# Comparison of dsDNA and ssDNA-based NGS library construction methods for targeted genome and methylation profiling of cfDNA

Jianchao Zheng [1,2], Zhilong Li[1,2], Xiuqing Zhang[1,3], Hongyun Zhang[2], Shida Zhu[2,4],

Jianlong Sun[2†], and Yuying Wang[2†]


[1]College of Life Sciences, University of Chinese Academy of Sciences, Beijing

100049，China.

[2]BGI Genomics, BGI-Shenzhen, Shenzhen 518083, China.

[3]BGI-Shenzhen, Shenzhen 518083, China.

[4]Shenzhen Engineering Laboratory for Innovative Molecular Diagnostics, BGI-Shenzhen, Shenzhen, 518120, China


[†]Correspondence

Dr. Jianlong Sun and Dr. Yuying Wang, BGI Genomics, BGI-Shenzhen, Shenzhen 518083, China. E-mail: sunjianlong@genomics.cn and wangyuying@genomics.cn

**Conflict of Interest Statement**

JZ, ZL, HZ, SZ, JS, and YW are employees of BGI Genomics, BGI-Shenzhen. XZ is an employee of BGI-Shenzhen.

**Abstract**

Cell-free DNA (cfDNA) profiling by next generation sequencing (NGS) has wide applications in cancer diagnosis, prognosis, and therapy response monitoring. One key step of cfDNA deep sequencing workflow is NGS library construction, whose efficiency determines effective sequencing depth, sequencing quality, and accuracy. In this study, we compared two different cfDNA library construction methods for the applications of mutation detection and methylation profiling: the conventional method which captures double-stranded DNA (dsDNA) molecules, namely the dsLib workflow, and an alternative method which captures single-stranded DNA (ssDNA), namely the ssLib workflow. Our results suggest that the dsLib method was preferrable for mutation detection while the ssLib method proved more efficient for methylation analysis. Our findings could help researchers choose more appropriate library construction method for corresponding downstream sequencing applications.

**Keywords**: cfDNA, NGS, library construction methods, target sequencing, methylation.

1 **Introduction**

2 Cell-free DNA (cfDNA), primarily derived from cell apoptosis, has been shown to be

3 an important biomarker of many physiological and pathological conditions such as

4 autoimmunity, infection, pregnancy, exercise, transplantation, and cancer [1-3]. It is

5 detectable in almost all body fluids including plasma, serum, and urine, with a peak

6 size at approximately 166 bp [4, 5]. In cancer patients, there is a subset of cfDNA

7 known as circulating tumor DNA (ctDNA), which originates from tumor cells and

8 carries genetic and epigenetic characteristics of the tumor [6, 7]. cfDNA has a short

9 half-life in circulation (reported to be between 16 and 150 minutes), can be repeatedly

10 sampled, and may potentially overcome intratumor heterogeneity compared to tissue

11 biopsies [8]. These unique characteristics make cfDNA-based liquid biopsy an ideal

12 approach for cancer diagnosis, prognosis, and therapy response monitoring [9, 10].

13 However, due to low concentration of cfDNA in plasma (~3 ng/ml in healthy

14 individuals) and very small fraction of ctDNA among abundant cfDNA that derived

15 from blood cells and normal tissues, accurate detection of ctDNA remains a

16 challenging task [11, 12].

17

18 Massively parallel sequencing (MPS), also known as next generation sequencing

19 (NGS), has been widely applied in both research and diagnostic fields [13, 14].

20 Millions of DNA fragments can be simultaneously sequenced and analyzed by NGS

21 [15]. Furthermore, targeted capture sequencing allows for deeper sequencing for

22 target regions of interest at a lower cost [16]. Remarkably, efficient library

1 construction before targeted capture sequencing determines effective sequencing

2 depth and remains indispensable to successful sequencing of the target regions, and is

3 particularly critical for sequencing of limited amount of cfDNA, and for identification

4 of variants with lower allele fractions [9, 17]. During library construction,

5 platform-specific adapters, which contain sample barcode sequence(s) and common

6 primer binding sites for subsequent amplification and sequencing, are ligated to both

7 ends of the original DNA fragments [18]. Various library construction methods have

8 been introduced, aiming to improve DNA conversion efficiency (defined as the

9 fraction of original DNA molecules that are successfully converted to the final library)

10 [19-23]. Conventional double-stranded library (dsLib) construction workflow (such as

11 what is used in the KAPA Hyper Prep kit) consists of following steps: (i) end repair

12 and dA-tailing of the double-stranded (dsDNA) templates; (ii) adapter ligation; (iii)

13 library amplification and purification [24]. On the other hand, single-stranded library

14 (ssLib) construction was usually initialized by adapter ligation to single-stranded

15 DNA (ssDNA) templates and followed by library amplification and purification [19,

16 22, 25]. The ssLib construction method was originally developed to recover ancient

17 and/or degraded DNA fragments which are usually poorly captured by conventional

18 dsDNA-based library preparation [26, 27]. Previously, researchers have compared

19 their performance in applications such as non-invasive prenatal testing (NIPT), which

20 is based on shallow whole-genome sequencing (WGS), and found no advantage for

21 ssDNA-based methodology [28]. However, there hasn't been systematic study to

22 compare the performance of these two methods when used for cfDNA sequencing for

1    cancer-related applications.

2

3    In this study, we compared the dsLib workflow and ssLib workflow for targeted deep

4    sequencing (for variant detection) and methylation sequencing (for detection of

5    cytosine methylation, an important form of epigenomic modification) of cfDNA, two

6    applications important for cancer diagnosis. We found that, for targeted deep

7    sequencing, the dsLib method achieved overall better performance and satisfactory

8    limit of detection (LOD). For methylation sequencing, we compared the dsLib and

9    ssLib workflow coupled with either bisulfite-based or enzyme-based cytosine

10   conversion methods, and found that ssLib coupled with bisulfite conversion showed

11   notably better performance.

12

13   **Materials and Methods**

14   **Ethical Compliance**

15   This study was approved by the institutional review board of BGI (NO. BGI-IRB:

16   19077).

17

18   **Sample collection and cfDNA isolation**

19   After obtaining informed consent, blood samples were collected from 37 healthy

20   volunteers and 2 lung cancer patients in 10 mL K2 EDTA BD Vacutainer tubes. Blood

21   was separated immediately by an initial centrifugation at $1,600 \times g$ for 10 min and then

22   by a second centrifugation at $16,000 \times g$ for 10 min. Plasma were pooled and split into

1    4ml per reaction for cfDNA isolation using MagPure Circulating DNA Maxi Kit

2    (Magen, China) per manufacturer's instruction. Extracted cfDNA samples from

3    healthy volunteers were pooled together to obtain sufficient homogeneous material for

4    subsequent analysis. cfDNA was quantitated by Qubit dsDNA High Sensitivity Kit

5    (Thermo Fisher Scientific, USA). The 1% Multiplex I cfDNA reference standards

6    HD778 (Horizon Discovery, UK) were spiked into healthy donor cfDNA at 0.1%,

7    0.25%, or 0.5% to simulate cfDNA samples with defined mutant allele frequencies

8    (MAFs) . Experiments were performed in triplicates.

9

10   **Double-stranded cfDNA library construction**

11   Duplex unique molecular identifier (UMI) adapters for MGISEQ-2000 sequencer

12   were designed according to principles described by Newman et al [21] with the

13   modification that 3-bp UMIs were chosen instead of 2-bp UMIs in order to

14   accommodate a higher library complexity. To avoid potential issues during

15   sequencing caused by low complexity at the T-A ligation position (constant base), 32

16   pairs of UMI adapters were incorporated with an additional base (G or C) before the

17   T-A    ligation    position.    Long    oligonucleotides    UMIxxL    (5'-

18   Phosphorylation-[C/G/-]-NNNAAGTCGGAGGCCAAGCGGTCTTAGGAAGACAA

19   -3')        and        short        oligonucleotides        UMIxxS

20   (5'-GACATGGCTACGATCCGACTNNN-[G/C/-]-T-3') were synthesized by BGI

21   tech solutions (Beijing Liuhe co.limited). Each oligo was dissolved to 100 μM using

22   TE buffer. For each pair of adapters, 5 μL UMIxxL and 5 μL UMIxxS oligos (100 μM)

1    were combined and brought up to 20 µL with TE buffer. Oligos were annealed for

2    more than 30 minutes at room temperature. 64 UMI adapters (25 µM) were mixed and

3    diluted to 5 µM, marked as UMI64M.

4

5    Double-stranded cfDNA libraries were prepared either by KAPA Hyper Prep kit

6    (Kapa Biosystems, cat. No. KK8504) per manufacturer's instruction or our custom

7    library construction protocol. For the latter, briefly, 1-10 ng cfDNA was mixed with

8    end-repair master mix consisting of T4 DNA polymerase (Enzymatics, cat. No.

9    P7080L), T4 polynucleotide kinase (Enzymatics, cat. No. Y9040L), rtaq DNA

10   polymerase (MGI, cat. No. 01E012MM), dNTP, and T4 DNA ligase buffer, and kept

11   at 20□ for 30 min followed by 65□ for 30 min. Then UMI64M adapter was added to

12   the end-repair reaction product and mixed by pipetting, followed by adding ligation

13   master mix consisting of golden T4 DNA ligase (MGI, cat. No. 02E004MM), 10× T4

14   DNA ligase buffer, and PEG6000 (Sigma Aldrich, 50%). The ligation reaction was

15   incubated at 16□ for 60 min. Adapter ligated DNA was purified using Agencourt

16   AMPure XP beads. Next, index PCR was then performed and purified using

17   Agencourt AMPure XP beads. The concentration of final library was determined by

18   Qubit dsDNA High Sensitivity Kit.

19

20   Double-stranded cfDNA methylation sequencing libraries were prepared according to

21   above library preparation workflow with following modifications: (i) 0.05 ng

22   fragmented lambda DNA was spiked into the 10 ng cfDNA to monitor bisulfite

conversion rate; (ii) all cytosines of adapter were methylated; (iii) after purification of the ligation product, bisulfite conversion was performed using EZ DNA Methylation Gold kit (Zymo Research, cat. No. D5006) or EM-seq Conversion Module (NEB, cat. No. E7125); (iv) index PCR was performed by 2×Golden U+ High-fidelity Readymix (MGI, cat. No. 01K01701MM).

**Single-stranded cfDNA library construction**

The single-stranded library preparation method was based on the ssDNA2.0 method [19] with the modification that T-A ligation was used to further improve ligation efficiency. Briefly, MyOne C1 beads carrying the extension product were resuspended in the A-tailing reaction mix consisting of Klenow (3'-5' exo-) (Enzymatics, cat. No. P7010-LC-L), 10× blue buffer, and dATP, and incubated at 37℃ for 30 min then at 75℃ for another 30min. The libraries were amplified by a specific number of PCR cycles based on cfDNA input amount, purified by Agencourt AMPure XP beads, and eluted in nuclease-free water.

Single-stranded cfDNA methylation sequencing libraries were prepared as above after input cfDNA was converted using either the EZ DNA Methylation Gold kit (Zymo Research, cat. No. D5006) or the EM-seq Conversion Module (NEB, cat. No. E7125). To monitor the conversion rate, 0.05 ng fragmented lambda DNA was spiked into 10 ng cfDNA.

**Target capture and sequencing**

A custom capture panel that spans 220 kb and covers 139 cancer driver genes was designed and synthesized by IDT technologies as previously described [29]. Targeted genome capture was performed using xGen® Lockdown® Reagents (IDT technologies) and BGI adapter-specific blockers (BGI). 6 or 8 Libraries were pooled (400ng each) and captured per manufacturer's instruction.

Targeted methylation capture was performed using a custom-designed 198kb panel of TargetCap methylation probes and reagents (BoKe Bioscience China, cat. No. MP121CD) and BGI adapter-specific blockers (BGI). 6 or 8 Libraries were pooled (400ng each) and captured per manufacturer's instruction.

The above captured cfDNA genome or methylation libraries were amplified and purified with AMPure XP beads. Library concentration was determined by Qubit dsDNA High Sensitivity Kit.

Captured libraries were sequenced on MGISEQ-2000 sequencer (MGI, China) using the $2 \times 100$ paired-end sequencing method per manufacturer's instruction.

**Preparation of two-human cfDNA blend sample**

White blood cells from the two donors were first sequenced to determine genotypes. 11 heterozygous from the "spike-in" donor and 58 homozygous single nucleotide

1 polymorphisms (SNPs) shared by the two donors covered by the IDT target capture

2 panel were then selected to measure the sensitivity and specificity of variant detection,

3 respectively. cfDNA samples of the two donors were mixed at a ratio of 1:200 to

4 simulate cfDNA with a 0.25% "spike-in" variant allele frequencies (VAFs) using the

5 heterozygous SNPs from the "spike-in" donor. Experiments were performed in

6 duplicates.

7

8 **Data analysis**

9 Adapter trimming and quality control of sequencing data were performed using Fastp

10 (v0.19.7) [30]. Paired-end reads of targeted sequencing and targeted methylation

11 sequencing were aligned to the hg19 reference human genome using bwa (v0.7.17)

12 and BitMapperBS (v1.0.0.8), respectively [31, 32]. Duplications were marked and

13 reads were deduplicated using sambamba (v0.6.8) [33]. Removal of sequencing errors

14 using duplex UMIs and variant calling were performed using custom python scripts.

15 Methylation rates of cytosines were calculated as #C/(#C+#T) for each CpG site with

16 at least 4x coverage, and M-bias of sequence reads was analyzed using MethylDackel

17 (v0.3.0) (https://github.com/dpryan79/MethylDackel). The cytosine conversion rate

18 was calculated using the methylation ratio of the spiked-in lambda DNA. GC-bias

19 metrics were analyzed using Picard Tools (v 2.10.10)

20 (http://broadinstitute.github.io/picard). Insert size distribution, base distribution of

21 reads, on-target rate, and sequencing depth were analyzed using custom Perl scripts.

22

**Data Access**

The data that support the findings of this study have been deposited in the CNSA
(https://db.cngb.org/cnsa/) of CNGBdb with accession number CNP0001331.


**Results**

**Comparison of dsLib *vs*. ssLib method for targeted deep sequencing of cfDNA**

We first compared double-stranded library (dsLib) preparation and single-stranded
library (ssLib) preparation methods for cfDNA mutation detection using deep
sequencing (Figure 1A, see Methods for more details). KAPA Hyper Prep kit, a
widely used NGS library construction kit which is based on the conventional dsDNA
library preparation methodology, was also included as a reference to evaluate
performance of our self-developed dsLib workflow. Duplex unique molecular
identifier (UMI)-based adapters were used to reduce noises that may derive from PCR
and/or sequencing errors [21] (see Methods for more details). Since in clinical
practice the amount of extracted cfDNA was often limited and highly variable [34],
we used 1 ng, 5 ng, and 10 ng cfDNA as inputs for library construction respectively
(Supplementary Table 1). Prepared libraries underwent hybridization-based target
enrichment procedure and captured libraries were sequenced to > 20000x raw average
depth (see Methods for more details). Results showed that library yields were similar
between dsLib and ssLib workflow (Supplementary Figure 1A). The two workflows
also achieved similar deduplicated depths (Figure 1B and Supplementary Table 1). Yet,
the ssLib workflow was more complicated and time-consuming than the dsLib

1   method (8h vs 3.5h, see Methods for more details). Remarkably, our self-developed

2   dsLib protocol showed significantly better performance than the commercial KAPA

3   workflow (Figure 1B and Supplementary Figure 1).

4

5   To further validate its ability to detect low abundance mutations in cfDNA and

6   confirm the limitation of detection (LOD), we applied our dsLib workflow on 40 ng

7   cfDNA spiked-in with cfDNA reference standards, simulating cfDNA samples with

8   defined variant allele frequencies (VAFs) (0.1%, 0.25%, and 0.5%, see Methods for

9   more details). 100% (24/24), 100% (24/24), 95.8% (23/24), and 91.7% (22/24)

10   mutations were detected in cfDNA samples with 1%, 0.5%, 0.25% and 0.1% expected

11   VAFs respectively, showing good correlation between the measured and expected

12   VAFs (Figure 1C). The analytical performance of our assay was also evaluated using

13   two-human cfDNA blend samples (see Methods for more details) to more closely

14   mimic cfDNA carrying low VAF mutations. Briefly, single nucleotide polymorphism

15   (SNP) sites where the "spike-in" donor carries heterozygous alleles while the

16   "background" donor carries homozygous alleles were used to evaluate assay

17   sensitivity; SNP sites where the "spike-in" donor and "background" donor carry the

18   same homozygous alleles were used to evaluate assay specificity. We obtained a

19   sensitivity of 95.5% (21/22 SNPs evaluated) and a specificity of 99.1% (115/116

20   SNPs evaluated) using the UMI error correction. Sensitivity was slightly lower

21   (86.4%; 19/22 SNPs evaluated) if only variants supported by at least one duplex UMI

22   family are considered true variants, while specificity was further improved to 100%

1    (116/116) (Supplementary Table 2). The results indicated that our custom dsLib

2    workflow provides satisfactory sensitivity for detection of low abundance variants in

3    cfDNA.

4

5    ctDNA has been proven to be shorter than cfDNA originated from normal cells [35,

6    36]. Theoretically, the ssLib workflow preferentially enriches short DNA molecules

7    and therefore may enrich ctDNA and improve its detection [28]. Copy number

8    variation (CNV) is a hallmark of cancer and could be used as a biomarker for ctDNA

9    [37]. Here, we compared CNV detectability of plasma cfDNA from lung cancer

10   patients using either dsLib or ssLib workflow to test the hypothesis that ssLib may

11   enrich for shorter ctDNA. We found no significant difference in CNV detection by

12   ssLib workflow *vs*. dsLib (Figure 1D), consistent with previous study which showed

13   that ssDNA-based workflow did not enrich for fetal DNA for NIPT, despite the

14   finding that it did enrich for shorter cfDNA fragments [38]. Taken together, our results

15   suggest that dsLib workflow is more preferable for ctDNA mutation detection.

16

17   **Comparison of dsLib *vs*. ssLib for cfDNA methylation sequencing**

18   Bisulfite sequencing has been a widely used sequencing technology for methylation

19   profiling, where methylation status of cytosines could be determined at

20   single-nucleotide resolution. This technology leverages the fact that methylated

21   cytosine remains unaffected when treated with sodium bisulfite, whereas

22   unmethylated cytosine is converted to uracil [39].

1

2　To compare performance of the single-stranded methylation sequencing library

3　construction (ssmLib) and the double-stranded methylation sequencing library

4　construction (dsmLib) (Figure 2A), we applied these two workflows on 1 ng, 5 ng,

5　and 10 ng cfDNA as inputs and captured the libraries with a 198 kb methylation

6　capture panel (Supplementary Table 3). Sequencing results showed that ssmLib

7　produced significantly higher library yields and deduplicated depths than dsmLib; the

8　on-target rates were also slightly higher in ssmLib libraries than dsmLib libraries

9　(Figure 2 B-C and Supplementary Figure 2). Notably, libraries produced by ssmLib

10　had more short insert fragments than those produced by dsmLib (Figure 2D). These

11　results can be attributed to DNA degradation caused by the bisulfite conversion

12　process, which involves high temperature and low pH conditions [40]: during ssmLib

13　workflow, the resulted short cfDNA fragments can still be captured by the

14　ssDNA-based adapter ligation; on the other hand, during dsmLib workflow, since

15　bisulfite was applied to the adapter-ligated dsDNA, excessive damage of the

16　templates will cause the libraries to lack paired adapters and lost during subsequent

17　amplification, resulting in much lower library yields and effective sequencing depths.

18　For measurements of CpG site methylation level, technical replicates showed good

19　correlation for both methods (Supplementary Figure 3) with various DNA input

20　amounts (Figure 2E).

21

22　Methylation bias (M-bias) is the term describing measured methylation levels that

1     deviate from true values, often observed near the 3' end of sequenced fragments due

2     to unmethylated cytosines introduced by the end-repair step during dsDNA-based

3     library preparation [41, 42]. Theoretically, libraries produced by ssmLib may show

4     less to no M-bias since there is no end-repair step involved (Figure 2A). Indeed, we

5     observed severe M-bias in Read 2 of dsmLib libraries, but not in ssmLib libraries

6     (Figure 2F). Taken together, these results suggest that ssmLib method is more

7     preferrable for the application of cfDNA methylation sequencing.

8

9     Recently, several enzyme-based cytosine conversion methods have been developed as

10    gentler substitutes for bisulfite conversion [43, 44]. We also compared performance of

11    a novel enzyme-based workflow (the NEB EM-seq Conversion Module) with the

12    conventional bisulfite conversion workflow (using the widely used ZYMO EZ DNA

13    Methylation Gold kit) (Figure 3A). EM-seq conversion module uses a two-step

14    enzymatic conversion process to detect modified cytosines: the first step uses TET2

15    and an oxidation enhancer to protect modified cytosines from downstream

16    deamination while converting 5-methylcytosine (5mC) to 5-carboxycytosine (5caC).

17    The second step uses APOBEC to enzymatically deaminate cytosine but does not

18    convert 5caC (the original 5mC). As expected, the ssmLib libraries produced by

19    bisulfite conversion had more short insert fragments than those produced by

20    enzyme-based conversion. Meanwhile, with ssmLib workflow, the bisulfite

21    conversion method generated significantly higher library yields and deduplicated

22    depths than enzymatic conversion, while similar cytosine conversion efficiencies were

1    observed for the two methods (Figure 3C-E and Supplementary Table 3). For dsmLib

2    workflow, however, there was no significant difference in either library yields,

3    deduplicated depths, or fragment size distributions between the two conversion

4    methods (Figure 3B, Supplementary Figure 4 and Supplementary Table 3). Among all,

5    bisulfite conversion coupled with ssmLib workflow still achieved the highest

6    deduplicated sequencing depth. Also, the enzymatic conversion is more

7    time-consuming (8h vs 3.5h) than the bisulfite conversion. Taken together, our results

8    favor bisulfite conversion coupled with ssLib workflow for cfDNA methylation

9    sequencing.

10

11   **Conclusion**

12   The double-stranded library preparation method is more advantageous for ctDNA

13   mutation detection thanks to the higher data quality and easy workflow. Meanwhile,

14   bisulfite conversion coupled with single-stranded library preparation showed overall

15   better performance for cfDNA methylation sequencing. Our results suggest that when

16   performing high-throughput sequencing for cfDNA, depending on the downstream

17   applications, these two library preparation methods should be chosen accordingly.

18

19   **Discussion**

20   In recent decades, thanks to the development of NGS technology, the cost of

21   high-throughput DNA sequencing had dropped dramatically, making it affordable for

22   researchers worldwide [45]. Library construction is a key step for successful NGS

1   workflow and high-quality data generation. In this study, we compared dsDNA and

2   ssDNA-based library construction methods for cfDNA deep sequencing (i.e., for

3   ctDNA variant detection) and methylation profiling.

4

5   A major difference between dsDNA and ssDNA based cfDNA library construction

6   methods is that cfDNA molecules harboring single-strand breaks (also called nicks) as

7   well as those existing as ssDNA form could be utilized by the ssLib (or ssmLib)

8   workflow but would not be ligatable when using the dsLib (or dsmLib) workflow

9   (Figure 1A and 2A). Naturally nicked and/or single-stranded cfDNA molecules may

10   only be a very small fraction hence this difference would be expected to be small and

11   may not cause significant impact on the effective sequencing depth. Indeed, we

12   observed similar deduplicated depth for cfDNA libraries generated using ssLib or

13   dsLib workflow (Supplementary Figure 1B); in fact, deduplicated depth of ssLib

14   libraries were even slightly inferior than dsLib, possibly due to the fact that ssLib

15   workflow is lengthier and requires more beads purification and therefore may cause

16   template loss.

17

18   Using detected CNV level as an indicator of ctDNA fraction, we also showed that

19   there was no significant enrichment of ctDNA by ssLib compared to dsLib workflow

20   (Figure 1D), consist with previous research conducted in the setting of NIPT which

21   showed that ssLib workflow does not enrich for shorter fetal DNA [28, 38]. It was

22   suggested that intrinsic biological differences between fetal DNA and maternal DNA

23   molecules might account for the failure of ssDNA workflow to enrich for fetal DNA

1    [28, 38], and similar mechanism may also explain our results for ctDNA. Further

2    study is needed to deepen our understanding of cfDNA/ctDNA generation processes

3    and/or to develop novel library construction methods for ctDNA enrichment.

4

5    Importantly, application of dsLib workflow further allows utilization of duplex UMIs,

6    which make it possible to recover original dsDNA fragments following paired-end

7    sequencing and utilize the information from complementary strands of DNA

8    molecules to correct possible PCR and/or sequencing errors, achieving an extra low

9    base error rate and higher specificity with variant detection [21]. Taken together, our

10    results demonstrate that current state-of-the-art dsDNA-based library preparation is

11    more preferable for the application of deep sequencing for ctDNA variant detection.

12

13    On the contrary, a clear advantage was observed for ssmLib libraries for bisulfite

14    sequencing compared to dsmLib (Figure 2B-C). This is because libraries were

15    constructed before bisulfite conversion during the dsmLib workflow (Figure 2A), and

16    the nicked DNA resulting from the bisulfite conversion won't be sequenced due to the

17    lack of paired adapters. During the ssmLib workflow, however, cfDNA ligation

18    happens after the bisulfite treatment, where the nicked and single-stranded DNA

19    molecules resulting from the bisulfite treatment can still be ligated with adapters,

20    therefore preserving more DNA templates for sequencing (Figure 2A and 2D),

21    eventually achieving a higher effective depth.

22

1    Theoretically, gentler enzyme-based cytosine conversion method would avoid the

2    assumed template loss caused by bisulfite treatment on the adapter-ligated library

3    fragments and may therefore greatly improve the results of dsLib workflow when

4    used for methylation profiling. Our results, however, still favored the bisulfite

5    conversion for both dsmLib and ssmLib workflow due to the higher library yields as

6    well as higher deduplicated depths, suggesting that there may be excessive loss of

7    templates during the enzyme-conversion workflow (Figure 3C-E and Supplementary

8    Figure 4). Indeed, this may be attributed to the two rounds of beads purification in the

9    enzyme-based conversion. Also, the current enzyme-based conversion workflow is

10    more labor- and time-consuming compared to the bisulfite conversion. Development

11    of more effective enzyme-based cytosine conversion methods may require

12    improvements in template recovery and further simplification of the workflow.

13

14    In addition, methylation bias (M-bias) was proposed to be an important library

15    preparation quality metric for methylation profiling, since its existence could cause

16    significant bias in measurements of methylation level [41, 42]. M-bias is caused by

17    the end-repair step in the conventional dsmLib workflow which typically recruits

18    unmethylated cytosines instead of methylated cytosines during the fill-in reaction

19    (Figure 2A). The filled-in cytosines were then converted to uracils regardless of the

20    original cytosine methylation status in the genome, resulting in incorrect methylation

21    level being assigned to the 3' end of the sequenced reads [41, 42]. The ssmLib method

22    could perfectly overcome this problem since it does not involve an end-repair step and

1   is a post–bisulfite conversion library construction method (Figure 2A and 2F), adding

2   another advantage to the ssmLib method. Collectively, our results favor the use of

3   ssmLib workflow for cfDNA methylation profiling. Our findings could help

4   researchers maximize the efficiency of NGS library preparation and produce better

5   quality sequencing data.

6

7

# 1 Reference

2 1.    Lui, Y.Y.N., et al., *Predominant Hematopoietic Origin of Cell-free DNA in Plasma and Serum*
3       *after Sex-mismatched Bone Marrow Transplantation.* Clinical Chemistry, 2002. **48**(3): p.
4       421-427.

5 2.    Snyder, Matthew W., et al., *Cell-free DNA Comprises an In Vivo Nucleosome Footprint*
6       *that Informs Its Tissues-Of-Origin.* Cell, 2016. **164**(1): p. 57-68.

7 3.    Aravanis, A.M., M. Lee, and R.D. Klausner, *Next-Generation Sequencing of Circulating Tumor*
8       *DNA for Early Cancer Detection.* Cell, 2017. **168**(4): p. 571-574.

9 4.    Lo, Y.M.D., et al., *Maternal Plasma DNA Sequencing Reveals the Genome-Wide Genetic and*
10      *Mutational Profile of the Fetus.* Science Translational Medicine, 2010. **2**(61): p. 61ra91.

11 5.   Chan, A.K., R.W. Chiu, and Y.M. Lo, *Cell-free nucleic acids in plasma, serum and urine: a new*
12      *tool in molecular diagnosis.* Ann Clin Biochem, 2003. **40**(Pt 2): p. 122-30.

13 6.   Schwarzenbach, H., D.S.B. Hoon, and K. Pantel, *Cell-free nucleic acids as biomarkers in cancer*
14      *patients.* Nature Reviews Cancer, 2011. **11**(6): p. 426-437.

15 7.   van der Pol, Y. and F. Mouliere, *Toward the Early Detection of Cancer by Decoding the*
16      *Epigenetic and Environmental Fingerprints of Cell-Free DNA.* Cancer Cell, 2019. **36**(4): p.
17      350-368.

18 8.   Bronkhorst, A.J., V. Ungerer, and S. Holdenrieder, *The emerging role of cell-free DNA as a*
19      *molecular marker for cancer management.* Biomolecular detection and quantification, 2019.
20      **17**: p. 100087-100087.

21 9.   Volik, S., et al., *Cell-free DNA (cfDNA): Clinical Significance and Utility in Cancer Shaped By*
22      *Emerging Technologies.* Molecular Cancer Research, 2016. **14**(10): p. 898.

23 10.  Corcoran, R.B. and B.A. Chabner, *Application of Cell-free DNA Analysis to Cancer Treatment.*
24      New England Journal of Medicine, 2018. **379**(18): p. 1754-1765.

25 11.  Gormally, E., et al., *Circulating free DNA in plasma or serum as biomarker of carcinogenesis:*
26      *Practical aspects and biological significance.* Mutation Research/Reviews in Mutation
27      Research, 2007. **635**(2): p. 105-117.

28 12.  Chabon, J.J., et al., *Integrating genomic features for non-invasive early lung cancer detection.*
29      Nature, 2020.

30 13.  Schuster, S.C., *Next-generation sequencing transforms today's biology.* Nature Methods, 2008.
31      **5**(1): p. 16-18.

32 14.  Shendure, J., et al., *DNA sequencing at 40: past, present and future.* Nature, 2017. **550**(7676):
33      p. 345-353.

34 15.  Tucker, T., M. Marra, and J.M. Friedman, *Massively parallel sequencing: the next big thing in*
35      *genetic medicine.* American journal of human genetics, 2009. **85**(2): p. 142-154.

36 16.  Green, E.D., E.M. Rubin, and M.V. Olson, *The future of DNA sequencing.* Nature, 2017.
37      **550**(7675): p. 179-181.

38 17.  Mardis, E.R., *DNA sequencing technologies: 2006–2016.* Nature Protocols, 2017. **12**(2): p.
39      213-218.

40 18.  Lundberg, D.S., et al., *Practical innovations for high-throughput amplicon sequencing.* Nature
41      Methods, 2013. **10**(10): p. 999-1002.

42 19.  Gansauge, M.T., et al., *Single-stranded DNA library preparation from highly degraded DNA*
43      *using T4 DNA ligase.* Nucleic Acids Res, 2017. **45**(10): p. e79.

20. Newman, A.M., et al., *An ultrasensitive method for quantitating circulating tumor DNA with broad patient coverage.* Nat Med, 2014. **20**(5): p. 548-54.

21. Newman, A.M., et al., *Integrated digital error suppression for improved detection of circulating tumor DNA.* Nat Biotechnol, 2016. **34**(5): p. 547-555.

22. Raine, A., et al., *SPlinted Ligation Adapter Tagging (SPLAT), a novel library preparation method for whole genome bisulphite sequencing.* Nucleic Acids Res, 2017. **45**(6): p. e36.

23. Olova, N., et al., *Comparison of whole-genome bisulfite sequencing library preparation strategies identifies sources of biases affecting DNA methylation data.* Genome Biol, 2018. **19**(1): p. 33.

24. Keats, J.J., et al., *Whole Genome Library Construction for Next Generation Sequencing*, in *Disease Gene Identification: Methods and Protocols*, J.K. DiStefano, Editor. 2018, Springer New York: New York, NY. p. 151-161.

25. Gansauge, M.T. and M. Meyer, *Single-stranded DNA library preparation for the sequencing of ancient or damaged DNA.* Nat Protoc, 2013. **8**(4): p. 737-48.

26. Glocke, I. and M. Meyer, *Extending the spectrum of DNA sequences retrieved from ancient bones and teeth.* Genome Res, 2017. **27**(7): p. 1230-1237.

27. Sanchez, C., et al., *Circulating nuclear DNA structural features, origins, and complete size profile revealed by fragmentomics.* JCI insight, 2021. **6**(7): p. e144561.

28. Vong, J.S., et al., *Single-stranded DNA library preparation preferentially enriches short maternal DNA in maternal plasma.* J Clinical chemistry, 2017. **63**(5): p. 1031-1037.

29. Chen, K., et al., *Non-invasive lung cancer diagnosis and prognosis based on multi-analyte liquid biopsy.* Molecular Cancer, 2021. **20**(1): p. 23.

30. Chen, S., et al., *fastp: an ultra-fast all-in-one FASTQ preprocessor.* Bioinformatics, 2018. **34**(17): p. i884-i890.

31. Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows–Wheeler transform.* Bioinformatics, 2009. **25**(14): p. 1754-1760.

32. Cheng, H. and Y. Xu, *BitMapperBS: a fast and accurate read aligner for whole-genome bisulfite sequencing.* bioRxiv, 2018: p. 442798.

33. Tarasov, A., et al., *Sambamba: fast processing of NGS alignment formats.* Bioinformatics, 2015. **31**(12): p. 2032-2034.

34. Bronkhorst, A.J., J. Aucamp, and P.J. Pretorius, *Cell-free DNA: Preanalytical variables.* Clinica Chimica Acta, 2015. **450**: p. 243-253.

35. Underhill, H.R., et al., *Fragment length of circulating tumor DNA.* J PLoS genetics, 2016. **12**(7): p. e1006162.

36. Jiang, P., et al., *Lengthening and shortening of plasma DNA in hepatocellular carcinoma patients.* Proceedings of the National Academy of Sciences of the United States of America, 2015. **112**(11): p. E1317-E1325.

37. Li, J., et al., *Cell-free DNA copy number variations in plasma from colorectal cancer patients.* Molecular oncology, 2017. **11**(8): p. 1099-1111.

38. Moser, T., et al., *Single-Stranded DNA Library Preparation Does Not Preferentially Enrich Circulating Tumor DNA.* Clinical Chemistry, 2017. **63**(10): p. 1656-1659.

39. Hayatsu, H., *Discovery of bisulfite-mediated cytosine conversion to uracil, the key reaction for DNA methylation analysis--a personal account.* Proceedings of the Japan Academy. Series B, Physical and biological sciences, 2008. **84**(8): p. 321-330.

1  40.  Grunau, C., S.J. Clark, and A. Rosenthal, *Bisulfite genomic sequencing: systematic*
2      *investigation of critical experimental parameters.* Nucleic Acids Res, 2001. **29**(13): p. E65-5.
3  41.  Hansen, K.D., B. Langmead, and R.A. Irizarry, *BSmooth: from whole genome bisulfite*
4      *sequencing reads to differentially methylated regions.* Genome Biology, 2012. **13**(10): p. R83.
5  42.  Ryan, D.P. and D. Ehninger, *Bison: bisulfite alignment on nodes of a cluster.* BMC
6      Bioinformatics, 2014. **15**(1): p. 337.
7  43.  Vaisvila, R., et al., *EM-seq: Detection of DNA Methylation at Single Base Resolution from*
8      *Picograms of DNA.* bioRxiv, 2019: p. 2019.12.20.884692.
9  44.  Liu, Y., et al., *Bisulfite-free direct detection of 5-methylcytosine and 5-hydroxymethylcytosine*
10     *at base resolution.* Nat Biotechnol, 2019. **37**(4): p. 424-429.
11 45.  Kchouk, M., et al., *Generations of sequencing technologies: from first to next generation.*
12     2017. **9**(3).

13

1  **Figure Legends**

2  **Figure 1: Comparison of dsLib and ssLib workflow for cfDNA mutation detection by**
3  **targeted deep sequencing. (A)** Schematic view of our self-developed dsLib and ssLib workflow.
4  See Methods for more details. **(B)** Deduplicated depths of libraries constructed by dsLib, ssLib,
5  and the KAPA kit. Duplicates were performed for each experimental condition. Data are presented
6  as mean ± SD. N.S, p>=0.05; *, 0.01<=p<0.05; **, 0.001<=p<0.01; ****, p<0.0001, as calculated
7  by Student's t-test. **(C)** Detection of low VAF mutations by dsLib workflow in simulated cfDNA
8  samples. Triplicates were performed for each experimental condition. Data are presented as mean
9  ± SD. The numbers in parentheses represent the number of detected mutations/total mutations. **(D)**
10 CNVs detected by dsLib and ssLib methods respectively, in plasma cfDNA samples from lung
11 cancer patients P1 and P2. X-axis, chromosome. Y-axis, CNV adjusted by GC content and
12 map-ability.

13

14 **Figure 2: Comparison of dsmLib and ssmLib workflow for cfDNA methylation profiling by**
15 **bisulfite sequencing. (A)** Schematic view of our dsmLib and ssmLib procedures. **(B)** On-target
16 rates and **(C)** deduplicated depths of libraries prepared by dsmLib and ssmLib workflow.
17 Duplicates were performed for each experimental condition. Data are presented as mean ± SD.
18 N.S, p>=0.05; *, 0.01<=p<0.05; **, 0.001<=p<0.01; ****, p<0.0001, as calculated by Student's
19 t-test. **(D)** Size distributions of library insert fragments. **(E)** Pearson correlation of methylation
20 levels between ssmLib (x-axis) and dsmLib (y-axis) libraries. **(F)** M-bias plots of libraries
21 prepared by dsmLib and ssmLib workflow. For each row from left to right: Read 1 ++ strand,
22 Read 1 -+strand, Read 2 +- strand, and Read 2 --strand. X-axis, position in read (bp). Y-axis,
23 methylation level (%).

24

25 **Figure 3: Comparison of the chemical and enzymatic cytosine conversion for cfDNA**
26 **methylation sequencing. (A)** Technical principles of bisulfite conversion and enzymatic
27 conversion. 5caC, 5-carboxylcytosine. T, thymine. **(B)** Size distribution of library insert fragments.
28 X-axis, fragment size (bp). Y-axis, frequency count. **(C)** CT conversion rates, **(D)** library yields,
29 and **(E)** deduplicated depths of ssmLib libraries. Duplicates were performed for each experimental
30 condition. Data are presented as mean ± SD. N.S, p>=0.05; *, 0.01<=p<0.05, as calculated by
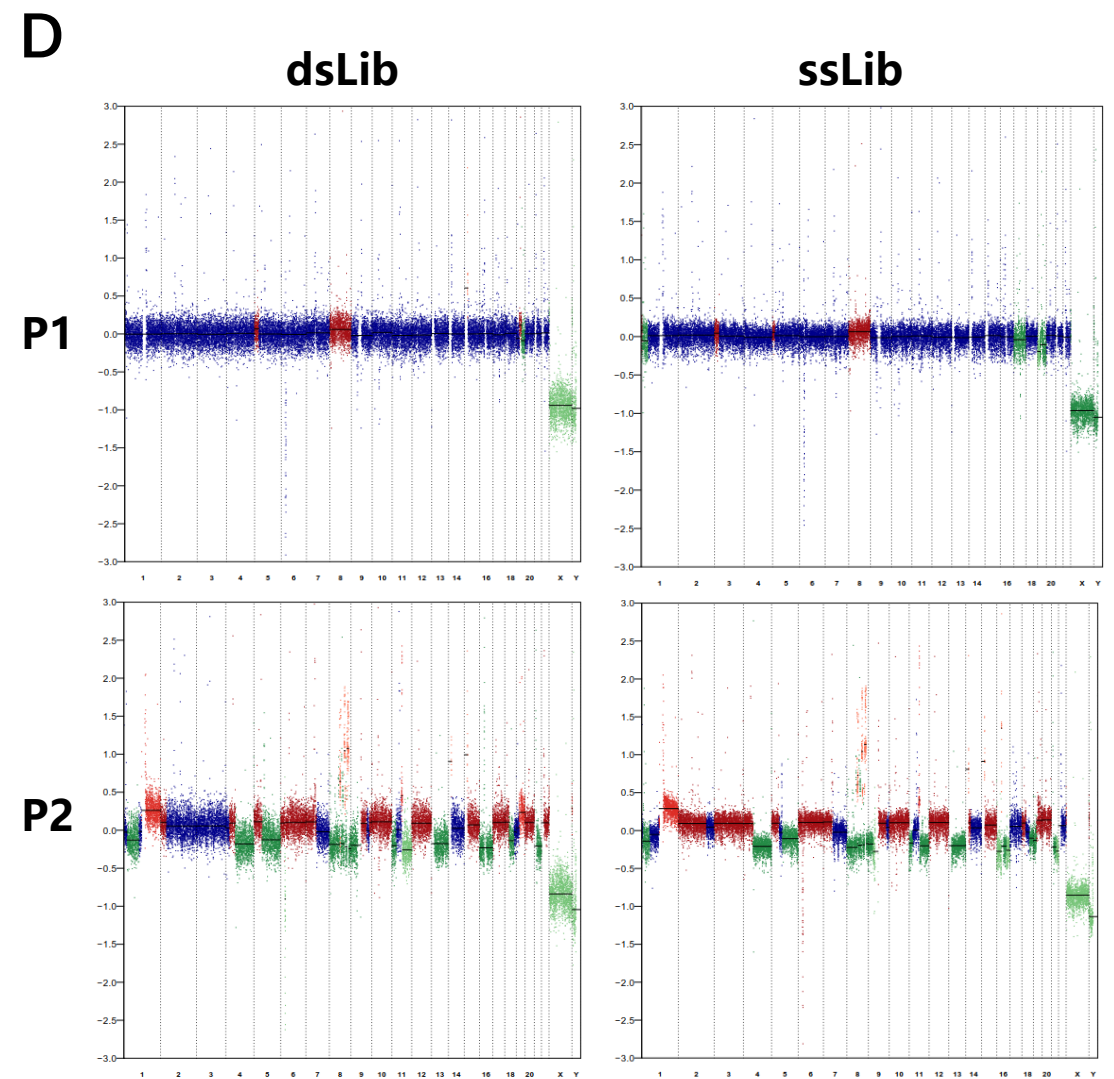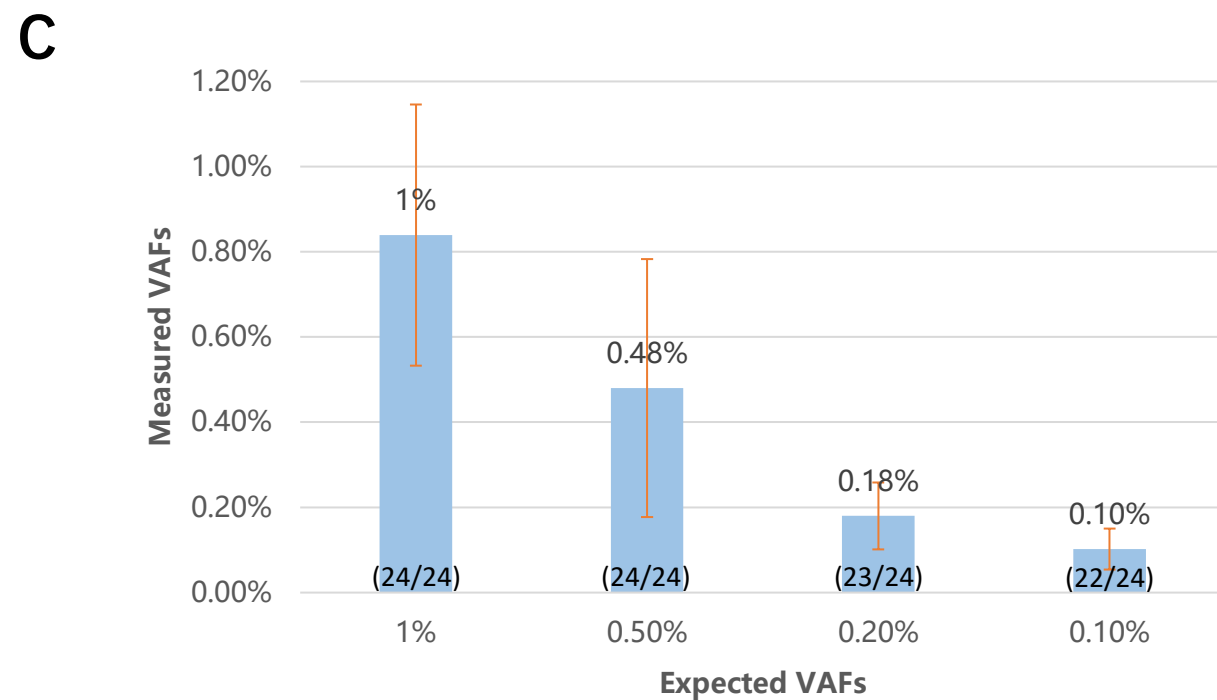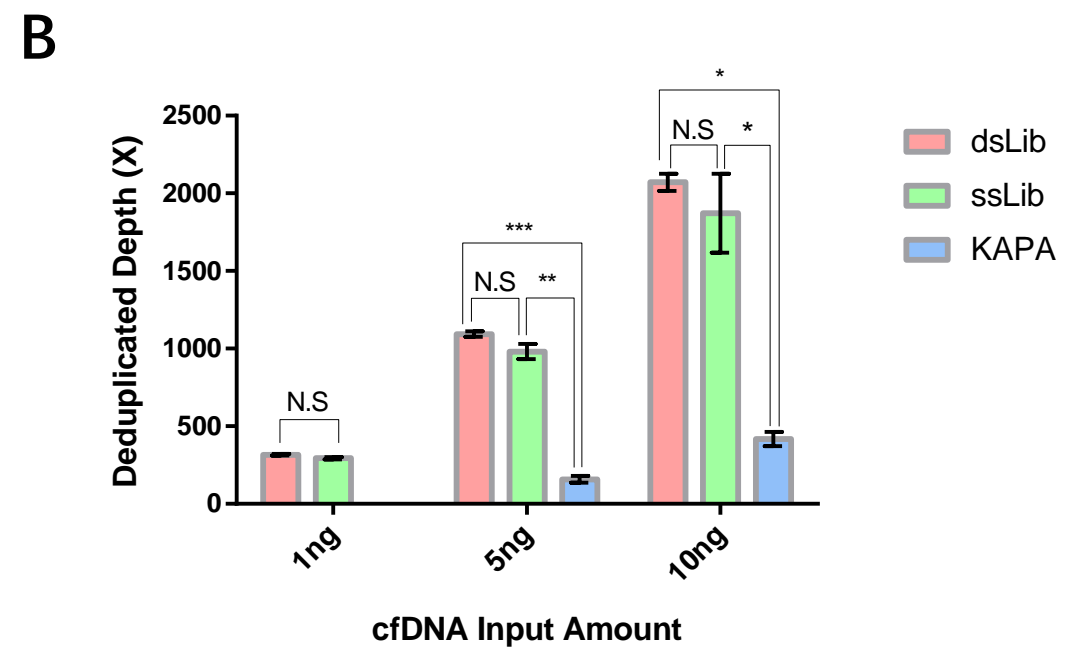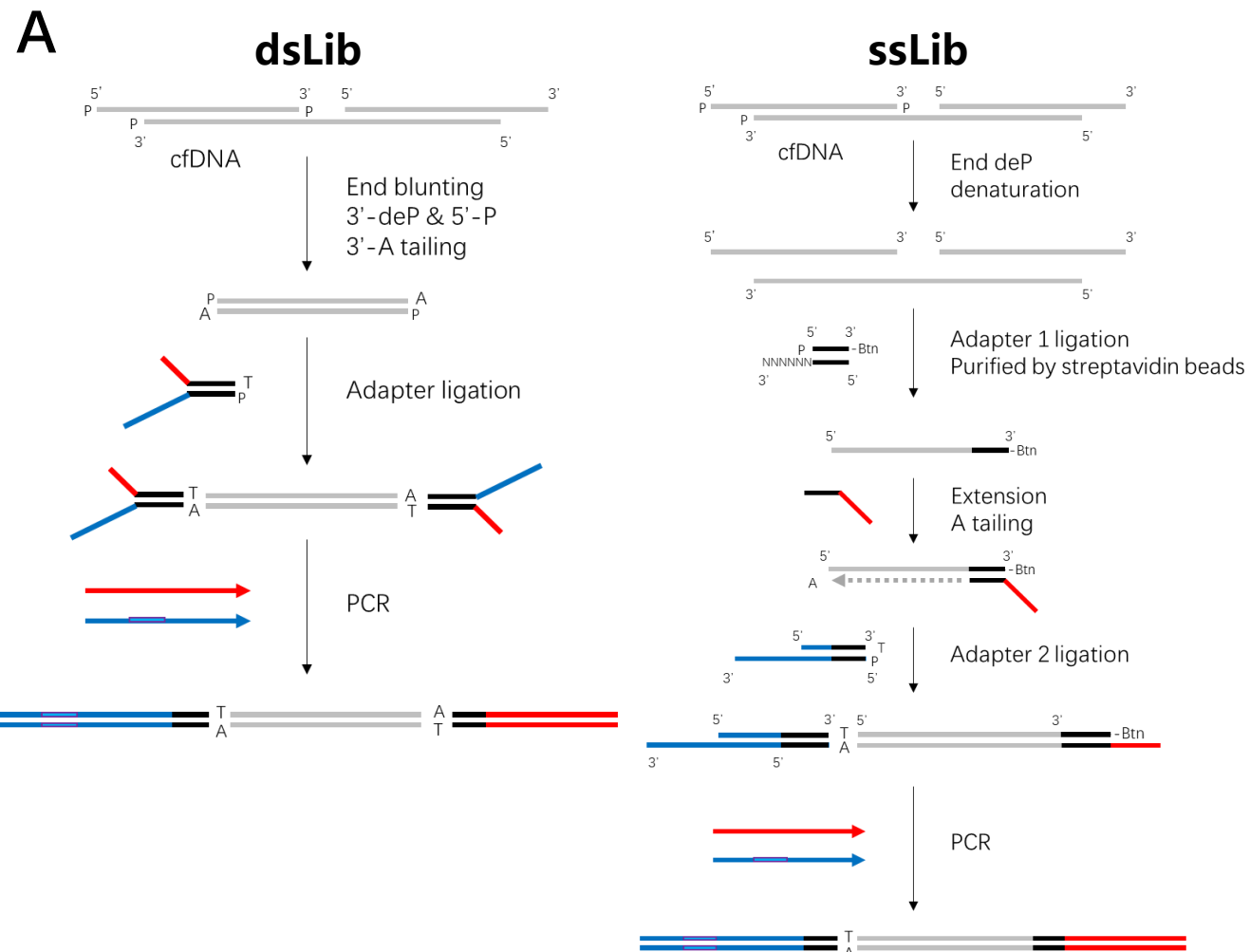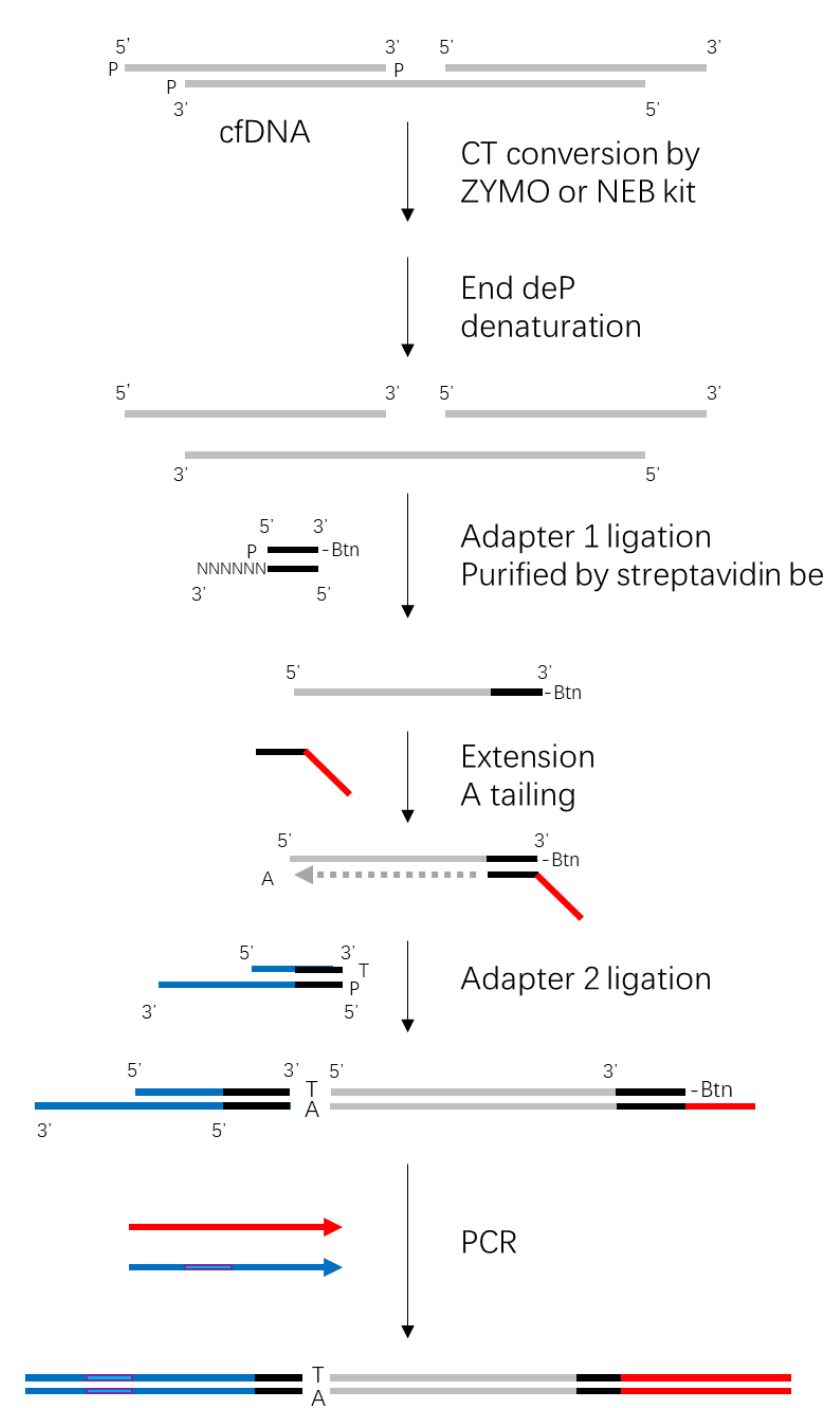31 Student's t-test.

32

1    **Supplementary Figure 1: (A) Library yields and (B) mean fragment lengths of**
2    **sequenced libraries constructed by dsLib, ssLib, and the KAPA workflow.**
3    Duplicates were performed for each experimental condition. Data are presented as
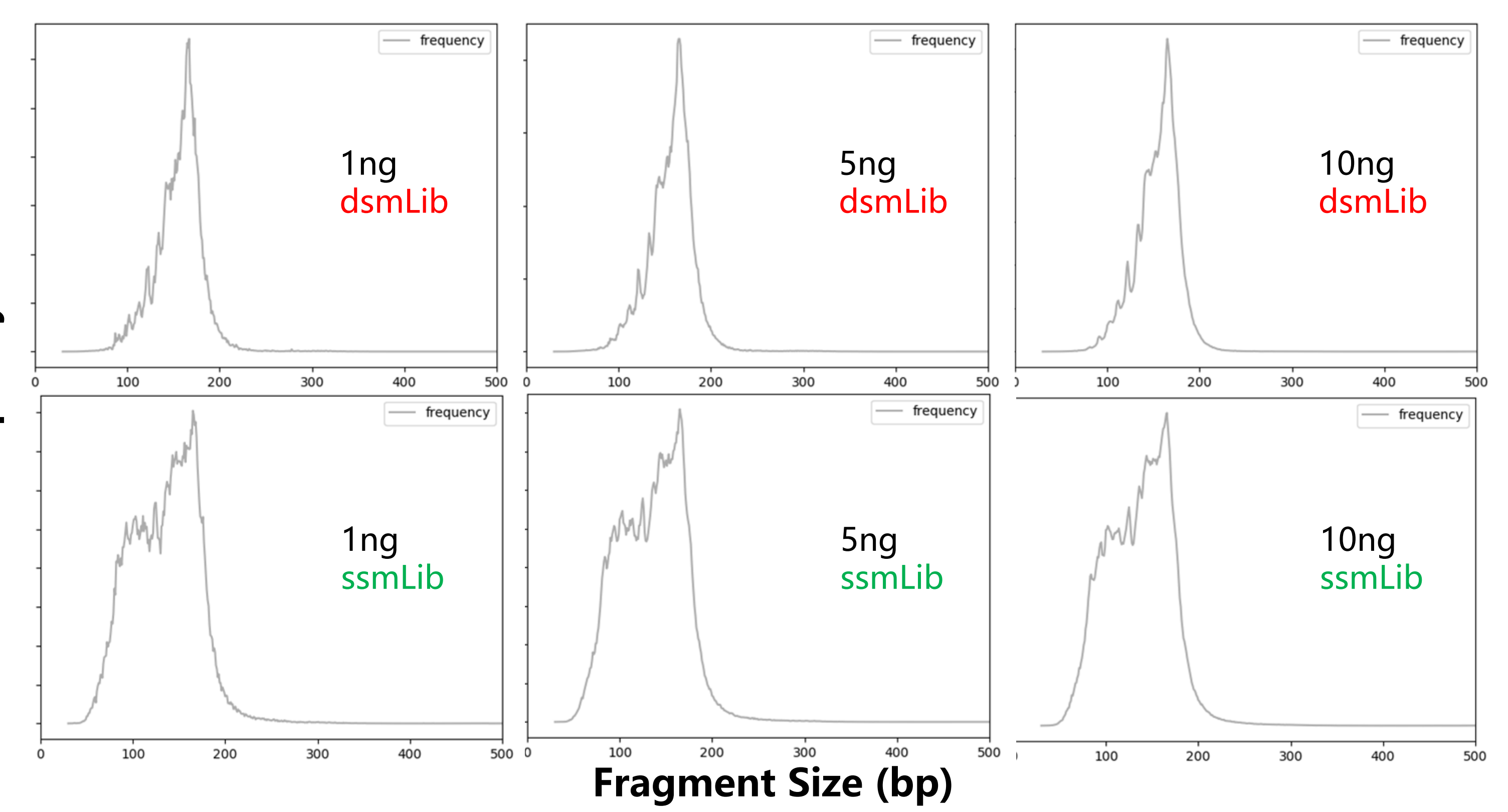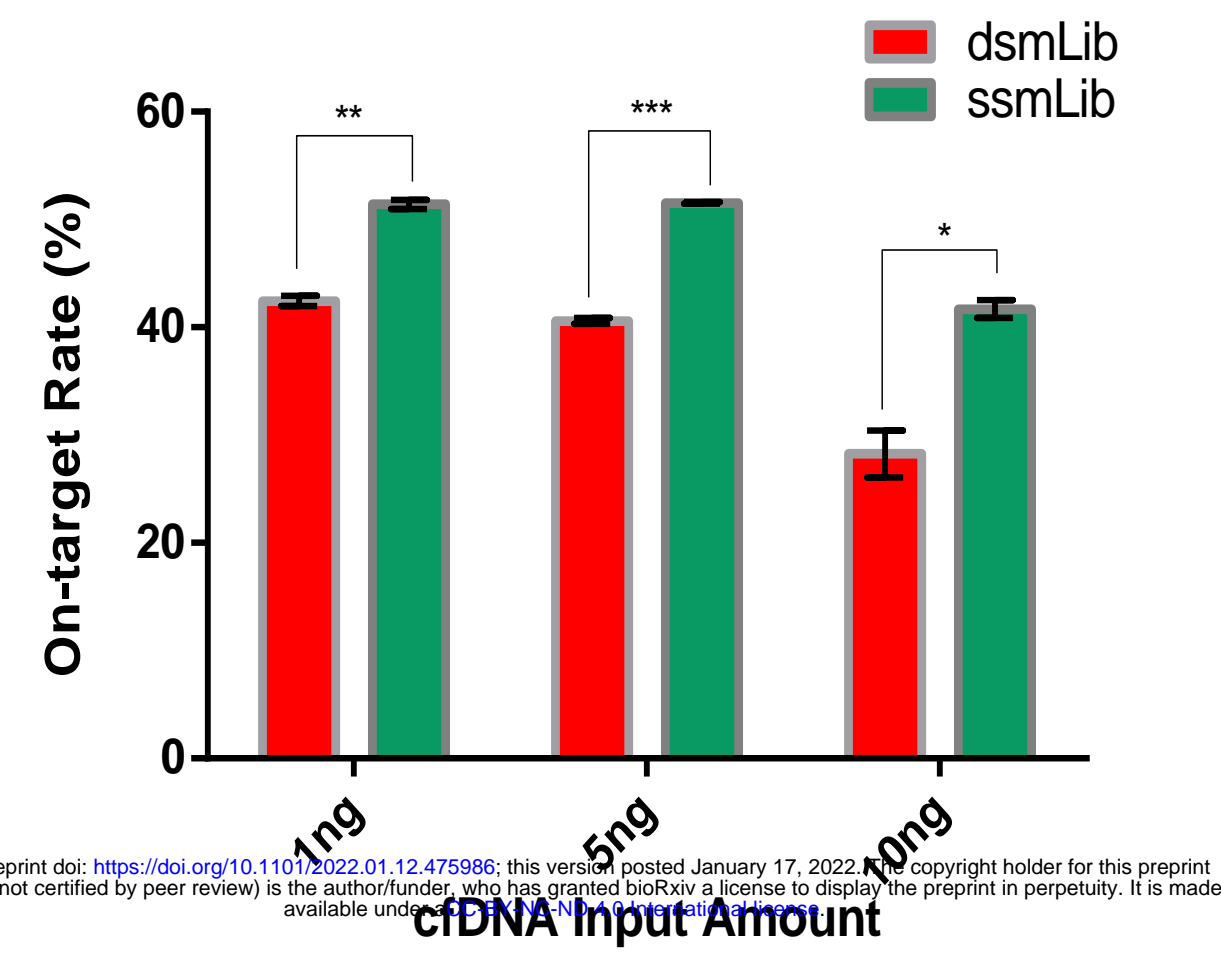4    mean ± SD. N.S, p>=0.05; *, 0.01<=p<0.05; **, 0.001<=p<0.01; ****, p<0.0001, as
5    calculated by Student's t-test.

6    **Supplementary Figure 2: Library yields by dsmLib and ssmLib.** Duplicates were
7    performed for each experimental condition. Data are presented as mean ± SD. N.S,
8    p>=0.05; *, 0.01<=p<0.05; **, 0.001<=p<0.01, as calculated by Student's t-test.

9    **Supplementary Figure 3: Pearson correlation of methylation levels between (A)**
10    **dsmLib or (B) ssmLib libraries.**

11    **Supplementary Figure 4: (A) CT conversion rates, (B) library yields, and (C)**
12    **deduplicated depths of dsmLib libraries using bisulfite and enzymatic conversion.**
13    Duplicates were performed for each experimental condition. Data are presented as
14    mean ± SD. N.S, p>=0.05; *, 0.01<=p<0.05, as calculated by Student's t-test.

15

A

**dsLib**

**ssLib**

B

C

D

**dsLib**                **ssLib**

P1

P2