in Fig.S4, ribonucleotide modifications were made on the three terminal bases at 3'-end of Adaptor A/B. And for rG Tn5, the last base at 3'-end of adaptors was modified.

The dsDNA tagmentation was performed in 1xTD buffer [10mM Tris-Cl (pH 7.6, ROCKLAND, Cat.No.MB-003), 5mM MgCl₂ (Invitrogen, Cat.No.AM9530G), 10% N,N-Dimethylformamide (Sigma, Cat.No.D4551)]. The reaction was incubated at 55°C for 30min.

As for RNA/DNA hybrid, tagmentation was performed in buffer containing 10mM Tris-Cl (pH 7.6), 5mM MgCl₂, 10% N,N-Dimethylformamide, 9% PEG8000 (VWR Life Science, Cat.No.97061), 0.85mM ATP (NEB, Cat.No.P0756). In SHERRY library preparation, the corresponding amount of pTXB1 Tn5 transposome used for different initial input of total RNA was list as below:

| Amount of total RNA input | Amount of Tn5 transposome |
|---------------------------|---------------------------|
| Single cell (~10 pg) | 0.003 µl |
| 100 pg | 0.003 µl |
| 10 ng | 0.006 µl |
| 200 ng | 0.050 µl |

The transposome could be diluted in 1xTn5 dialysis buffer [50mM Hepes (pH 7.2, Leagene, Cat.No.CC064), 0.1M NaCl (Invitrogen, Cat.No. AM9759), 0.1 mM EDTA (Invitrogen, Cat.No.AM9260G), 1 mM DTT, 0.1% Triton X-100, 10% glycerol].The reaction was incubated at 55°C for 30min.

Commercial Tn5 transposomes were available in Nextera XT DNA Library Prep Kit (Illumina, Cat.No.FC-131-1024) and TruePrep DNA Library Prep Kit V2 for Illumina (Vazyme, Cat.No.TD501).

SHERRY library preparation and sequencing

To construct scSHERRY library, the single cell tagmentation product was mixed well with 4 units of Bst 3.0 DNA Polymerase (NEB, Cat.No.M0374) and indexed common primers (Vazyme, Cat.No.TD202) in 1 x Q5 High-Fidelity Master Mix (NEB, Cat.No.M0492). Then index PCR was performed as follow: 72°C 15min, 98°C 30s, 10 cycles of [98°C 20s, 60°C 20s, 72°C 2min], 72°C 5min. The PCR product was purified with 0.85:1 ratio by VAHTS DNA Clean Beads (Vazyme, Cat.No.N411) and eluted in 30µl nuclease-free water (Invitrogen, Cat.No.AM9937) for another 18

cycles of PCR. When performing high-throughput experiment, each sample could be amplified by 15 cycles, then merged for beads purification and library quality check.

As for purified 10ng or 200ng total RNA input, the tagmentation product was firstly gap-filled with 100 units of Superscript II and 1 x Q5 High-Fidelity Master Mix at 42°C for 15min, then Superscript II was inactivated at 70°C for 15min. When inputting 100pg total RNA, the extension enzyme was replaced with 4 units of Bst 2.0 Warmstart DNA Polymerase (NEB, Cat.No.M0538). Correspondingly, the reaction temperature was upregulated to 72°C and inactivation was performed at 80°C for 20min. After that, indexed common primers were added to perform PCR. PCR cycles were listed as below:

| Amount of total RNA input | Index PCR cycles |
|---------------------------|------------------|
| 100 pg | 25 |
| 10 ng | 15 |
| 200 ng | 12 |

The resulting library was purified with 1:1 ratio by VAHTS DNA Clean Beads. Quantification was done by Qubit 2.0 and quality check was done by Fragment Analyzer Automated CE System. The sequencing platform we used was Illumina NextSeq 500 or HiSeq 4000.

NEBNext and SmartSeq2 library preparation

NEBNext RNA-Seq library preparation starting from 10ng and 200ng total RNA was performed using NEBNext Ultra II RNA Library Prep Kit for Illumina (NEB, Cat.No.E7770). Single cell SmartSeq2 library was constructed as previously reported [2].

Ligation Test and Strand Test

In Ligation Tests, tagmentation products from 200ng HEK293T total RNA were purified by DNA Clean & Concentrator-5 (Zymo Research, Cat.No.D4013) and eluted with 20µl nuclease-free water. After gap-filling with Superscript II, Ligation Test 1 was processed directly to index PCR while product in Ligation Test 2 was digested by 12.5 units of RNase H (NEB, Cat.No.M0297) at 37°C for 20min before index PCR.

In Strand Tests, we used dUTP (Thermo Scientific, Cat.No.R0133), dATP, dCTP, dGTP (NEB, Cat.No.N0446) mix, each of them at equal concentration, to incorporate in cDNA during reverse transcription step. 0.15µl Tn5 transposome was used to tagment the resulting hybrid. Fragments were then column-purified and gap-filled by Bst 2.0 Warmstart DNA Polymerase, and column purification was again applied. For Strand Test 1, 3 units of USER enzyme (NEB, Cat.No.M5505) and 40units of recombinant rnase inhibitor was added into elution products and incubated at 37°C for 20min for DNA strand digestion. Indexed common primers added with digestion product in 1 x Q5 High-Fidelity Master Mix were then reacted at 85°C for 30s, followed by 60°C for 2min and temperature went down to 4°C slowly. After that, reverse transcription was performed with 200 units of Superscript II added at 42°C for 30min, then transferred to index PCR program. For Strand Test 2, the USER digestion product was directly performed index PCR with 1 x KAPA HiFi HotStart Uracil+ ReadyMix (Kapa Biosystems, Cat.No.KK2801). Protocol of Strand Test 3 was almost same as Strand Test 2, except replacing USER enzyme with RNase H. For dU-SHERRY, the USER enzyme digestion was omitted compared with Strand Test 2 workflow.

Docking model

To generate the substrate-transposon DNA-Tn5 structure model, a 32bp dsDNA or DNA/RNA hybrid were generated by 3D-NuS (3-Dimensional Nucleic Acid Structures) web server. Then the substrate dsDNA or DNA/RNA was manually docked to the transposon DNA-Tn5 structure, PDB ID: 1MUS, based on charge and shape complimentary.

Data analysis

Sequencing adaptors or poly(T/A) positioned at end of paired reads were recognized and removed by Cutadapt v1.15 [3]. The trimmed reads which length was shorter than 20bp were filtered. Remaining reads were down sampled to 2 million (except that library with 200ng total

RNA input used for differential gene expression analysis was 10 million) total reads, and aligned with index built from human(hg38) genome and known transcript annotations by Tophat2 v2.1.1 [4]. The mapped reads were then used to calculate FPKM value for each known gene (annotation acquired from UCSC) by Cufflinks v2.2.1 with multi-mapped reads correction. Gene with FPKM more than 1 was considered to be detected. The exonic rate, duplicate rate and insert size of library were all calculated by Picard Tools v2.17.6.

General coverage across known transcripts was plotted by RSeQC v.2.6.4 [5]. For specific transcript, depth of mapped reads overlapped with transcript position was calculated by Samtools v1.3.1.

We used DESeq2 v1.22.2 [6] to perform differential gene expression analysis with raw count-matrix acquired by HTSeq v.0.11.0 [7]. Differentially expressed genes should meet following criteria: 1) FPKM value >1; 2) significant p-value <5x10⁻⁶; 3) absolute value of log₂(Fold Change) >1. Counts in correlation plot were normalized mainly according to DESeq2 normalization method [6], which considered library size and library compensation. Correlation efficient R² and slope of linear fitting equations were calculated by least square method. The slope-skewness between two replicates was defined as |k₁-1|+|k₂-1|, k₁ or k₂ was slope when one of the replicates was conducted as X or Y axis.

GC content distribution of detected genes was plotted by custom Perl script. GC content was binned by 2% and gene number at each bin was normalized by dividing maximum gene number of one bin. And for each gene, only the longest transcript isoform was calculated.

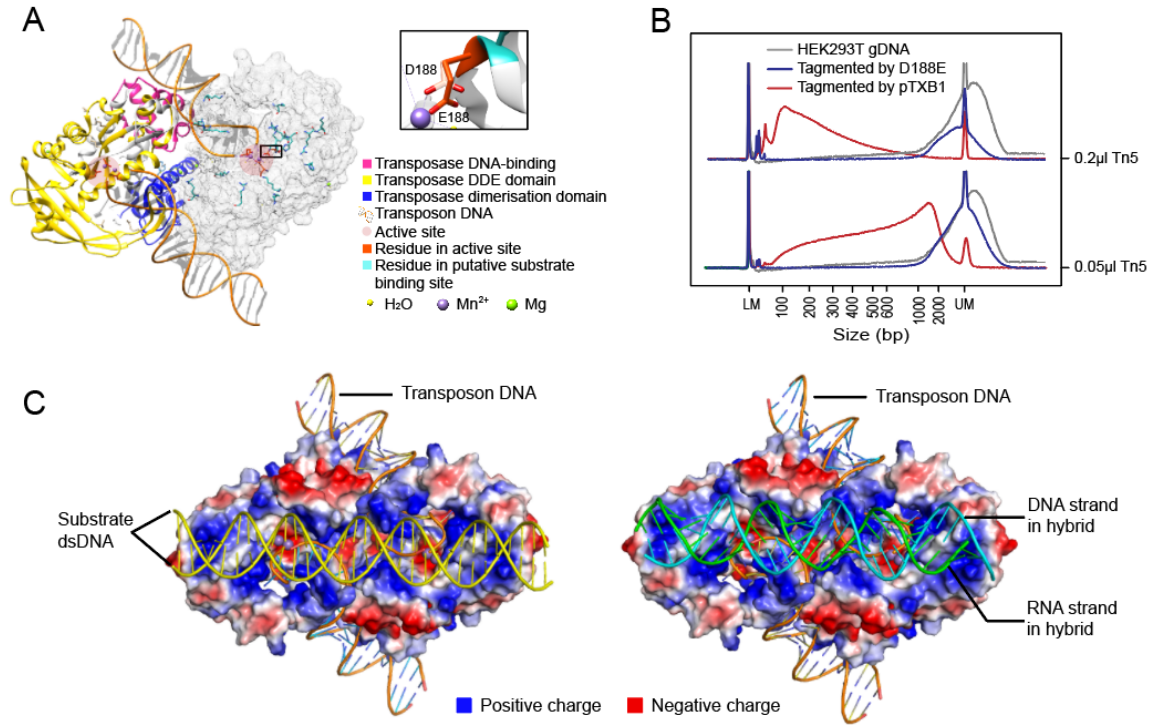


Fig. S1. Structure of Tn5 and D188E mutation in Tn5. **(A)** Structure of pTXB1 Tn5 (PDB ID: 1MUS). Left monomer marked domains in different colors. Right monomer marked residues of catalytic core and putative substrate binding site in atom form. [Referred to D. R. Davies, I. Y. Goryshin, W. S. Reznikoff, I. Rayment, Three-dimensional structure of the Tn5 synaptic complex transposition intermediate. *Science* 289, 77-85 (2000).] Black box showed D188E mutation. **(B)** Size distribution of genomic DNA with no treatment or tagmented by D188E mutant Tn5 or tagmented by pTXB1 Tn5. **(C)** Model of docking double-stranded DNA (Left) or RNA/DNA heteroduplex (Right) in predicted substrate binding sites of Tn5.

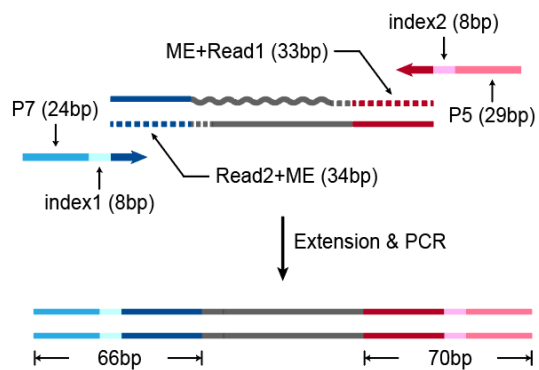


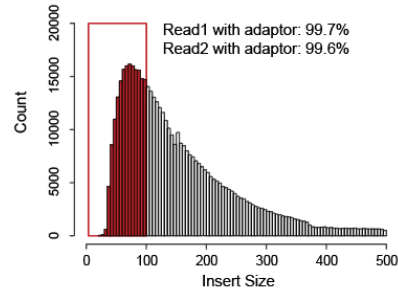
Fig. S2. Composition of products amplified from tagged RNA/DNA hybrid. Gray wavy line and straight line represent RNA and DNA separately. Dotted lines represent the track of extension step.

A

| | Mapping Rate | Exonic Rate | Gene Detected |
|-----------------|--------------|-------------|---------------|
| Ligation Test 1 | 92.20±0.20% | 77.76±0.18% | 11,825±7 |
| Ligation Test 2 | 93.30±0.00% | 78.98±0.10% | 11,784±1 |
| SHERRY | 89.77±0.05% | 80.32±0.04% | 11,784±8 |
| Strand Test 1 | 11.25±0.05% | 54.92±0.06% | 7,384±13 |
| Strand Test 2 | 4.05±0.15% | 61.21±0.32% | 5,566±81 |
| Strand Test 3 | 78.20±0.50% | 77.38±0.10% | 11,470±14 |
| dU-SHERRY | 90.90±0.10% | 76.31±0.00% | 11,335±40 |

B

Ligation Test 1 Insert Size Histogram



C

Ligation Test 2 Insert Size Histogram

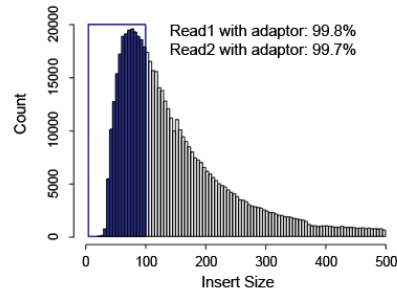


Fig. S3. Comparison within Ligation Tests and Strand Tests. **(A)** Sequencing indicators of Ligation Tests and Strand Tests. Each test consisted of two replicates of 200 ng HEK293T total RNA. **(B-C)** Insert size distribution of Ligation Tests. The colored bars marked reads which insert size is shorter than 100bp. Adaptors detected in these reads are counted.

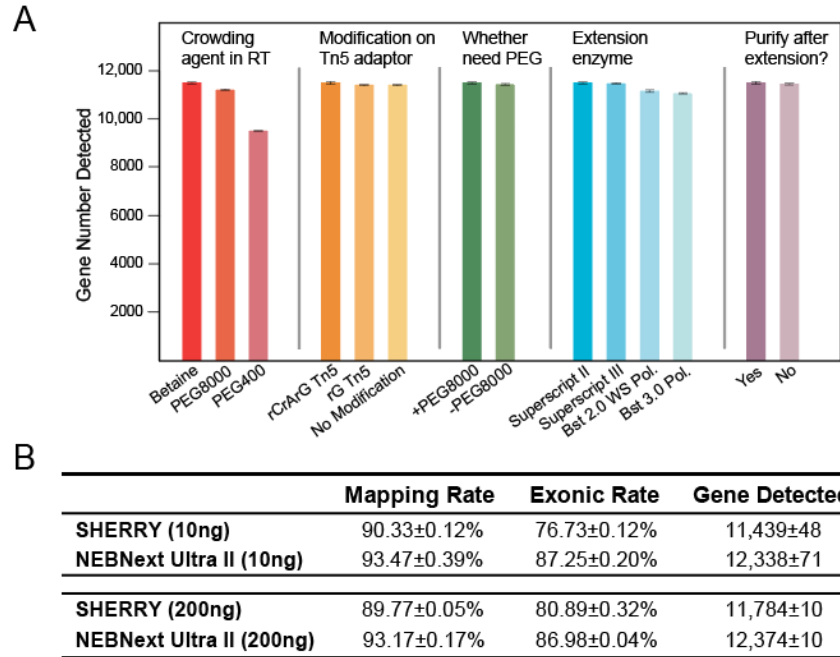


Fig. S4. Optimization of SHERRY and comparison with NEBNext. **(A)** Gene number detected by SHERRY under various experimental conditions. Each condition consisted of three replicates of 10 ng HEK293T total RNA. **(B)** Comparison of sequencing indicators between SHERRY and NEBNext with 10 ng and 200 ng HEK293T total RNA input. Each condition consisted of three replicates and down-sampled to 2 million total reads.

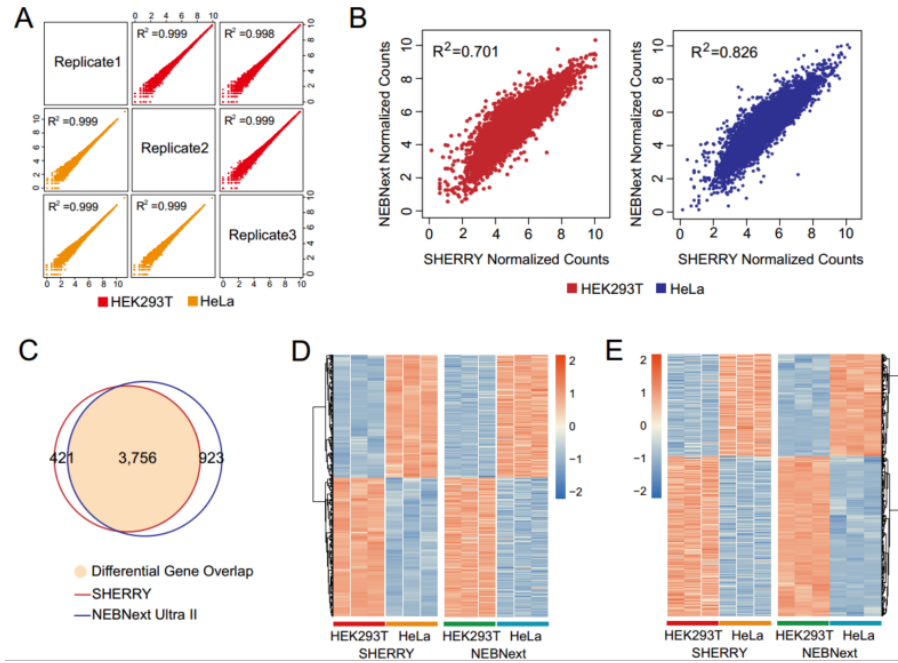


Fig. S5. Functional comparison between SHERRY and NEBNext. **(A)** Correlation of normalized gene counts among duplicates of SHERRY, which start from 200 ng HEK293T total RNA input. **(B)** Correlation of normalized genes counts (average of three replicates) between SHERRY and NEBNext within the two cell types. The input was 200 ng total RNA. **(C)** Differentially expressed genes of HeLa and HEK293T detected by SHERRY and NEBNext kit (200 ng input) are plotted into Venn Diagram. Colored area represents genes identified by both methods. Gene numbers are listed on corresponding part. **(D)** Heatmap of differentially expressed genes detected by SHERRY while missed by NEBNext kit. The Color bar indicates Z-score. **(E)** Heatmap of differentially expressed genes detected by NEBNext kit while missed by SHERRY. The Color bar indicates Z-score.

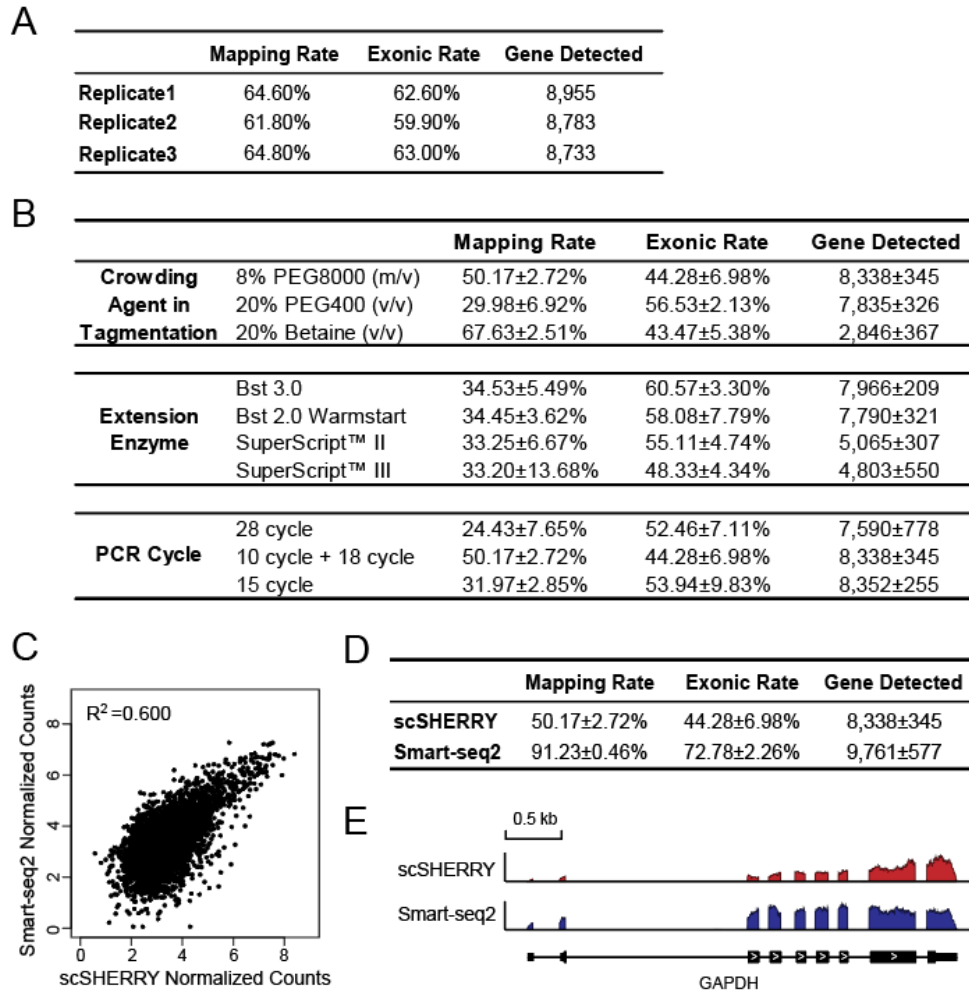
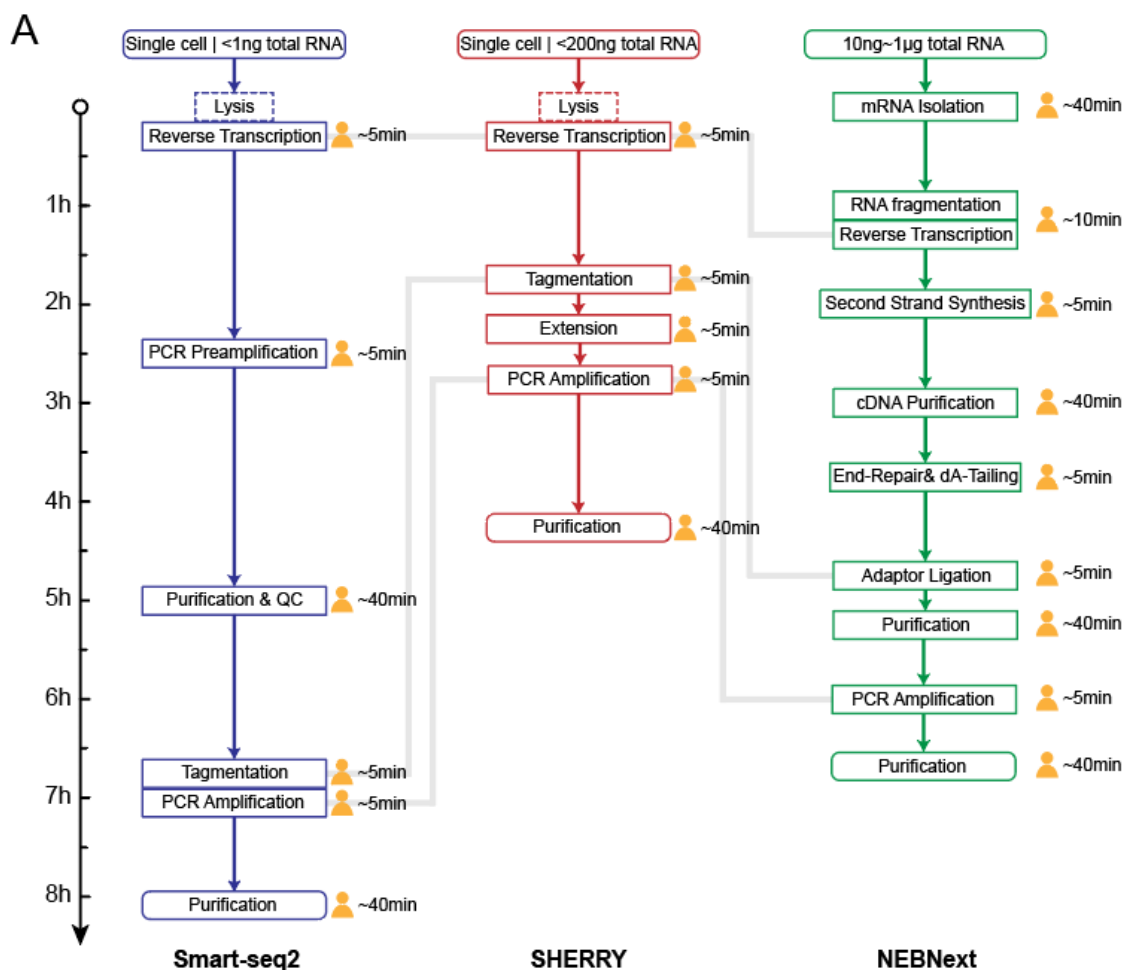


Fig. S6. Optimization of micro-input SHERRY and comparison with Smart-seq2. **(A)** Sequencing indicators of SHERRY (n=3) starting from 100 pg HEK293T total RNA input. **(B)** Comparison of scSHERRY library quality under various experiment conditions. Each condition used single HEK293T cell as input and consisted of 3-4 replicates. **(C)** Correlation of normalized gene counts (average of three replicates) between scSHERRY and Smart-seq2. **(D)** Comparison of sequencing indicators between scSHERRY (n=3) and Smart-seq2 (n=4). Both of them used single HEK293T cells as input. **(E)** The coverage of GAPDH transcript calculated from scSHERRY and Smart-seq2.



B

| | Smart-seq2 | SHERRY | | NEBNext |
|-------------------|------------|----------|----------------|----------|
| Lysis+RT | ~\$4.04 | ~\$3.94 | mRNA Isolation | ~\$2.79 |
| Pre-amplification | ~\$2.09 | - | - | - |
| Library Prep | ~\$42.10 | ~\$6.46 | Library Prep | ~\$45.42 |
| Total | ~\$48.23 | ~\$10.40 | Total | ~\$48.21 |

Fig. S7. Comparison of workflow and cost among Smart-seq2, SHERRY and NEBNext kit. **(A)** Workflow of Smart-seq2, SHERRY and NEBNext kit. Length of arrow indicates time consumed for each step. The human-shaped icon indicated hands-on time. Dotted box means this step is alternative. Gray line connects corresponding key step in each method. **(B)** Cost list of Smart-seq2, SHERRY and NEBNext kit.

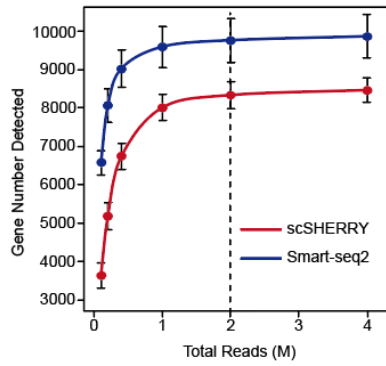


Fig. S8. Saturation curve of scSHERRY (n=3) and Smart-seq2 (n=4). Both methods used single HEK293T cell as input. Dotted line indicated 2 million total reads.

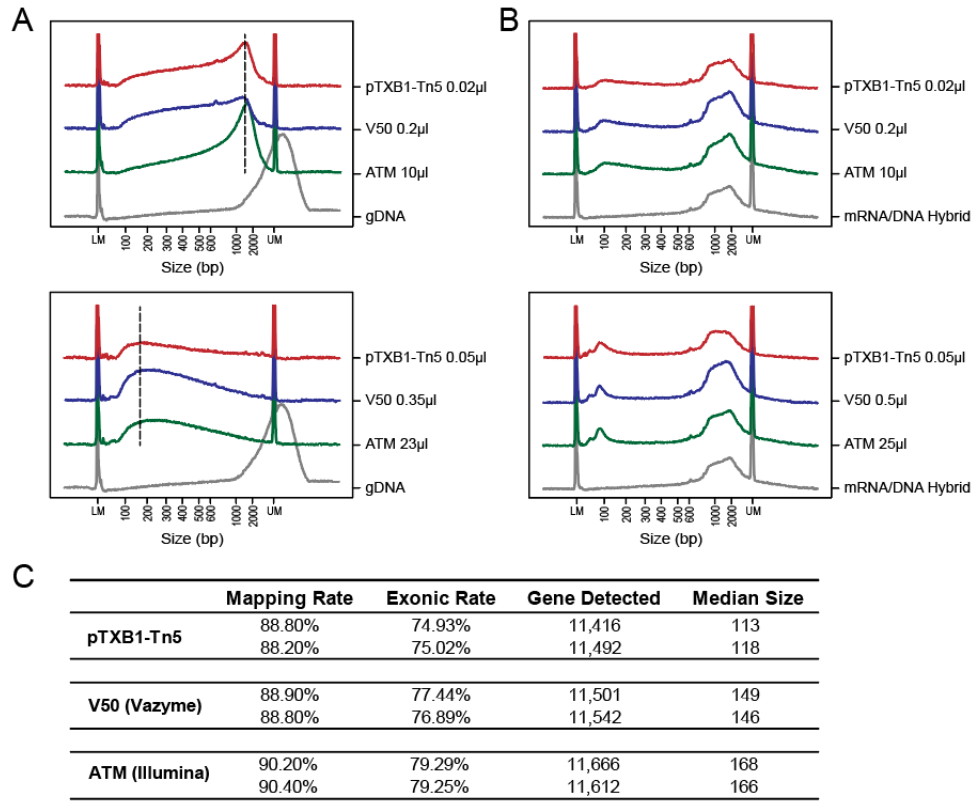


Fig. S9. Hybrid tagmentation activity of commercial Tn5. **(A)** Size distribution of genomic DNA with no treatment or tagmented by different volumes of pTXB1 Tn5, V50 or ATM. The dotted black line indicates peak of fragment size. **(B)** Size distribution of mRNA/DNA hybrid with no treatment or tagmented by different volumes of pTXB1 Tn5, V50 or ATM. **(C)** Sequencing indicators of SHERRY library constructed by three Tn5.

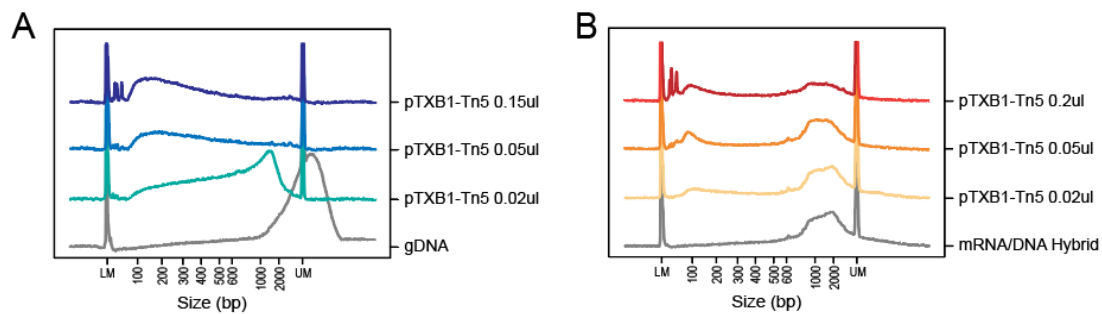


Fig. S10. Titration of Tn5 transposome for tagmentation. **(A-B)** Size distribution of 5ng genomic DNA or mRNA/DNA hybrid with no treatment or tagmented by different gradients of pTXB1 Tn5.

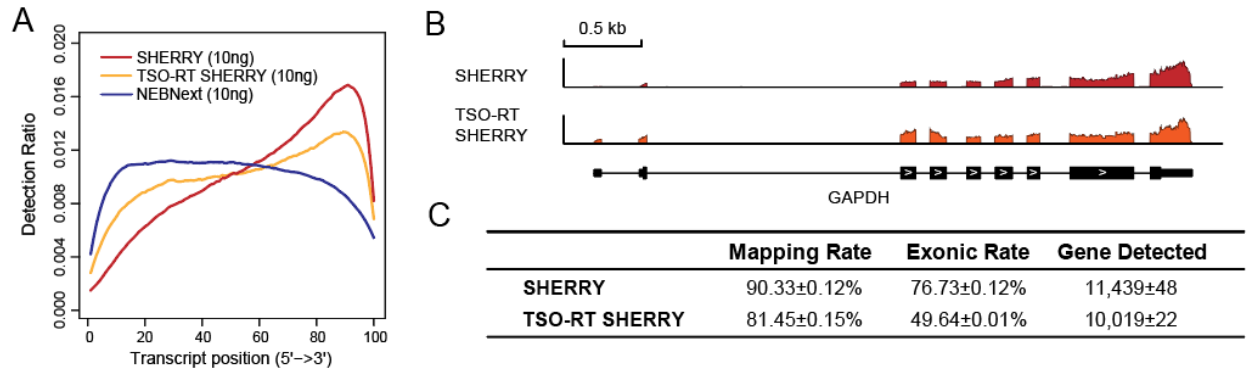


Fig. S11. Coverage evenness optimization of SHERRY (10 ng HEK293T total RNA input). **(A)** Normalized transcript coverage of standard SHERRY, SHERRY using TSO-RT method and NEBNext kit. **(B)** The coverage of GAPDH transcript calculated from SHERRY and TSO-RT SHERRY. **(C)** Comparison of sequencing indicators between SHERRY (n=3) and TSO-RT SHERRY (n=2).

References

- [1] S. Picelli, A. K. Björklund, B. Reinius, S. Sagasser, G. Winberg, R. Sandberg, Tn5 transposase and tagmentation procedures for massively scaled sequencing projects. *Genome Res.* 24, 2033-2040 (2014).
- [2] S. Picelli, O. R. Faridani, A. K. Björklund, G. Winberg, S. Sagasser, R. Sandberg, Full-length RNA-seq from single cells using Smart-seq2. *Nat. Protoc.* 9, 171-181 (2014).
- [3] M. Martin, Cutadapt removes adapter sequences from high-throughput sequencing reads. *Embnet Journal* 17, 200 (2011).
- [4] D. Kim, G. Pertea, C. Trapnell, H. Pimentel, R. Kelley, S. L. Salzberg, TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol.* 14, R36 (2013).
- [5] L. Wang, S. Wang, W. Li, RSeQC: quality control of RNA-seq experiments. *Bioinformatics* 28, 2184-2185 (2012).
- [6] M. I. Love, W. Huber, S. Anders, Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* 15, 550 (2014).
- [7] S. Anders, P. T. Pyl, W. Huber, HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* 31, 166-169 (2015).