1

2  &lt;Long title&gt;

3  # A new targeted capture method using bacterial artificial

4  chromosome (BAC) libraries as baits for sequencing relatively

5  large genes

6

7  &lt;Short title&gt;

8  # A new targeted capture method using BAC baits

9

10

11  Kae Koganebuchi[1#a], Takashi Gakuhari[2], Hirohiko Takeshima[3], Kimitoshi Sato[4],

12  Kiyotaka Fujii[4], Toshihiro Kumabe[4], Satoshi Kasagi[5], Takehiro Sato[6], Atsushi Tajima[6],

13  Hiroki Shibata[7], Motoyuki Ogawa[1,8], Hiroki Oota[1,8]*

14

15

16  [1] Department of Biological Structure, Kitasato University Graduate School of Medical
17  Sciences, Sagamihara, Kanagawa, Japan

18  [2] Center for Cultural Resource Studies, Kanazawa University, Kanazawa, Ishikawa, Japan

19    [3] Department of Marine Biology, School of Marine Science and Technology, Tokai

20    University, Shizuoka, Shizuoka, Japan

21    [4] Department of Neurosurgery, Kitasato University School of Medicine, Sagamihara,

22    Kanagawa, Japan

23    [5] School of Marine Biosciences, Kitasato University, Sagamihara, Kanagawa, Japan

24    [6] Department of Bioinformatics and Genomics, Graduate School of Advanced Preventive

25    Medical Sciences, Kanazawa University, Kanazawa, Ishikawa, Japan

26    [7] Medical Institute of Bioregulation, Kyushu University, Higashi, Fukuoka, Japan

27    [8] Department of Anatomy, Kitasato University School of Medicine, Sagamihara,

28    Kanagawa, Japan

29    [#a] Current Address: Department of Human Biology and Anatomy, Graduate School of

30    Medicine, University of the Ryukyus, Nishihara, Okinawa, Japan

31

32    * Corresponding author

33    E-mail: hiroki_oota@med.kitasato-u.ac.jp (HO)

34
35
36

## Author Contributions

38    Conceptualization: KK, TG, HO

39    Formal Analysis: KK

40    Investigation: KK, TS, HS

41    Resources: KS, KF, TK, HO

42    Supervision: TG, SK, HT, AT, MO

43    Writing - Original Draft Preparation: KK

44    Writing - Review & Editing: KK, HO

45

## Funding

52

## Competing interests

54    The authors have declared that no competing interests exist.

55

56

# **Abstract**

57

58    To analyze a specific genome region using next-generation sequencing technologies, the

59    enrichment of DNA libraries with targeted capture methods has been standardized. For

60    enrichment of mitochondrial genome, a previous study developed an original targeted

61    capture method that use baits constructed from long-range polymerase chain reaction

62    (PCR) amplicons, common laboratory reagents, and equipment. In this study, a new

63    targeted capture method is presented, that of bacterial artificial chromosome (BAC)

64    double capture (BDC), modifying the previous method, but using BAC libraries as baits

65    for sequencing a relatively large gene. We applied the BDC approach for the 214 kb

66    autosomal region, *ring finger protein 213*, which is the susceptibility gene of moyamoya

67    disease (MMD). To evaluate the reliability of BDC, cost and data quality were

68    compared with those of a commercial kit. While the ratio of duplicate reads was higher,

69    the cost was less than that of the commercial kit. The data quality was sufficiently the

70    same as that of the kit. Thus, BDC can be an easy, low-cost, and useful method for

71    analyzing individual genome region with substantial length.

72

73

74

# Introduction

75

76    The high-throughput sequencing technology, next-generation sequencing (NGS), has

77    made a striking impact on genomic research and the entire biological field. The NGS

78    technology is often called massively parallel sequencing because it effectively conducts

79    whole-genome sequencing in a relatively short time [1]. NGS enables researchers to

80    analyze the whole human genome of about 3 Gbp and identify all of the 30,000 genes in

81    only 1 week [2]. To analyze specific regions (e.g., whole exons, and already known

82    disease-related genes) using NGS, enrichment of DNA libraries with targeted capture

83    methods are standardized.

84        In capture methods, probes for enriching the targeted genomic regions are called

85    "baits" that attract molecules of interest as in fishing. There are two major approaches

86    for relatively large-scale genomic-region enrichment, the "on-array" and "in-solution"

87    methods. Both of these approaches target sequences up to several hundred kbp. In the

88    on-array capture method (Roche NimbleGen products), microarrays immobilize baits

89    that hybridize with the targeted regions and are used to enrich the genomic region of

90    interest. Meanwhile, the in-solution method (e.g., Roche, Illumina, Agilent

91    Technologies, and MYcroarray products) use biotinylated DNA or RNA baits to enrich

92    targeted region. Because DNA-RNA hybrids show higher efficiency than do DNA-DNA

5

93    hybrids, RNA baits are used in some systems [3]. The targeted DNA is recovered using

94    streptavidin-labeled magnetic beads. The in-solution approach has advantages compared

95    with the on-array method: the reagent cost is lower, less DNA is required, and it is

96    easily scaled because the capture method can be conducted entirely in small tubes [4].

97       Small-scale targeted capture methods have been proposed for the enrichment of

98    the complete mitochondrial genome (mtDNA) [5–7]. Maricic et al. (2010) presented a

99    capture method for the mtDNA molecules that used biotinylated polymerase chain

100   reaction (PCR) amplicons as baits. Human mtDNA is approximately 16.6 kbp long.

101   When constructing the baits, two primer sets of long-range PCR that amplify >9 kbp

102   regions are sufficient to cover whole mtDNA genome sequencing. The long-range PCR

103   amplicons are sheared with sonicators. The sheared amplicons are biotinylated and used

104   for the enrichment. This targeted capture method is cited by approximately 200 previous

105   studies that analyzed genomes of modern or ancient organisms (e.g., humans,

106   pathogens, animals, and fishes). The commercial targeted capture kits for small regions

107   cost approximately 250–900 USD per reaction. The method provided by Maricic et al.

108   (2010) is approximately 50 USD per reaction and much less expensive than commercial

109   methods. If the Maricic's method can be applied for large genes, then it definitely saves

110   the cost.

111    In order to enrich larger genomic regions than mtDNA, we conceived of using

112    bacterial artificial chromosome (BAC) libraries as baits, instead of PCR amplicons.

113    BAC is a vector that can carry DNA fragments of >300 kbp [8], and can be amplified by

114    culturing *E. coli* harboring the BAC. Human BAC libraries constructed in previous

115    studies [9–11] are available and distributed through resource centers. We named the

116    novel approach presenting in this study as "BAC double capture (BDC) method." Here

117    we show the conditions optimized for the BDC method, and the satisfactory efficiency

118    evaluated in comparison to the commercial enrichment kit in the NGS output data.

119

120

121    **Materials and Methods**

122

123    **Preparation of indexed libraries for testing experimental**

124    **conditions by BAC single capture (BSC) and BAC double**

125    **capture (BDC) with PrimeSTAR**

126    A DNA solution purchased from the Health Science Research Resources Bank (Osaka,

127    Japan) was used in the present study. The DNA concentration was measured using a

128    NanoPhotometer (Implen; CA, USA). The total amount of 5 μg of DNA was sheared

129    using a Covaris S2 sonicator (Covaris; MA, USA). The target peak was set at 400 bp. A

130    total of 50 ng of DNA was used to produce the indexed library, using an NEBNext Ultra

131    DNA Library Prep Kit and Multiplex Oligos for Illumina (New England BioLabs; MA,

132    USA). Sheared DNA was end-repaired, dA-tailed, and ligated to Illumina specific

133    adaptors. Sizes of the adaptor-ligated DNAs are selected to an approximate insert size,

134    400–500 bp, by Agencourt AMPure XP beads (Beckman Coulter; CA, USA). The

135    genomic DNA shotgun library was amplified with 4 PCR reactions using a primer pair,

136    Sol_bridge_P5 and Sol_bridge_P7, which was as presented in Maricic et al. (2010). We

137    used 500 pmol of the library as a template for PCR in a 50-μL solution containing

138    deoxynucleotide (dNTP) 0.2 mM, 0.2 μM of each primer, 1.25 U of PrimeSTAR GXL

139    DNA Polymerase (Takara Bio; Shiga, Japan). PCR was carried out following the

140    cycling reaction: 15 cycles of denaturation at 98°C for 10 sec, annealing at 60°C for 15

141    sec, extension at 68°C for 50 sec. Those PCR products were pooled and the solution was

142    purified using a MinElute PCR Purification Kit (Qiagen; Hilden, Germany), and it was

143    then eluted it into 23 μL buffer EB (Qiagen). The concentration of the solution was

144    measured using a NanoPhotometer (Implen). The total amount of 2μg per capture

145    reaction was obtained.

146

8

## Bait production

148    The BAC from the CHORI-17 library (ID number: CH17-24F19) included *RNF213* and

149    nearby four genes with the intergenic regions (Fig 1). The total length of the BAC, from

150    the BACPAC Resources Center (https://bacpacresources.org), was 213,477 bp.

151    NucleoBond BAC 100 (Macherey-Nagel; Düren, Germany) was used to purify the

152    BAC. The concentration was measured using a NanoPhotometer (Implen). The total

153    amount of 5 μg of BAC was sheared using a Covaris S2 sonicator (Covaris). Because

154    Maricic et al. (2010) recommended that smear DNA band of a gel electrophoresis

155    should be brightest at a size smaller than 1 kbp, and no fragment longer than 5 kbp

156    should be visible, four default settings of the peaks were selected: 150 bp, 300 bp, 500

157    bp, and 800 bp. The seven baits we obtained showed different peaks: 151 bp, 340 bp,

158    456 bp, 492 bp, 522 bp, 619 bp, 735 bp, and 882 bp that were included in the range of

159    that Maricic et al. (2010) showed. The sheared BACs were purified using a MinElute

160    PCR Purification Kit (Qiagen). Subsequently, 1.5-μg sheared BACs per capture reaction

161    were prepared, and the products were then biotinylated according to the protocol used in

162    the previous study [7]. To evaluate the effects of the baits' lengths, the five baits (150

163    bp, 340 bp, 619 bp, 735 bp, and 882 bp) were used; for the number of the captures, two

164    baits (340 bp or 456 bp peak) were used; and for the hybridization temperature, baits

165    that showed a 492 bp peak were used. Baits of 522 bp peak were used for an initial

166    BDC (with PrimeSTAR), and baits of 492 bp peak were used for an additional BDC.

167

## 168    BAC single capture (BSC) for testing experimental conditions

169    This enrichment was conducted according to the protocol of Maricic et al. (2010) using

170    BAC baits. We named it "BAC single capture (BSC)." Concentrations of enriched

171    libraries were assessed using a KAPA Quantification Kit (Kapa Biosystems, Cape

172    Town, South Africa). The size distributions of enriched libraries were verified using a

173    2100 Bioanalyzer (Agilent Technologies; CA, USA). To determine the technical

174    variability in targeted captures, each capture was performed in duplicate.

175

## 176    Initial protocol of BAC double capture (BDC) with

## 177    PrimeSTAR

178    This enrichment was conducted following the modified protocol of the NimbleGen

179    technical note "Double Capture: High Efficiency Sequence Capture of Small Targets for

180    use in SeqCap EZ Library, Applications on 454 Sequencing Systems" (Fig 2). SeqCap

181    EZ Hybridization and Wash Kit (Roche; Basel, Switzerland) were used according to the

182    technical note. The protocol was named, "BAC double capture (BDC) with

10

183    PrimeSTAR." Blocking oligonucleotide solutions and human Cot-1 DNA were added to

184    the library solution. The solution was dried out using a heat block at 95°C.

185    Hybridization buffer and formamide added to the dried DNA, and the mixture was

186    suspended by vortex mixing. The suspended mixture was single-stranded using a heat

187    block at 95°C for 10 min. Biotinylated BAC baits (500 ng) eluted by 4.5 μL PCR grade

188    water was added to the single-stranded DNA mixture and mixed by pipetting. Then the

189    solution was heated in a thermal cycler to 95°C for 10 min and incubated at 65°C

190    overnight (12–16 h). Following incubation, Dynabeads M-270 Streptavidin (Invitrogen;

191    CA, USA) was added to the hybridization mixture. Bound DNA fragments were washed

192    and eluted using NGS MagnaStand (Nippon Genetics, Tokyo, Japan). After the wash

193    and the elution, PCRs were run of the enriched library (26 μL) before removing

194    magnetic beads. The 5 μL of the eluted library was used as a template for PCR in a 50

195    μL solution containing deoxynucleotide (dNTP) 0.2 mM, 0.2 μM of each primer,

196    Sol_bridge_P5 and Sol_bridge_P7 in Maricic et al. (2010), 1.25 U of PrimeSTAR GXL

197    DNA Polymerase (Takara Bio). The 1st post-capture PCR was carried out following the

198    cycling reaction: 16 cycles of denaturation at 98°C for 10 sec, annealing at 60°C for 15

199    sec, extension at 68°C for 50 sec into the plateau phase according to Maricic et al.

200    (2010). The PCR amplicon was purified using a MinElute PCR Purification Kit

11

201    (Qiagen). Then, the 2nd capture was conducted using the enriched and purified library

202    using the same steps as in the 1st capture. After that, the 2nd post-capture PCR of the

203    2nd captured library was run into the plateau phase (20 cycles) using the same cycling

204    condition as in the 1st post-capture PCR. The amplified library was purified using the

205    methods described above. Quantification of the amplified capture library was conducted

206    with a KAPA Library Quantification Kit for Illumina NGS platforms (Kapa Biosystems)

207    and a 2100 Bioanalyzer (Agilent Technologies). To determine the technical variability

208    in targeted captures, each capture was performed in duplicate.

209

## Sequencing for BSC and BDC with PrimeSTAR

211    The enriched libraries by BSC were sequenced on a MiSeq (Illumina; CA, USA) using

212    Illumina MiSeq reagent kit v2 ($2 \times 25$ cycles) or v2 nano ($2 \times 150$ cycles) or v3 ($2 \times 75$

213    cycles). Fastq files were processed using Trimmomatic (version 0.35) in the paired-end

214    palindrome mode to remove TruSeq adapter sequences, low-quality reads (average:

215    <Q20), and nucleotides after the 5'-end from the 26th base and following bases of each

216    read, regardless of quality, to minimize the differences among the three reagent kits.

217          The enriched libraries by BDC with PrimeSTAR were sequenced on a MiSeq

218    (Illumina) using the Illumina MiSeq Reagent Kit v3 ($2 \times 75$ cycles). Fastq files were

12

219     processed using Trimmomatic (version 0.35) in the paired-end palindrome mode to

220     remove TruSeq adapter sequences and low-quality (average: <Q20) reads.

221

## Alignment for BSC and BDC with PrimeSTAR

223     The quality-controlled reads were aligned with the Burrows-Wheeler Aligner (BWA)

224     software version 0.7.12-r1039 [12] to the human genome (GRCh37) with default

225     parameters. Duplicate reads per sample were marked using the MarkDuplicates tool

226     from the Picard software version 1.128 (https://broadinstitute.github.io/picard/) and

227     local realignments around indels were performed on per sample basis using the

228     IndelRealigner tool from the Genome Analysis Toolkit (GATK) software version 3.4-46

229     [13]. The reads mapped to a large tandem repeat, chr17: 78234665-78372586, were

230     removed. Coverage and average depth per sample of targeted regions were calculated

231     using GATK's DepthOfCoverage analysis. The number of mapped and duplicated reads

232     were obtained using SAMtools version1.2 flagstat analysis [14].

233

## Production of indexed libraries for BDC and MB

235     DNA was extracted from bloods of 24 moyamoya disease (MMD) cases collected at

236     Kitasato University Hospital using DNA Extractor WB Kit (Wako Pure Chemical

13

237    Industries; Osaka, Japan). All the patients included in this study provided written

238    informed consent. This project was approved by the ethics committee at Kitasato

239    University School of Medicine. The concentrations of DNA extracts were measured

240    using a NanoPhotometer (Implen) and Qubit 3.0 Fluorometer (ThermoFisher Scientific;

241    MA, USA). Using a Covaris S220 sonicator (Covaris), 2 µg of DNA was sheared. The

242    target peak was set at 300 bp. To produce an indexed library using NEBNext Ultra DNA

243    Library Prep Kit and Multiplex Oligos for Illumina (New England BioLabs), 500 ng of

244    DNA was used. Sheared DNA was end-repaired, dA-tailed, ligated to Illumina specific

245    adaptors, size selected to an approximate insert size of 400–500 bp by Agencourt

246    AMPure XP beads (Beckman Coulter), and amplified by 6 or 7 cycles of PCR. The

247    libraries were purified using Agencourt AMPure XP beads (Beckman Coulter).

248

## 249    The final protocol of BDC

250    The protocol of BDC with PrimeSTAR was modified with KAPA HiFi DNA

251    Polymerase. The following is the final protocol of BDC. The PCR amplification process

252    of the protocol of BDC was improved with KAPA HiFi DNA Polymerase. DNA

253    libraries of eight MMD cases were used for the final protocol. The 1st and 2nd post-

254    captured libraries were used as templates for PCR in a 50 µL solution containing

14

255    deoxynucleotide (dNTP) 0.3 mM, 0.5 μM of each primer, 1.0 U of KAPA HiFi DNA

256    Polymerase (Kapa Biosystems). PCR was carried out using the following protocol: an

257    initial denaturing step at 98ºC for 2 min, 8 cycles for 1st post-capture libraries or 13

258    cycles for the 2nd post-capture libraries of denaturation at 98ºC for 20 s, annealing at

259    60ºC for 30 s, extension at 72ºC for 40 s, and a final extension step at 72ºC for 5 min.

260    To determine the technical variability in targeted captures, each capture was performed

261    in duplicate.

262

## MYbaits double capture (MB)

264    DNA libraries of 24 MMD cases were used for the targeted capture experiment. The

265    captures were performed using the MYbait Custom Kit (MYcroarray; MI, USA)

266    constructed for enrichment of the same region of the BAC twice, following the

267    manufacturer's instructions (http://www.mycroarray.com/pdf/MYbaits-manual-v3.pdf).

268    The baits of MYbaits were uniquely designed to map to the human reference genome.

269    The designed baits covered nearly 80% of the target region. The libraries were

270    hybridized to half of an aliquot of the RNA baits per reaction. The 1st post-capture PCR

271    was 8 cycles. The 2nd capture was conducted using the whole quantity of the 1st post-

272    capture PCR amplicons with the same protocol as that in the 1st capture. The 2nd post-

15

273    capture PCR was 11 cycles. After the enrichment, the libraries were purified using

274    MinElute PCR Purification kit (Qiagen), quantified using TapeStation (Agilent

275    Technologies) and Qubit 3.0 Fluorometer (ThermoFisher Scientific). The method is

276    called MYbaits double capture (MB).

277

## 278    **Sequencing and alignment for BDC and MB**

279    The enriched libraries were sequenced on a MiSeq (Illumina) using Illumina MiSeq

280    reagent kit v3 ($2 \times 75$ bp chemistry). Fastq files were processed using Trimmomatic

281    (version 0.35) [15] in the paired-end palindrome mode to remove TruSeq adapter

282    sequences and low-quality (average: <Q20) reads. The method of alignment for BDC

283    and MB was the same as for that of BSC and BDC with PrimeSTAR.

284

## 285    **Valiant calling for BDC and MB**

286    The resulting data was analyzed with the GATK version 3.4-46, according to GATK Best

287    Practices recommendations [13,16,17]. Following the guidelines for experiments of

288    small-targeted regions, this workflow included calling variants and producing the

289    genomic variant call format (gVCF) files in target regions individually per subject using

290    a HaplotypeCaller, followed by joint genotyping data to produce a multisample raw VCF

291    file using GenotypeGVCFs. Default settings were used for both tools. After variant

292    calling, the following annotations and thresholds were used to remove low-confidence

293    SNPs, based on GATK recommendations for hard filtering: QD <2.0, FS >60.0,

294    HaplotypeScore >13.0, MQ <40, MQRankSum $<-12.5$, ReadPosRankSum $<-8.0$.

295    Similarly, the following filters were applied to remove low-confidence indels: QD <2.0,

296    FS >200.0, ReadPosRankSum $<-20.0$. We extracted variants information from the

297    filtered VCF file using VCFtools [18].

298

299

## 300    Results

301

## 302    Experimental conditions of BDC

303    The effects of baits lengths were evaluated. BSC was performed with the five baits with

304    different lengths and the on-target rates were graphed (Fig 3). The approximate 350–750

305    bp peak baits showed more stability and higher on-target rates (0.28–0.43%) than did

306    the 150 bp peak bait (0.16% and 0.33%). The rates of the 150 bp peak bait showed a

307    larger difference (0.17%) between the duplicates than did that of the longer baits (0.04–

308    0.08%).

17

309     The effects of numbers of captures were evaluated. BSC and BDC were

310     performed with PrimeSTAR and the on-target rates were graphed (Fig 4). There was

311     one capture of BSC and two of BDC with PrimeSTAR, which showed higher on-target

312     rates (16.53% and 12.21%) than did that of BSC (0.28% and 0.32%). Therefore, two

313     captures were more efficient than one.

314     The effects of hybridization temperature were evaluated. Hybridizations were

315     performed at 45$^{\circ}$C and 65$^{\circ}$C, and the on-target rates were graphed (Fig 5).

316     Hybridization at 65$^{\circ}$C showed higher on-target rates (21.12% and 25.01%) than did

317     those hybridized at 45$^{\circ}$C (3.51% and 5.78%).

318

## Comparison between BDC and a commercial targeted capture

## method

321     PrimeSTAR GXL DNA polymerase for BDC was used in the initial protocol of BDC

322     (see Materials and Methods). However, in the final protocol of BDC, it was converted

323     to KAPA HiFi DNA polymerase (see Materials and Methods) because the polymerase

324     showed more high yield than did the PrimeSTAR GXL DNA polymerase (S1 Figure, S1

325     Protocol). The quality of the NGS data of BDC with KAPA HiFi DNA polymerase and

326     the targeted capture method was compared with that of MB.

18

327    Averages of rates of "unique reads" and "duplicate reads" were calculated (Table 1).

328    "Unique reads" mean reads that mapped uniquely to a reference genome. "Duplicate

329    reads" mean reads that mapped to a reference genome at the same position with the

330    other reads and had the same length and the same variation. When raw NGS datasets are

331    processed, duplicate reads are removed from the dataset. BDC with KAPA HiFi DNA

332    polymerase showed a higher average of duplicate-reads rate (46.9%) than did that of

333    MB (16.1%). The averages of depths between BDC and MB were compared (Table 1).

334    "Depth" means the number of reads that mapped at a genomic position. The averages of

335    total depth calculated by adding depth of each genomic position were almost the same,

336    approximately 30 M. The averages of depth were also almost the same, approximately

337    140. The averages of on-target rates were also compared (Table 1): that of BDC (22.5%)

338    was similar to that of MB (24.3%).

339

340    **Validation of variant sites**

341    The validation of BDC was evaluated using sequence data of 8 samples from MMD

342    cases that were conducted both BDC and MB. The called SNP sites of BDC (572 sites)

343    were larger than those of MB (549 sites) (Table 2). The numbers of SNPs registered in

344    dbSNP were larger in BDC (540 sites) than in MB (517 sites). The number of SNPs not

345 registered in dbSNP of BDC was the same as that of MB. The concordant rates of

346 genotypes of SNP sites between BDC and MB were calculated (Table 3). The SNPs

347 registered in dbSNP were 98.4%. The SNPs registered in dbSNP were 97.3%. In

348 twenty-eight SNPs, at least 1 out of 8 MMD samples were called different genotypes

349 between BDC and MB. Those SNPs were placed at genomic regions where was difficult

350 to map reads and call variants correctly in (poly A or G regions: 7 sites, CNV: 11 sites,

351 retro transposons: 8 sites, low complexity regions: 2 sites).

352

## Comparisons of the average values of between BDC and MB

354 The data qualities and costs of BDC and MB were evaluated comparing nine categories

355 (Table 4). The required genomic DNA for each method was higher weight in BDC (1.5

356 μg) than that in MB (0.5 μg). BDC required baits which we constructed myself, while

357 MB required manufactured baits included in a targeted capture kit. BDC took 10 days to

358 prepare the baits because of the time required to order *E. coli* harboring BAC that covered

359 *RNF213* and amplified and purified it. MB took 60 days from design of the baits to its

360 arrival. The period of a targeted capture experiment in BDC was the same as MB, 3 days.

361 The experimental cost without sequencing using MiSeq of BDC (USD 55) was lower

362 than that of MB (USD 270). We examined four points regarding the quality of the data.

20

363    That of MB was the required library weight. The duplicate read rate of BDC (46.9%) was

364    higher than that of MB (16.1%). The average of depth of BDC (140.5) was the almost

365    same as that of MB (141.9). The on-target rate of BDC (22.5%) was also quite close to

366    that of MB (24.3%). Therefore, the data quality of BDC was close in on MB, except only

367    for the duplicate read rate.

368

369

370    **Discussion**

371

372    We compared the three experimental conditions, the baits lengths, the number of

373    captures, and the hybridization temperatures, and found appropriate conditions

374    exclusively for the BDC method. These might not be the best conditions, but better ones

375    for the method.

376        We first found that the baits of 350–750 bp peak obtained a higher on-target rate

377    than did the peak baits around 150 bp and 900 bp (Fig 3). The magnetic beads

378    (Dynabeads M-270) that immobilize baits reduce the binding capacity for large DNA

379    fragments due to the likelihood of steric hindrance. Twice as many copies of a 500 bp

380    DNA fragment bind to the beads than in the case of a 1,000 bp DNA fragment.

21

381    Therefore, around 900 bp baits would show lower on-target rates. The present study

382    showed that about 150 bp baits were the lowest and most unstable on-target rates.

383    Commercial targeted capture kits have uniformly about 100 bp baits (e.g., Agilent: 120

384    bp, Illumina: 100 bp, MYcroarray: 80–120 bp). The protocols of such kits optimized

385    many conditions (e.g., bait lengths, the bait densities, the hybridization temperatures,

386    and the hybridization reagents). Those results indicate that the experimental conditions

387    of the protocol in the present study are not applicable for the short baits. To optimize the

388    baits lengths for BDC, we should examine more patterns of baits lengths in the next

389    step.

390        We observed much higher on-target rates with double captures than that with a

391    single capture (Fig 4), suggesting that the double capture definitely enriches more target

392    libraries. The targeted captures of small regions are especially more difficult than those

393    of larger regions, because genome libraries have smaller volumes of narrow targeted

394    regions (e.g., hundreds kbp) than those of broad targeted regions (e.g., several Mbp).

395    Double captures seem to enable researchers to enrich more DNA of interest than do

396    single captures.

397        We found that the hybridization temperature at 65℃ was suitable for effective

398    library enrichment (Fig 5), suggesting that the hybridization at 65℃ gave higher

22

399 specificity than that at 45°C. From these results, overall, we would propose that the

400 conditions are the baits with 350–750 bp peak, twice capturing, and the 65°C

401 hybridization temperature.

402  The average of duplicate read rate of BDC was higher than that of MB (Table 1).

403 In the BDC protocol, PCR includes one more cycle than that in the MB protocol for

404 obtaining libraries of adequate quantities. This could be the reason why the higher rate

405 of duplication in BDC was observed. It might be important for improving BDC to

406 reduce the number of PCR cycles.

407  The on-target rate of BDC (Table 1) was higher than the that of the previous

408 manual methods that enriched mtDNA using baits constructed from long-range PCR

409 amplicons or genomic regions using BAC-based baits [7,19]. A previous study [7]

410 proposed a single capture method. Another previous study [19] also proposed a single

411 capture method and used non-sheared BACs that were affected by their own steric

412 hindrance. Thus, we would claim that BDC is improved from the previous manual

413 capture methods.

414  BDC showed more called SNP sites than MB (Table 2). The baits used in MB

415 were synthetic oligonucleotide probes. Non-unique baits in a human-reference-genome

416 sequence were excluded in the probes designed. Therefore, the baits of MB covered

23

417     approximately 80% of the target region. On the other hand, the baits used in BDC were

418     constructed from BAC, which covered the whole target regions. The differences of

419     processes constructing probes between BDC and MB could affect the numbers of

420     enriched libraries and the called SNP sites.

421         The concordance of SNP genotypes between BDC and MB was >97% (Table 3).

422     The SNPs that were called the different genotypes placed at poly A, poly G, repetitive

423     regions and transposons. Those genomic regions that showed low uniqueness in

424     genomes are, in general, more likely to have errors in PCRs and mapping reads [20].

425     Thus, the SNP sites that showed discordance of genotypes locate genomic regions

426     where were difficult to call SNPs correctly.

427         BDC allows efficient capture of the genomic library of NGS for large genes. First,

428     the approach is cost-effective in that it only requires standard laboratory equipment and

429     reagents that cost USD55 per reaction (Table 4). Second, it is fast, i.e., it enables

430     researchers to perform captures immediately without designing and constructing baits

431     (Table 4). Third, the on-target rates are almost the same for BDC and MB (Tables 1 and

432     4). Those features enable more laboratories to start easily targeted captures. More

433     adjustments of capture conditions make a better BDC. When BDC is used for the other

434     human genomic regions and any organism's genome, the conditions of targeted captures

24

435     may require adjustment. E.g., if the GC (guanine and cytosine) contents of targeted

436     regions are different from those of the BAC used in the present study, suitable

437     hybridization temperatures can be changed. BDC enables the recovery of targeted

438     genomic regions like a large gene from most such ancient samples. Because a large

439     amount of non-targeted DNA and bacterial DNA extracted from those bones, those

440     samples are needed to retrieve only endogenous genomes. Our new approach helps to

441     conduct paleogenomic studies.

442

443

## 444     **Acknowledgments**

445

453     Tokyo) gave us equipment making a genome library for next-generation sequencing. Dr.

454     Alison Devault at MYcroarray provided a lot of advice to construct the bait set. And we

455     thank Robert E. Brandt, Founder, CEO, and CME, of MedEd Japan, for editing and

456     formatting the manuscript.

457

458

## 459 **References**

460     1.    Zhang J, Chiodini R, Badr A, Zhang G. The impact of next-generation

461           sequencing on genomics. J Genet Genomics. Elsevier Limited and Science Press;

462           2011;38: 95–109. doi:10.1016/j.jgg.2011.02.003

463     2.    Shimada J, Mogi S. Necessity of Research and Development of the Structural

464           Genome Analysis Technology. Sci Technol Trends Q Rev. 2004;10: 11–18.

465     3.    Carpenter ML, Buenrostro JD, Valdiosera C, Schroeder H, Allentoft ME, Sikora

466           M, et al. Pulling out the 1%: Whole-Genome Capture for the Targeted

467           Enrichment of Ancient DNA Sequencing Libraries. Am J Hum Genet. 2013;93:

468           852–864. doi:10.1016/j.ajhg.2013.10.002

469     4.    Bainbridge MN, Wang M, Burgess DL, Kovar C, Rodesch MJ, D'Ascenzo M, et

470           al. Whole exome capture in solution with 3 Gbp of data. Genome Biol. 2010;11:

471    R62. doi:10.1186/gb-2010-11-6-r62

472    5.    Briggs AW, Good JM, Green RE, Krause J, Maricic T, Stenzel U, et al. Targeted

473          Retrieval and Analysis of Five Neandertal mtDNA Genomes. Science. 2009;325:

474          318–321. doi:10.1126/science.1174462

475    6.    Briggs AW, Good JM, Green RE, Krause J, Maricic T, Stenzel U, et al. Primer

476          Extension Capture: Targeted Sequence Retrieval from Heavily Degraded DNA

477          Sources. J Vis Exp. 2009; 1573. doi:10.3791/1573

478    7.    Maricic T, Whitten M, Pääbo S. Multiplexed DNA Sequence Capture of

479          Mitochondrial Genomes Using PCR Products. Fleischer RC, editor. PLoS One.

480          Public Library of Science; 2010;5: e14004. doi:10.1371/journal.pone.0014004

481    8.    Shizuya H, Birren B, Kim UJ, Mancino V, Slepak T, Tachiiri Y, et al. Cloning

482          and stable maintenance of 300-kilobase-pair fragments of human DNA in

483          Escherichia coli using an F-factor-based vector. Proc Natl Acad Sci U S A.

484          1992;89: 8794–8797. doi:10.1073/pnas.89.18.8794

485    9.    Osoegawa K, Woon PY, Zhao B, Frengen E, Tateno M, Catanese JJ, et al. An

486          improved approach for construction of bacterial artificial chromosome libraries.

487          Genomics. 1998;52: 1–8. doi:10.1006/geno.1998.5423

488    10.   Osoegawa K, Mammoser AG, Wu C, Frengen E, Zeng C, Catanese JJ, et al. A

27

489     bacterial artificial chromosome library for sequencing the complete human

490     genome. Genome Res. 2001;11: 483–496. doi:10.1101/gr.169601

491   11.  Cheung VG, Nowak N, Jang W, Kirsch IR, Zhao S, Chen XN, et al. Integration

492     of cytogenetic landmarks into the draft sequence of the human genome. Nature.

493     2001;409: 953–958. doi:10.1038/35057192

494   12.  Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler

495     transform. Bioinformatics. 2010;26: 589–595. doi:10.1093/bioinformatics/btp698

496   13.  McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al.

497     The Genome Analysis Toolkit: A MapReduce framework for analyzing next-

498     generation DNA sequencing data. Genome Res. 2010;20: 1297–1303.

499     doi:10.1101/gr.107524.110

500   14.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The

501     Sequence Alignment/Map format and SAMtools. Bioinformatics. 2009;25:

502     2078–2079. doi:10.1093/bioinformatics/btp352

503   15.  Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina

504     sequence data. Bioinformatics. 2014;30: 2114–2120.

505     doi:10.1093/bioinformatics/btu170

506   16.  DePristo MA, Banks E, Poplin R, Garimella K V, Maguire JR, Hartl C, et al. A

507  framework for variation discovery and genotyping using next-generation DNA

508  sequencing data. Nat Genet. Nature Publishing Group; 2011;43: 491–498.

509  doi:10.1038/ng.806

510 17. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, del Angel G, Levy-

511  Moonshine A, et al. From FastQ Data to High-Confidence Variant Calls: The

512  Genome Analysis Toolkit Best Practices Pipeline. Current Protocols in

513  Bioinformatics. Hoboken, NJ, USA: John Wiley & Sons, Inc.; 2013. p. 11.10.1-

514  11.10.33. doi:10.1002/0471250953.bi1110s43

515 18. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The

516  variant call format and VCFtools. Bioinformatics. 2011;27: 2156–2158.

517  doi:10.1093/bioinformatics/btr330

518 19. Alvarado DM, Yang P, Druley TE, Lovett M, Gurnett C a. Multiplexed direct

519  genomic selection (MDiGS): a pooled BAC capture approach for highly accurate

520  CNV and SNP/INDEL detection. Nucleic Acids Res. 2014;42: e82.

521  doi:10.1093/nar/gku218

522 20. Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing:

523  computational challenges and solutions. Nat Rev Genet. 2012;46: 36–46.

524  doi:10.1038/nrg3164

525

526

527

528

# Figure legends

**Fig 1. A genomic position of the BAC (ID: CH17-24F19).** The BAC contains five genes: *caspase recruitment domain family member 17* (*CARD17*), *solute carrier family 26 member 11* (*SLC26A11*), *ring finger protein 213* (*RNF213*), *N-sulfoglucosamine sulfohydrolase* (*SGSH*), *CTD-2047H16.4* (uncharacterized gene).

**Fig 2. An overview of the BAC double capture (BDC) method, which we modified (Maricic et al., 2010).** On the left, the bait construction from the BAC is shown; on the right, the production of indexed libraries that are used in the library enrichment (center). Those colored light red are the BAC-based baits, dark red represents targeted DNA molecules in the libraries, black represents non-targeted DNA molecules in the libraries, green and pink represent indexes, gray represents adapters, and blue and yellow represent biotinylated adapters. Thick lines represent double stranded DNA, and thin lines represent single stranded DNA.

**Fig 3. On-target rates depending on BAC baits length.** On-target rate equals the reads mapped to the target region divided by the reads mapped to the whole reference genome.

**Fig 4. On-target rates depending on the number of captures.** The formula to calculate the on-target rate is the same as in Fig 3.

**Fig 5. On-target rates of the hybridization temperature.** The formula to calculate the on-target rate is the same as in Fig 3.
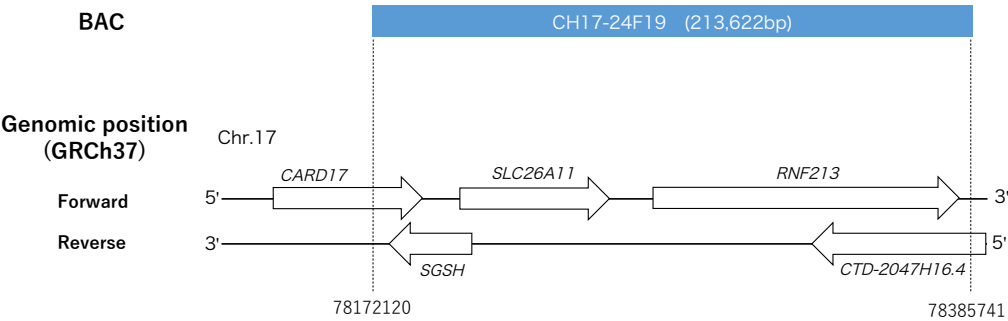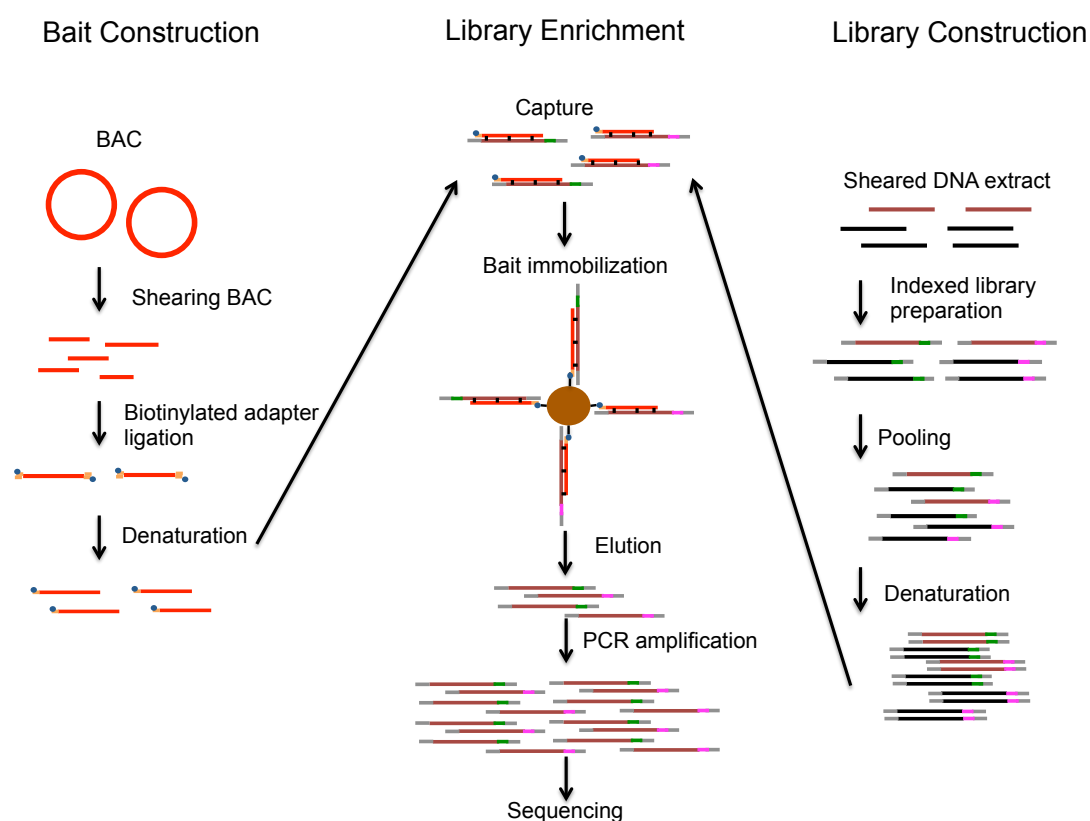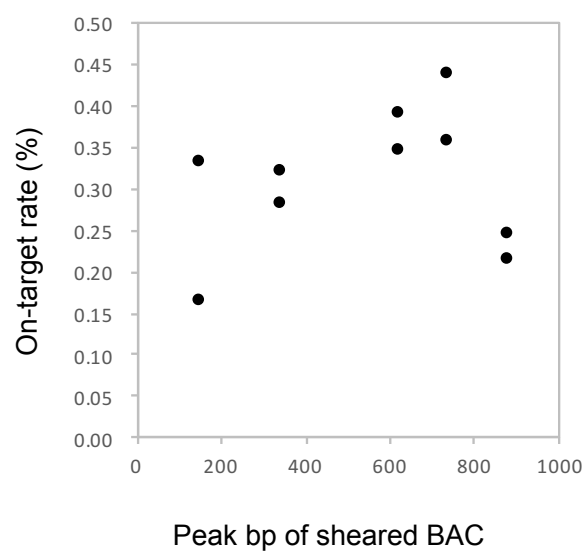
31

557 **Figures and Tables**



558

559 **Fig 1**

560

561

562    **Fig 2**

563

564

565    **Fig 3**

566

567

568    **Fig 4**

569

35

570

571 **Fig. 5**

572

573 **Table 1.  Average numbers of BDC and MB.**

| Method | BDC | MB |
|---|---|---|
| **Duplicate reads (%)** | 46.9 | 16.1 |
| **Unique reads (%)** | 53.1 | 83.9 |
| **Total depth** | 30,020,317 | 30,321,836 |
| **Depth** | 140.5 | 141.9 |
| **On-target rate (%)** | 22.5 | 24.3 |

574

575

576 **Table 2.   Numbers of validated SNPs.**

| Method | All SNPs | dbSNP (registered) | dbSNP (non-registered) |
|--------|----------|--------------------|------------------------|
| **BDC** | 572 | 540 | 32 |
| **MB** | 549 | 517 | 32 |

577

578

579 **Table 3.   Concordant rates of genotypes.**

| | **dbSNP (registered)** | **dbSNP (non-registered)** |
|---|---|---|
| **Number of sites*** | 478 | 32 |
| **Concordant rate of Genotypes** | 98.4% | 97.3% |

580 **\***Sites that genotyped all eight samples with both BDC and MB.

581

582 **Table 4.** **Summary of comparisons between BDC and MB.**

| Category | | Method | |
|---|---|---|---|
| **1** | **2** | **BDC** | **MB** |
| **Cost** | **Required library weight (μg)** | 1.5 | 0.5 |
| | **Baits construction** | Self-making | Outsource |
| | **Preparation period of baits (day)** | 10 | 60 |
| | **Period of targeted capture (day)** | 3 | 3 |
| | **Cost per reaction (USD)** | 55 | 270 |
| **Data quality** | **Duplicate reads (%)** | 46.9 | 16.1 |
| | **Average of depth** | 140.5 | 141.9 |
| | **On-target rate (%)** | 22.5 | 24.3 |
| | **SNP concordance* (%)** | 98.4 | |

583 *A calculated value from SNPs registered in dbSNP138

584

# Supporting information

585

586

587 **Supplementary Figure Legends**

588

589 **S1 Fig. Comparison of PCR efficiency between PrimeSTAR GXL DNA Polymerase**

590 **and KAPA HiFi DNA Polymerase.**

591

592

593

594    **S1 Fig**

595

596

597     **S1 Protocol    Amplification efficiency of two polymerases.**

598     A 2 μL of 1st post-captured library solution before removing magnetic beads

599     concentrated using a BDC method was used as a template for PCR in a 20 μL solution

600     containing 0.5 U of PrimeSTAR GXL DNA Polymerase (Takara Bio), deoxynucleotide

601     (dNTP) 0.2 mM, 0.2 μM of each primer, Sol_bridge_P5 and Sol_bridge_P7 in Maricic

602     et al. (2010). The PCR and purification were carried out using the same method as BDC

603     with PrimeSTAR. The same volume of the 1st post-captured library solution was used

604     as a template for PCR in a 20 μL solution containing, 0.4 U of KAPA HiFi DNA

605     Polymerase (Kapa Biosystems), deoxynucleotide (dNTP) 0.3 mM, 0.5 μM of each

606     primer, Sol_bridge_P5 and Sol_bridge_P7 in Maricic et al. (2010). PCR was carried out

607     using the following protocol: an initial denaturing step at 98°C for 2 min, 16 cycles for

608     the 1st post-capture library of denaturation at 98°C for 20 s, annealing at 60°C for 30 s,

609     extension at 72°C for 45 s, and a final extension step at 72°C for 5 min. To determine

610     the technical variability in targeted captures, each PCR was performed in duplicate. The

611     PCR amplicons were quantified using Qubit 3.0 Fluorometer (ThermoFisher Scientific).

612

613

43