1  **Single molecule targeted sequencing for cancer gene mutation detection**

2  Yan Gao[1], Liwei Deng[1], Qin Yan[1], Yongqian Gao[2], Zengding Wu[1], Jinsen Cai[1], Daorui Ji[1],

3  Gailing Li[1], Ping Wu[1], Huan Jin[1], Luyang Zhao[3], Song Liu[4], Michael W. Deem[5], Jiankui He[6,1,*]

4  [1]Direct Genomics Co., Ltd. Shenzhen, Guangdong, China
5  [2]Institute of Advanced Materials, Nanjing Tech University, Nanjing, Jiangsu, China
6  [3]Chemistry Department, North Carolina State University, Raleigh, NC, USA
7  [4]Clinical Medical Research Center, the Second Clinical Medical College of Jinan University
8  (Shenzhen People's Hospital), Shenzhen, Guangdong, China
9  [5]Departments of Bioengineering and Physics & Astronomy, Rice University, Houston, TX, USA
10  [6]Department of Biology, South University of Science and Technology of China, Shenzhen,
11  Guangdong, China

12  [*]Corresponding author, Email: hejk@sustc.edu.cn

## Abstract

14  With the rapid decline cost of sequencing, it is now clinically affordable to examine multiple
15  genes in a single disease-targeted test using next generation sequencing. Current targeted
16  sequencing methods require a separate step of targeted capture enrichment during sample
17  preparation before sequencing, and the library preparation process is labor intensive and time
18  consuming. Here, we introduced an amplification-free Single Molecule Targeted Sequencing
19  (SMTS) technology, which combined targeted capture and sequencing in one step. We
20  demonstrated that this technology can detect low-frequency mutations of cancer genes. SMTS
21  has several advantages, namely that it requires little sample preparation and avoids biases and
22  errors introduced by PCR reaction. SMTS can be applied in cancer gene mutation detection,
23  inherited condition screening and noninvasive prenatal diagnosis.

## Introduction

25  In the past few years, the cost of large-scale DNA sequencing has been dramatically driven down
26  by the tremendous advances in next-generation sequencing (NGS)[1]. Nonetheless, the cost of
27  human whole genome sequencing and bioinformatics interpretation is still significant. In clinical
28  practice, NGS is used to examine specific gene panels such as cancer genes and inherited
29  conditions, sample numbers are high and data volume per sample is relatively small. It is often
30  more cost-effective and time-efficient to target, capture, and sequence only the genomic regions
31  of interest[2]. For example, there are several cancer gene panels commercially available, targeting
32  as few as 50 to many hundreds of genes that are frequently mutated in cancer patients[3]. The
33  cancer gene panel targeted sequencing has been proved to be useful in hereditary cancers
34  diagnosis, and disease management.

35  Current NGS based targeted sequencing methods require a separate step of capture enrichment
36  during sample preparation before sequencing[4, 5]. The two most commonly used custom-capture

37  methods are based on hybridization or on highly multiplexed PCR. In the solution-based
38  hybridization method, biotinylated DNA or RNA complementary probes are designed bind to
39  gene targets, which are then purified using streptavidin-labeled magnetic beads. In the
40  multiplexed PCR method, hundreds or thousands of PCR primer pairs are mixed to amplify the
41  targeted genes.

42  In this report, we demonstrated a technology and platform to perform Single Molecule Targeted
43  Sequencing (SMTS), which combined targeted capture and sequencing in one step. We used a
44  combination of Total Internal Reflection Fluorescence (TIRF) microscope and single molecule
45  fluorescence dyes to reject unwanted background noise and get single molecule resolution
46  images[6]. The gene-specific flow cell was constructed with capture primers for gene regions of
47  interest and the target genes can thus be sequenced without copying the DNA or enrichment
48  before sequencing. Compared to current targeted sequencing methods with separate capture steps,
49  SMTS has significant advantages, including little sample preparation and avoidance of biases
50  and errors introduced by PCR amplification[7]. SMTS can be applied in cancer gene mutation
51  detection, inherited condition screening, and high-resolution human leukocyte antigen (HLA)
52  typing.

## Results

### *Single molecule detection*

55  The fundamental limitation of detection of single molecule fluorescence signals stems from the
56  intrinsic qualities of the fluorophore. The key challenge is to reduce the background interference,
57  which may arise from Raleigh scattering, Raman scattering, and contaminant fluorescence.
58  Various single-molecule fluorescence microscopy techniques have been developed in the last
59  two decades to overcome the difficulty in detecting single molecules with high signal to noise
60  ratios in the presence of optical background[8].

61  We applied Total Internal Reflection Fluorescence (TIRF) microscopy in this study. The optical
62  setup is shown in Fig. 1. When light strikes an interface going from coverslip glass to fluid in the
63  flow cell chamber at an angle greater than a critical angle, it undergoes a total internal reflection.
64  This generates an exponentially decaying light field called the "evanescent wave" above the
65  surface of glass. The evanescent wave excites fluorescent molecules within about 150-200
66  nanometers of the surface. The fluorescence from the labeled DNA molecules anchored on the
67  glass surface is detected through a microscope objective and fluorescence filters by high
68  sensitivity Electron-Multiplying CCD (EMCCD) cameras. As only the vicinity of the surface is
69  illuminated, the noise from the bulk fluids of flow cell chamber is dramatically reduced. Single
70  DNA molecules anchored on the surface can thus be monitored with high signal to noise (Fig. S1,
71  S2 and S3).

72    The choice of fluorescent dyes to label nucleotides is also critical for single molecule detection.
73    Many common fluorescent labels show rather low photostability if high-intensity laser excitation
74    is used and processes are to be observed over long periods of time. We choose the ATTO 647N
75    dyes to label the nucleotides, which fluoresces twice as strong as cyanine 5 in aqueous solution.
76    Meanwhile, we optimized the imaging buffer to increase the photostability up to five times (Fig.
77    S5).

78    Single-step photobleaching is used as a quality control to distinguish single molecule from
79    multiple molecules. In an ideal situation, each DNA molecule is separately binding to the flow
80    cell surfaces and the minimal distance between two DNA molecules is larger than the diffraction
81    limit of light. In a random attachment cenari(as used in the present study) drive by Poisson
82    statistics,  two or more DNA molecules may bind to the surface at a distance less than the
83    Rayleigh criterion. We quantified the amount of single DNA molecules to aggregated DNA
84    molecules binding to the surface by observing the photobleaching patterns. The single molecules
85    photobleached in single steps, while aggregated molecules photobleached in multiple steps (Fig.
86    2). We observed that 38% of spots are real single molecules, where 36% of spots are aggregated
87    molecules. Only the sequences from the real single molecule spots will be used for analysis.

88    *Targeted hybridization and sequencing*

89    The EGFR, KRAS, BRAF genes were selected for sequencing in this studies. In particular, we
90    aimed to sequence the 8 genetic variants that are related to drug response, including six point
91    mutations and two short deletions (Table 1). Eight capture probe sequences were designed in the
92    upstream of drug response related mutations. The capture probes are synthesized and anchored to
93    the flow cell surface by a expoxy-NH2 bond. We synthesized two sets of target DNA templates
94    for sequencing. The first set was wild type sequence and the second set contained mutations and
95    short deletions (Table 1). Each target DNA template contained a Cy3 fluorescence dye at the 3'
96    end. Excitation of 3' Cy3 fluorescent dyes was used to mark positions of annealed templates on
97    the flow cell surfaces. Synthetic target DNA templates were hybridized to the flow cell with
98    surface-attached capture probes (Fig 3a).

99    The sequencing reaction began with locating the target DNA templates, which are randomly
100   hybridized to capture probes (Fig. 3a). The Cy3 fluorescent dyes attached to target DNA
101   templates are excited by a 532nm green laser and the images were collected to locate the
102   positions of target DNA templates. Then, disulfide linked Atto647N labeled reversible
103   terminators and DNA polymerases were added to the flow cell. The reversible terminators were
104   nucleotide analogs modified to contain a cleavable liner, which allowed only one reversible
105   terminator to be incorporated into the DNA molecule at one time. The polymerase synthesis
106   reaction was carried out at temperature 37 ℃, with one of four types of reversible terminators and
107   necessary cofactors.  Unincorporated reversible terminators were washed way. The Atto647N
108   dyes are excited by a 640nm red laser in an optimized imaging buffer mixture with oxygen
109   scavenging, free radical scavenging, and triplet quenching components. The images were

110    processed using a custom written computer program to automatically locate the spot, determine
111    image noise, and filter out false-positive spots. After imaging, the Atto647N fluorescence dyes
112    were cleaved from the reversible terminators, and the system is ready for a second round of
113    adding reversible terminators and polymerases. The sequencing cycle are repeated many times to
114    achieve the desired length of read (Fig. 3b).

115    *Sequencing coverage depth*

116    To demonstrate the performance of SMTS, we sequenced the wild-type EGFR/KRAS/BRAF
117    DNA templates. The synthesized DNA templates were hybridized to the flow cell with surface-
118    attached capture probes. We sequenced DNA for 19-30 cycles, which enable to cover all
119    mutation/deletion loci. 300 fields of view were imaged for each cycle. In each field of view,
120    there are approximate 2200-2500 reads on average. The sequencing reads were aligned to
121    reference sequences with customized program of Smith-Waterman algorithm (Table 2).   We
122    observed that the coverage depth varies among different DNA templates (Figure 4a). The
123    possible explanation is that the hybridization efficiency for DNA templates is sequence-
124    dependent and the secondary structures that involve the target region can also affect
125    hybridization efficiency. The average coverage depth was 1954-fold. Higher coverage depth can
126    be achieved by capturing images for more fields of view.

127    *Sequencing accuracy*

128    The accuracy was calculated by comparing the reference sequences with the consensus
129    sequences. Consensus sequences were calculated as the most frequent bases at each position in
130    the sequence alignment (Table 2). By comparing the consensus sequence to the reference
131    sequence base-by-base, the consensus sequence is 100% identical to the reference sequence in
132    our four repeated experiments. We performed sampling-subsampling to the sequence data to get
133    low-coverage data, and recalculated the consensus sequences at different coverage depth. If each
134    base was covered only one time, which means the coverage depth is 1 fold, the accuracy was 95%
135    on average. If each base was covered with 5 times or more on average, the consensus accuracy is
136    approaching 100% accuracy (Fig. 4c). We performed multiple repeated experiments to estimate
137    the errors in the raw sequencing data. The reads from each template were separately aligned to
138    the DNA reference. Each position in the reference was mapped by multiple reads. The error rate
139    of a position was the ratio of reads disagreeing with the reference divided by the total number of
140    reads mapped to the reference. The overall error rate was an average of error rate of all positions.
141    The error of raw sequencing reads was dominated by deletion (Fig. 4b). The substitution error is
142    relatively small, in four repeated experiment, the average substitution rate is 0.52% per base (Fig.
143    S6).

144    *Detecting low frequency of mutations*

145  The wild type DNA was mixed with mutant type DNA at 10:1 and 97:3 ratios (Table 1). The
146  DNA mixture was hybridized to the flow cell and sequenced. Each raw sequence read was
147  aligned to reference sequences to determine whether it originated from wild type or mutant type
148  DNA. As a control, we also sequenced pure wild type DNA with the same condition. We found
149  that the percentage of mutant DNA detected in the DNA mixture was significantly higher than
150  that in pure wild type DNA control (Fig. 5). In this experiment, SMTS can detect mutant
151  sequences with frequency 3%.

## Discussion

153  We here demonstrated a method of capturing and sequencing DNA in a single step, which
154  provides a much simpler approach to targeted sequencing. We have shown that the mutations
155  and short deletions can be accurately detected at low frequency.

156  We have included several mutations of EGFR/KRAS/BRAF genes in this study. These mutations
157  are actionable and can be therapeutically target. Somatic mutations in EGFR in exon 18, 19, 21
158  and the T790M point mutation in exon 20 are predictive of a clinical response to the EGFR
159  tyrosine kinase inhibitor drugs gefitinib and erlotinib[9, 10]. Somatic mutations in KRAS (codons
160  12, 13) and BRAF (V600E) in colorectal cancer that predict poor prognosis and nonresponse to
161  anti-EGFR antibodies. BRAF V600E is predictive of a positive response to the BRAF V600-
162  specific inhibitor vemurafenib in melanoma[11].

163  SMTS has several advantages over the more traditional Sanger sequencing and other NGS
164  platforms commonly used for the detection of mutations. Firstly, there is little required in the
165  way of sample preparation. Only sonication of the genomic DNA is needed. In the case of
166  nucleic acids from sources such as FFPE or cfDNA, it is possible that even the sonication is not
167  needed. Other high throughput sequencing technology such as Illumina requires days of labor
168  work on sample preparation, which contains multiple steps such as sonication, end repairing, dA
169  tailing, adaptor ligation, PCR amplification and target enrichment. Therefore, the SMTS
170  technology has the potenial of  reducing cost, turn-around time and the risk of errors in sample
171  handling. Secondly, SMTS technology directly sequences original individual molecules, not
172  PCR products. This should provide increased sensitivity for the detection of low prevalence
173  mutations and avoid PCR biases[12], which are essential features in the sequencing of a
174  heterogeneous cancer sample[13].

175  We observed that the coverage depth was not uniform among different positions. Some
176  sequences appeared to be difficult to be sequenced. The uniformity of coverage could be
177  improved by carefully designing the capture probes, in particularly, to avoid the secondary
178  structure. We also observed that only one third of fluorescence spots were from single molecules.
179  Under the random attachment scenario described in this study, a large portion of spots came from
180  two or more molecules binding closer than the diffraction limited resolution of the system. The
181  ratio of single molecules could be increased by optimizing the hybridization condition and/or

182  controlling the density of capture probes. The overall error rate of raw sequences was still
183  significant [14]. To reduce the error rate, we need to further optimize the chemical reaction
184  conditions for incorporating reversible terminators and cleaving the fluorescence dye after
185  imaging. Meanwhile, by modeling the error profiling, a better base calling algorithm could be
186  developed. The four reversible terminators (A, T, C and G) used in current study were labeled
187  with the same fluorescence dye. In future, we can modify the reversible terminators and label
188  each of four nucleotides with unique fluorescence dyes[15]. By doing so, the speed and accuracy
189  will be improved.

190  For the foreseeable future, the high cost and complexity of data analysis will limit the application
191  of whole-genome sequencing for the detection of mutations in a clinical setting. Targeted
192  resequencing of areas of interest will therefore remain key to determining mutational status.
193  SMTS is a stride forward in putting this into practice. Although currently only a few loci of a
194  few genes are screened, there is clearly scope for the creation of multi-gene capture arrays,
195  allowing large numbers of loci to be analyzed rapidly and cost-effectively with low DNA input
196  requirements. The single-step capturing and sequencing whole exome is also possible in future.
197  In its simplicity, this approach provides an opportunity to truly begin integrating the vast
198  quantity of genomic data generated in this next-generation era with clinical practice.

199

## Methods

200  ### Methods

201  *Optical setup*

202  A custom-engineered sequencer prototype contained a Total Internal Reflection Fluorescence
203  (TIRF) microscope with 60X oil objective (Nikon Ti-E, Japan), EMCCD camera with a
204  resolution of 512X512 (Andor, Belfasst, UK) and 2 color laser powers, 532nm (100mW) and
205  640nm (100mW). A motorized stage (ASI, Eugene, OR) was installed on the TIRF microscope
206  to hold and control the motion of the flow cell (Bioptechs, Bulter, PA) during sequencing. The
207  heater (Bioptechs, Bulter, PA) for both flow cell and objective was installed and can maintain the
208  temperature in chamber of the flow cell at 37 ℃.

209  *Flow cell and liquid handing*

210  The FCS2 flow cell contained the chemical functionalized coverslip with epoxy layer (Schott,
211  Jena, Germany), 0.175mm thick and 40mm in diameter. A gasket was assembled between the
212  coverslip and an aqueduct slide which forms the chamber (3mm X 23 mm X 0.25 mm) for
213  chemical reaction. The sandwiched structure part was fixed by a top with stainless steel tube
214  inside (inlet port and outlet port) and metal base. A Titan EZ valve with 12 channels (IDEX
215  Health & Science, Oak Harbor, WA, USA) was connected between the inlet of the flow cell and

216    sequencing reagents. The outlet of the flow cell was connected with a syringe pump (Tecan,
217    Männedorf, Swiss) to drive the fluidic in the system by suction.

*Surface chemistry*

219    Synthesized capture probes (oligonucleotides) were covalently coupled to the epoxy coated
220    coverslip surface.  The capture probes were firstly incubated at 95℃, then the coverslip was
221    immerged into a capture probe solution at 1 nM in 150mM $K_2HPO4$, pH 8.5 at 37℃ for 2 hours.
222    Then the coverslip was rinsed by 3X SSC with 0.1% Triton X-100 and 3X SSC, 150mM
223    $K_2HPO4$, pH 8.5 in sequence.

*Imaging processing*

225    Images are processed using a custom written spot localization algorithm (Fig. S4). Firstly, stage
226    drifts between different imaging cycles were corrected by calculating the peak position of two
227    images by Phase-Only Correlation (POC) function. After correcting all cycles with the
228    corresponding first cycle, the corrected images were convolved with a Gaussian kernel. The
229    correlation images were then subjected to the threshold determined by the noise measurement on
230    those images. All contiguous groups of pixels above the threshold were grouped as spots. After
231    that, each spot was fitted with a Gaussian function. This step allowed an accurate determination
232    of the centroid position for single molecules and both members of closely standing molecule
233    pairs. At the same time, clusters of three or more molecules were filtered out. A spot that
234    appeared twice at a same time point but under different wavelength lasers was considered as a
235    base incorporation event. Thus, the spot was renamed as an incorporation spot and marked on the
236    incorporation image. A set of incorporation spot centroids falling within a 1.6 pixel radius is
237    called a "track". Comparing with the order of adding reversible terminators, these "tracks" were
238    converted to the final sequences on the position of each incorporation spot.

*Target template of EGFR/KRAS/BRAF*

240    Eight mutation sites in three genes（EGFR, KRAS and BRAF）were covered by the target
241    templates, including six point mutations（G719A in EGFR exon 18, T790M in EGFR exon 20,
242    L858R and L861Q in EGFR exon 21, G12S and G13D in KRAS exon 2 and V600E in BRAF
243    exon 15）  and two short deletions（ΔE746-A750 deletions and ΔE747-A753 deletions in
244    EGFR exon 19）. We designed two target sequences for each genetic variant, which are wild
245    type and mutant type. The length of each target template was 70 bp, with a Cy3 fluorescence dye
246    attached to the 3' end. Synthetic target templates were hybridized with capture probes attached
247    on the surface of flow cell according to complementary matching principle.

*Capture probe design*

249    A 60nt capture probe sequence with 10 dT bases and an amine labeled 5' end was designed
250    according to the upstream gene sequence of mutation sites. The 50nt target-specific sequence at

251    the 3' end of capture probe sequence were designed according to the program BatchPrimer3, with

252    specified conditions: 20%-80% GC and Tm's >65℃. Capture probes and target templates were

253    synthesized by Sangon Biotech(Shanghai).

254    *Reversible terminators*

255    The modified reversible terminators are composed of nucleotide triphosphates, modified with a

256    detectable label (Atto647N) by disulfide linker and an inhibitor group(SeqLL, Woburn, MA,

257    USA). The inhibitor region has multiple negative charged groups (carboxyl group) allowing

258    incorporation of one nucleotide into the DNA duplex while prohibiting the second or third or

259    more nucleotide incorporation. The detectable label and inhibitor group were cleavable.

260    *Sequencing cycle*

261    The coverslip was incubated in synthesized templates labeled with Cy3 solution at 5nM in 3X

262    SSC, pH7, at 55 ℃ for 2 hours to form a DNA duplex. Then the surface was rinsed with 150mM

263    HEPES, 1X SSC and 0.1% SDS, followed by 150mM HEPES and 150mM NaCl. Finally the

264    coverslip was assembled into the follow cell.

265    The sequencing process was controlled automatically by the fluidic system. Two different types

266    of reagents containing nine pre-prepared reagents were used and stored at different temperatures.

267    One type is the chemical or biochemical reaction reagents, including four nucleotide (dNTP-

268    Atto647N) and DNA polymerase mixtures, cleavage reagent (TCEP, 50mM), cap reagent

269    (50mM idoacetamide), and imaging buffer (50mM Trolox, 20mM glucose and 5mM glucose

270    oxidase in HEPES buffer) stored at 4℃. The other is rinse buffer including rinse buffer 1

271    (150mM HEPES, 1X SSC and 0.1% SDS, pH 7.0) and rinse buffer 2 (150mM HEPES and

272    150mM NaCl, pH 7.0) stored at room temperature.

273    First, 0.25 μM reversible terminators (one of G, C, T and A) and 20nM polymerase mixture was

274    introduced into the flow cell, incubated for 4 minutes at 37 ℃ and washed out by rinse buffer1

275    and 2. Then imaging buffer (50mM Trolox, 20mM glucose and 5mM glucose oxidase in HEPES

276    buffer) was pumped in the flow cell. Then, the images of 300FOVs were taken. Typically, 4

277    exposures of 0.1 second were taken in each field of view (FOV, 54.6μm ×54.6μm). After

278    imaging, the flow cell was washed by rinse buffer. The cleave reagent was introduced into the

279    flow cell and reacted for 5 minutes flowed by the cap reagent under reaction for another 5

280    minutes. Finally the flow cell was washed by rinse buffer and finished the first cycle of

281    sequencing. The sequencing cycle was repeated with the same procedure, except changing the

282    reversible terminators. In this paper, the terminators were added into the system as the repeated

283    order of G, C, T, A.

284    *Bioinformatics*

285    Quality control on the sequence reads was first performed. Firstly, reads with length less than 5

286    bases were filtered out. Then, sequencing reads that appeared less than 4 times were filtered out.

287    Secondly, sequencing reads that could not be aligned to reference sequences are not included for

288    further analysis.

289    The alignment described above was performed with Smith-Waterman algorithm, which performs

290    local sequence alignment. By using a custom definition scoring system (which included the

291    substitution matrix and the gap-scoring scheme), the chosen algorithm could guarantee

292    identification find of the optimal local alignment. In this setup, the penalty for a deletion in a

293    read was -1, for an insertion -1, for a match 2, and for a substitution -2.

## References

294

295    1.      Metzker, M.L. Sequencing technologies - the next generation. *Nature reviews. Genetics* **11**, 31-
296           46 (2010).

297    2.      Dewey, F.E., Pan, S., Wheeler, M.T., Quake, S.R. & Ashley, E.A. DNA sequencing: clinical
298           applications of new DNA sequencing technologies. *Circulation* **125**, 931-944 (2012).

299    3.      Rehm, H.L. Disease-targeted sequencing: a cornerstone in the clinic. *Nature reviews. Genetics* **14**,
300           295-300 (2013).

301    4.      Clark, M.J. et al. Performance comparison of exome DNA sequencing technologies. *Nat Biotech*
302           **29**, 908-914 (2011).

303    5.      Mamanova, L. et al. Target-enrichment strategies for next-generation sequencing. *Nat Meth* **7**,
304           111-118 (2010).

305    6.      Braslavsky, I., Hebert, B., Kartalov, E. & Quake, S.R. Sequence information can be obtained from
306           single DNA molecules. *Proceedings of the National Academy of Sciences of the United States of*
307           *America* **100**, 3960-3964 (2003).

308    7.      Thompson, J.F. et al. Single-step capture and sequencing of natural DNA for detection of BRCA1
309           mutations. *Genome research* **22**, 340-345 (2012).

310    8.      Xie, X.S. & Trautman, J.K. Optical studies of single molecules at room temperature. *Annual*
311           *Review of Physical Chemistry* **49**, 441-480 (1998).

312    9.      Sakai, K. et al. Detection of epidermal growth factor receptor T790M mutation in plasma DNA
313           from patients refractory to epidermal growth factor receptor tyrosine kinase inhibitor. *Cancer*
314           *Science* **104**, 1198-1204 (2013).

315    10.    Paez, J.G. et al. EGFR Mutations in Lung Cancer: Correlation with Clinical Response to Gefitinib
316           Therapy. *Science* **304**, 1497-1500 (2004).

317    11.    Davies, H. et al. Mutations of the BRAF gene in human cancer. *Nature* **417**, 949-954 (2002).

318    12.    Aird, D. et al. Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries.
319           *Genome Biology* **12**, R18 (2011).

320    13.    Milos, P.M. Emergence of single-molecule sequencing and potential for molecular diagnostic
321           applications. *Expert review of molecular diagnostics* **9**, 659-666 (2009).

322    14.    Harris, T.D. et al. Single-molecule DNA sequencing of a viral genome. *Science* **320**, 106-109
323           (2008).

324    15.    Chen, F. et al. The History and Advances of Reversible Terminators Used in New Generations of
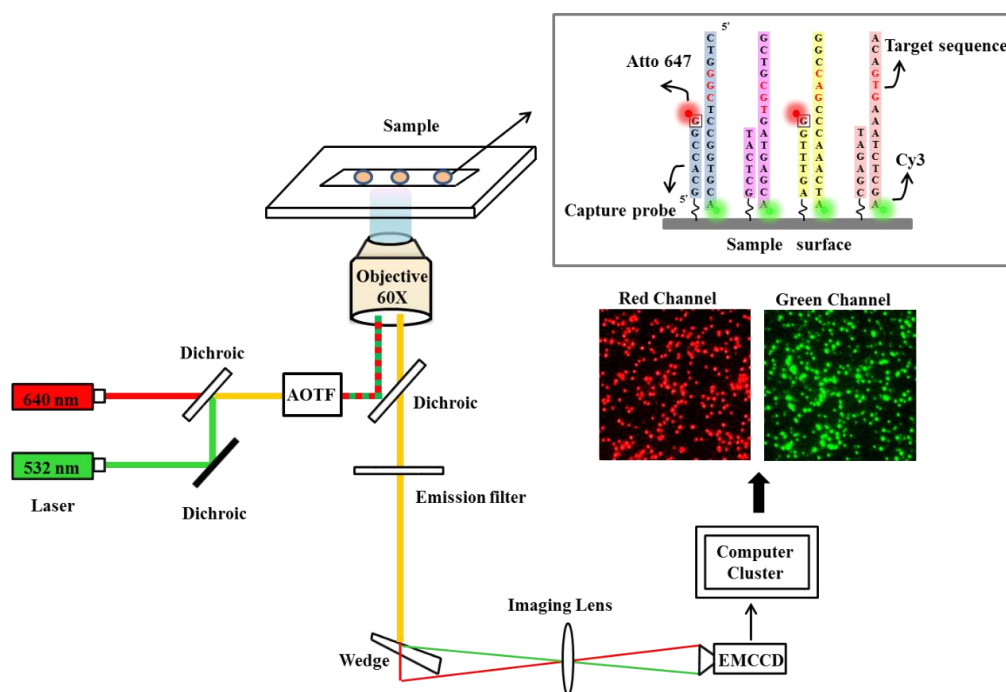325           Sequencing Technology. *Genomics, Proteomics & Bioinformatics* **11**, 34-40 (2013).
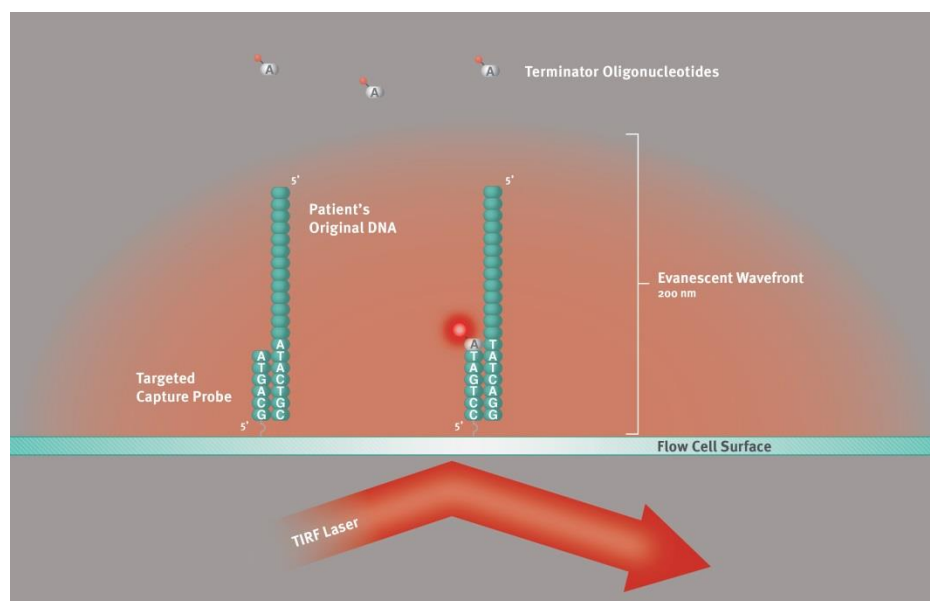
326

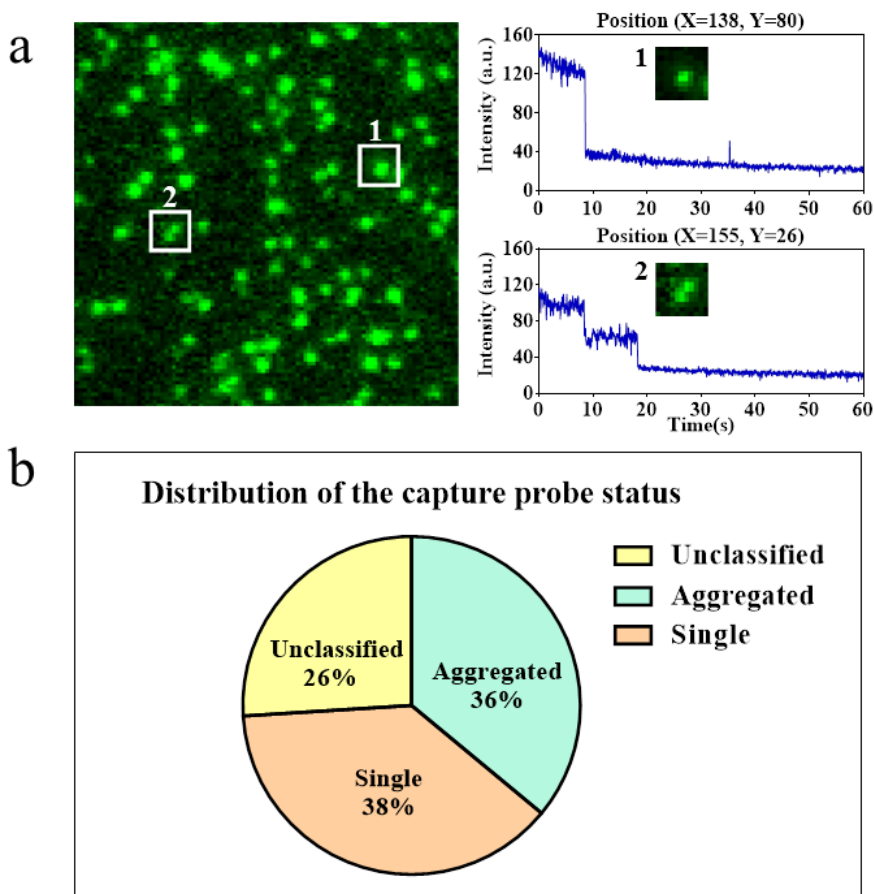# Acknowledgement

333    **Figures**

334    a



335

336    b



337

338    **Figure 1. Schematic drawing of single molecule sequencing platform.**  (a) Schematic drawing of the
339    optical setup. The green laser illuminates the Cy3 dyes which are attached to 3' end of the target DNA
340    template. The Cy3 dyes are non-cleavable. The red laser illuminates the cleavable Atto647N dyes which

341  are attached to reversible terminators. Both Cy3 and Atto647N fluorescence spectra are recorded
342  independently by an EMCCD. (b) Schematic of primed DNA templates attached to epoxy coated
343  coverslip surface. The capture probes are covalently attached to the coverslip surface, and the target DNA
344  templates are hybridized to the capture probes. The evanescent wave of TIRF illuminated the area within
345  200nm above the flow cell surfaces. The DNAs attached to the surfaces are within the range of
346  evanescent wave.



347

348  **Figure 2. Quantifying the ratio of single molecules of capture probes.** (a), the photobleaching of
349  single molecules are in a single step. Here, single spot was traced and its intensity was recorded. A single-
350  step photobleaching indicated that this spot was composed only one Cy3 molecule, i.e the spot #1. Spot
351  #2 was composed of two molecules binding together and therefore displaying two steps of
352  photobleaching. (b), The composition of single molecules, aggregated molecules and unclassified cases in
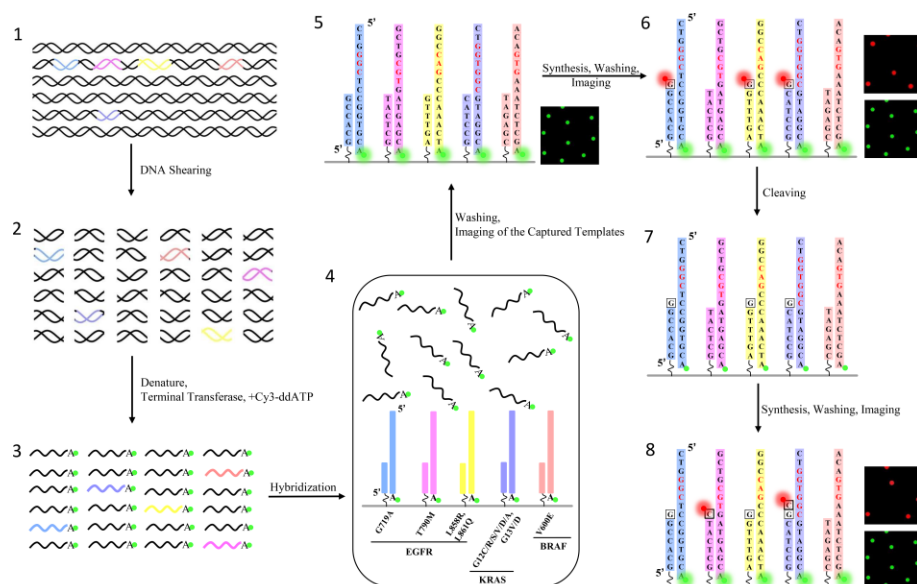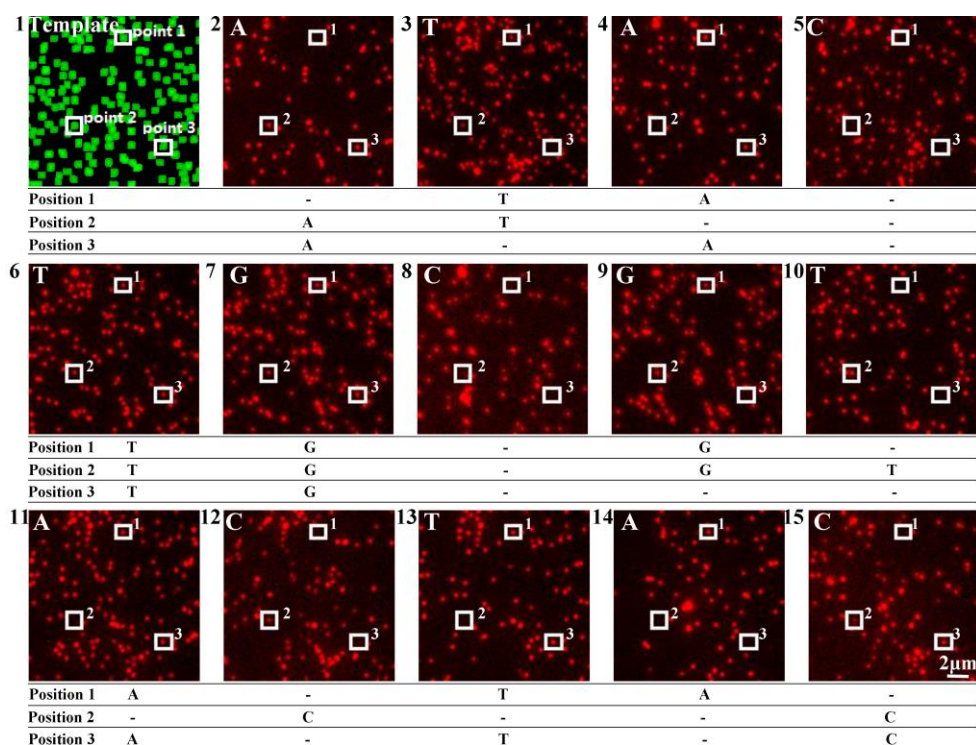353  one field of view.

354

355

356

357 **a**



358

359 **b**



| Base# | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|-------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|
| **Oligo#** | | | | | | | | | | | | | | | |
| 1 | T | A | T | G | G | A | T | A | C | C | C | T | C | A | C |
| 2 | A | T | T | G | G | T | C | C | T | G | T | T | T | A | A |
| 3 | A | A | T | G | A | T | C | A | C | G | G | T | A | A | A |

360

13

**Figure 3. (a)**, The sequencing procedure. DNA template with Cy3 attached at 3' end was hybridized to the flow cell anchored with capture probes (step 1-4). The capture probes are designed complementary to the genes of interest. Unhybridized DNA templates were washed away. The green laser excited the Cy3 fluorescence dye to locate the position of target DNA templates. (step 5). One of four types of reversible terminators labeled with red fluorescence and polymerases mixture were added to the flow cell. The DNA molecule extended a base if the reversible terminator matched complementary to the next base in the DNA molecule. Unincorperated reversible terminator was washed out. The red laser excited the Atto647 fluorescence dyes of reversible terminators (step 6). The fluorescence dyes in the reversible terminators were cleaved and wash away (step 7). A new cycle of sequencing began (step 8). (b), Multiple sequencing cycles, imaging and base calling. We traced a part of one field of view in multiple sequencing cycles. In the beginning, the image of Cy3 green fluorescence dyes were used to locate the position of target templates. Three positions were circled out and were traced. In the first cycle, reversible terminators A (nucleotide analogs) were flowed in for reaction. Position 2 and 3 successfully incorporated a base. In the second cycle, the reversible terminators T were flowed in for reaction. Position 1 and 2 successfully incorporated a base. The sequencing continued and the sequence of DNA template extended. The sequence of each DNA template in position 1, 2 and 3 can be reconstructed.

| Oligo Name | Probe Sequence (5'-3') | Gene |
|---|---|---|
| EKB-1P | TTTTTTTTTTCAGAGGCCTGTGCCAGGGACCTTACCTTATACACCGTGCCGAACGCACCG | EGFR |
| EKB-2P | TTTTTTTTTTTAGCAAAGCAGAAACTCACATCGAGGATTTCCTTGTTGGCTTTCGGAGATG | EGFR |
| EKB-3P | TTTTTTTTTTTCTGGATCCCAGAAGGTGAGAAAGTTAAAATTCCCGTCGCTATCAAGGAAT | EGFR |
| EKB-4P | TTTTTTTTTTTCACGTGTGCCGCCTGCTGGGCATCTGCCTCACCTCCACCGTGCAGCTCAT | EGFR |
| EKB-5P | TTTTTTTTTTTCAGGAACGTACTGGTGAAAACACCGCAGCATGTCAAGATCACAGATTTTG | EGFR |
| EKB-6P | TTTTTTTTTTTCTCCTTACTTTGCCTCCTTCTGCATGGTATTCTTTCTCTTCCGCACCCAG | EGFR |
| EKB-7P | TTTTTTTTTTTCCACAAAATGATTCTGAATTAGCTGTATCGTCAAGGCACTCTTGCCTAC | KRAS |
| EKB-8P | TTTTTTTTTTTAAATGGATCCAGACAACTGTTCAAACTGATGGGACCCACTCCATCGAGAT | BRAF |

| Oligo Name | Wild Sequence (5'-3') | Gene |
|---|---|---|
| EKB-1T-n | AATTCAAAAAGATCAAAGTGCTGGGCTCCGGTGCGTTCGGCACGGTGTATAAGGTAAGGTCCCTGGCACAGGCCTCTG | EGFR |
| EKB-2T-n | GTCGCTATCAAGGAATTAAGAGAAGCAACATCTCCGAAAGCCAACAAGGAAATCCTCGATGTGAGTTTCTGCTTTGCT | EGFR |
| EKB-3T-n | TTGGCTTTCGGAGATGTTGCTTCTCTTAATTCCTTGATAGCGACGGGAATTTTAACTTTCTCACCTTCTGGGATCCAG | EGFR |
| EKB-4T-n | GAGGCAGCCGAAGGGCATGAGCTGCGTGATGAGCTGCACGGTGGAGGTGAGGCAGATGCCCAGCAGGCGGCACACGTG | EGFR |
| EKB-5T-n | TCTTCCGCACCCAGCAGTTTGGCCAGCCCAAAATCTGTGATCTTGACATGCTGCGGTGTTTTCACCAGTACGTTCCTG | EGFR |
| EKB-6T-n | GATCACAGATTTTGGGCTGGCCAAACTGCTGGGTGCGGAAGAGAAAGAATACCATGCAGAAGGAGGCAAAGTAAGGAG | EGFR |
| EKB-7T-n | TAAACTTGTGGTAGTTGGAGCTGGTGGCGTAGGCAAGAGTGCCTTGACGATACAGCTAATTCAGAATCATTTTGTGGA | KRAS |
| EKB-8T-n | TAGGTGATTTTGGTCTAGCTACAGTGAAATCTCGATGGAGTGGGTCCCATCAGTTTGAACAGTTGTCTGGATCCATTT | BRAF |

| Oligo Name | Mutation Sequence (5'-3') | Gene | Amino Acid Variant |
|---|---|---|---|
| EKB-1T-m | AATTCAAAAAGATCAAAGTGCTGGCCTCCGGTGCGTTCGGCACGGTGTATAAGGTAAGGTCCCTGGCACAGGCCTCTG | EGFR | G719A |
| EKB-2T-m | AAAGTTAAAATTCCCGTCGCTATCAAGACATCTCCGAAAGCCAACAAGGAAATCCTCGATGTGAGTTTCTGCTTTGCT | EGFR | △E746-A750del |
| EKB-3T-m | ACATCGAGGATTTCCTTGTTGGCTTTCAATTCCTTGATAGCGACGGGAATTTTAACTTTCTCACCTTCTGGGATCCAG | EGFR | △E747-A753del |
| EKB-4T-m | GAGGCAGCCGAAGGGCATGAGCTGCATGATGAGCTGCACGGTGGAGGTGAGGCAGATGCCCAGCAGGCGGCACACGTG | EGFR | T790M |
| EKB-5T-m | TCTTCCGCACCCAGCTGTTTGGCCCGCCCAAAATCTGTGATCTTGACATGCTGCGGTGTTTTCACCAGTACGTTCCTG | EGFR | L858R |
| EKB-6T-m | GATCACAGATTTTGGGCGGGCCAAACAGCTGGGTGCGGAAGAGAAAGAATACCATGCAGAAGGAGGCAAAGTAAGGAG | EGFR | L861Q |
| EKB-7T-m | TAAACTTGTGGTAGTTGGAGCTTCTGACGTAGGCAAGAGTGCCTTGACGATACAGCTAATTCAGAATCATTTTGTGGA | KRAS | G12S, G13D |
| EKB-8T-m | TAGGTGATTTTGGTCTAGCTACAGAGAAATCTCGATGGAGTGGGTCCCATCAGTTTGAACAGTTGTCTGGATCCATTT | BRAF | V600E |

**Table 1. The capture probe and target DNA sequence information.** The top block is the capture probe sequences we synthesized. The capture probes were designed to capture EGFR/KRAS/BRAF genes. The middle block is the target DNA sequence designed for testing. These sequences are design based on the wild type of EGFR/KRAS/BRAF genes. Nucleotide bases in red color are drug related mutation sites. The bottom block is the target DNA sequence designed based on the mutant type.

383
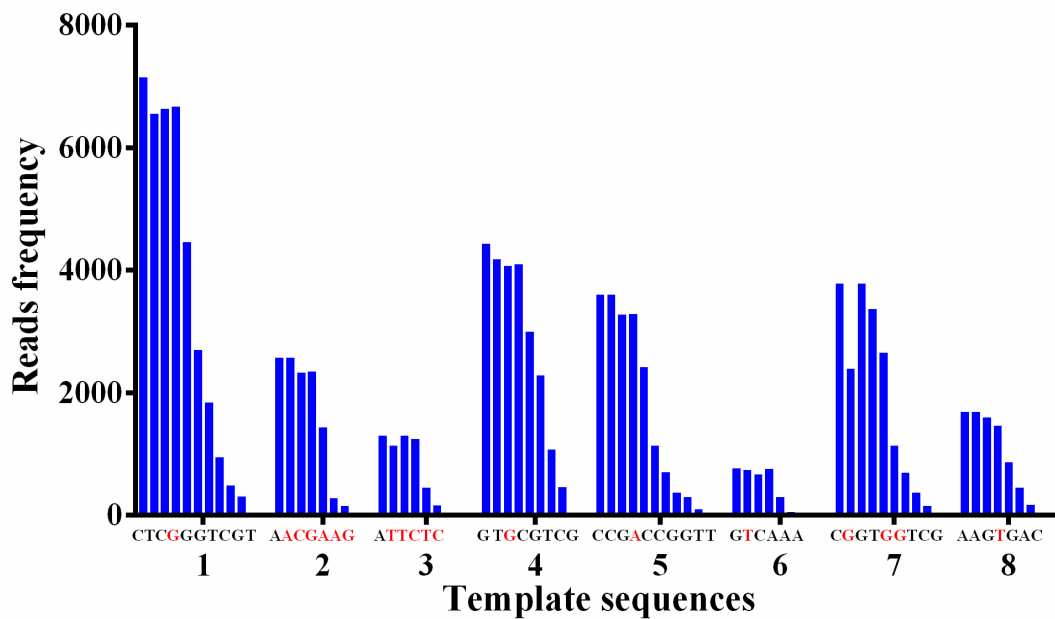
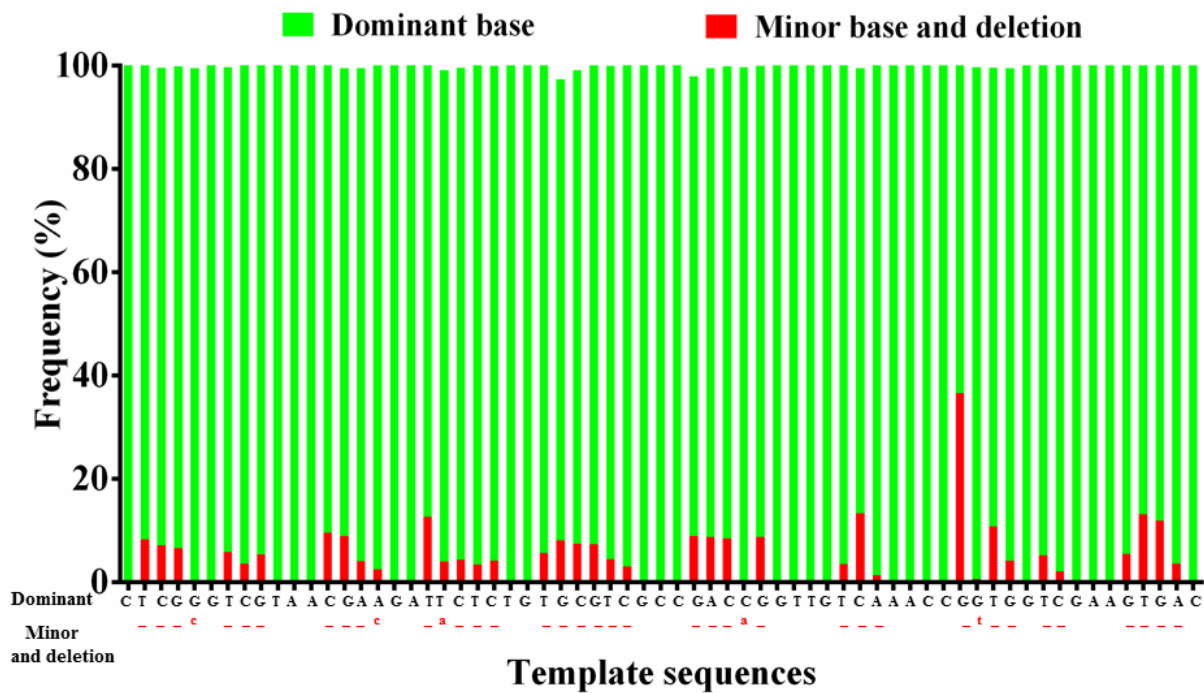| Reference | EGFR | | | | | | KRAS | BRAF |
|---|---|---|---|---|---|---|---|---|
| | CTCGGGTCGT | AACGAAG | ATTCTCT | GTGCGTCG | CCGACCGGTT | GTCAAAC | CGGTGGTCG | AAGTGAC |
| Alignment | CTCGGGTCGT | AACGAAG | ATTCTC | GT_CGTCG | CCGACCGGTT | GTCAA | C_GTGGT | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | ATTCTCT | GTaCGTCG | CCGACCGGTT | GTCAAAC | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | ATTC | GT_CGTCG | CCGACCGGTT | GT_AA | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCG | AACGAAG | ATaCTC | GT_CGTCG | CCGACCGGTT | GTCAA | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | ATaCTC | GTaCGTCG | CCGACCGGT | GTCA | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | A_TCTC | GTaCGTCG | CCGACCGGT | GTCA | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCGT | AACGAcG | A_TCTCT | GT_CGTCG | CCGACCGGTT | GTCA | CGGTGGTCG | AAGTGA |
| | CTCGGGTCGT | AACG_AG | A_TCTCT | GT_CGTCG | CCGACCGGTT | GTCAA | CGGTGGTCG | AAGTGA |
| | CTCGGGTCGT | AACG_AG | ATTC | GTGCGTCG | CCGACCGGTT | GTCA | CGGTGGTC | AAGTGA |
| | CTCGGGTCGT | AA_GAAG | ATTC | GTGCGTCG | CCGACCGGTT | G_CAAAC | CGGT | AAGTGA |
| | CTCGGGTCGT | AACGAA | ATTC | GTGCGTCG | CCGACCGGT | GTCAAAC | CGGTGGTC | AAGTGA |
| | CTCGGGTCGT | AACGAA | ATTC | GTGCGT | CCGACCGGT | GTCAAAC | C_GTGGTCG | AAGTGA |
| | CTCGGGTCGT | AACGAA | ATTCTC | GT_CGTCG | CCGACCGGT | GTCA | CGGTGGTC | AAGTGAC |
| | CTCGGGTCGT | AACGA | ATTCT | GTGCGT | CCGACCGGT | GT_AA | CGGTGGTC | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | ATTCT | GTaCGTCG | CCGACCGGT | GT_AA | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | ATaCT | GT_CGTCG | CCGACCGGT | GTCA | CGGTGGTCG | AAGTGAC |
| | CTCGGGTCGT | AACGAAG | ATaCTC | GT_CGTCG | CC_ACCGGTT | GTCA | CGGTGGTCG | AAGTGA |
| | CTCGGGT | AACGAAG | ATTCT | GT_CGTCG | CCGACCGGT | GTCA | CGGTGGTCG | AAGTGA |
| | CTCGGGTC | AACG_AG | A_TCTCT | GTGCGTC | CCGACCGGT | GTCA | CGGTGGTC | AAGTGA |
| Consensus sequence | CTCGGGTCGT | AACGAAG | ATTCTCT | GTGCGTCG | CCGACCGGTT | GTCAAAC | CGGTGGTCG | AAGTGAC |

384

385 **Table 2. Sequencing alignment of raw reads.** The top row is the reference sequence. Insertion errors
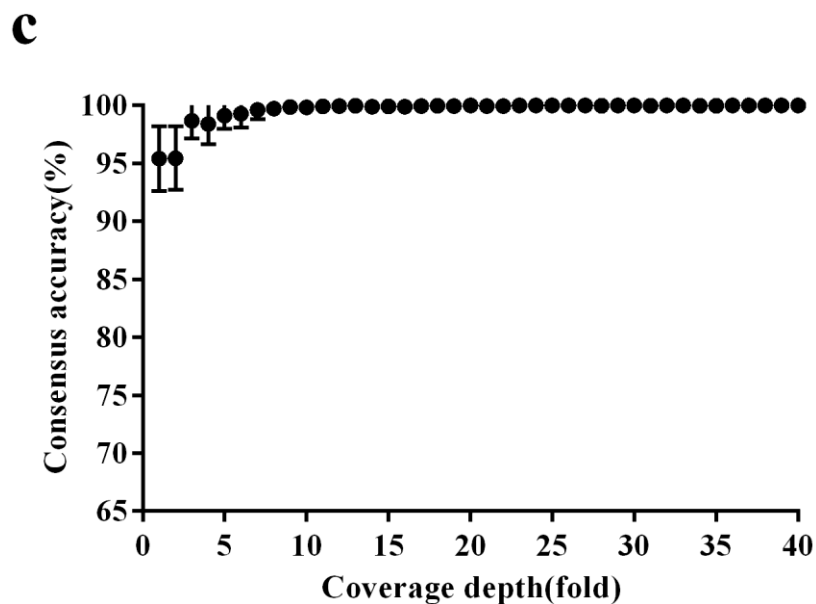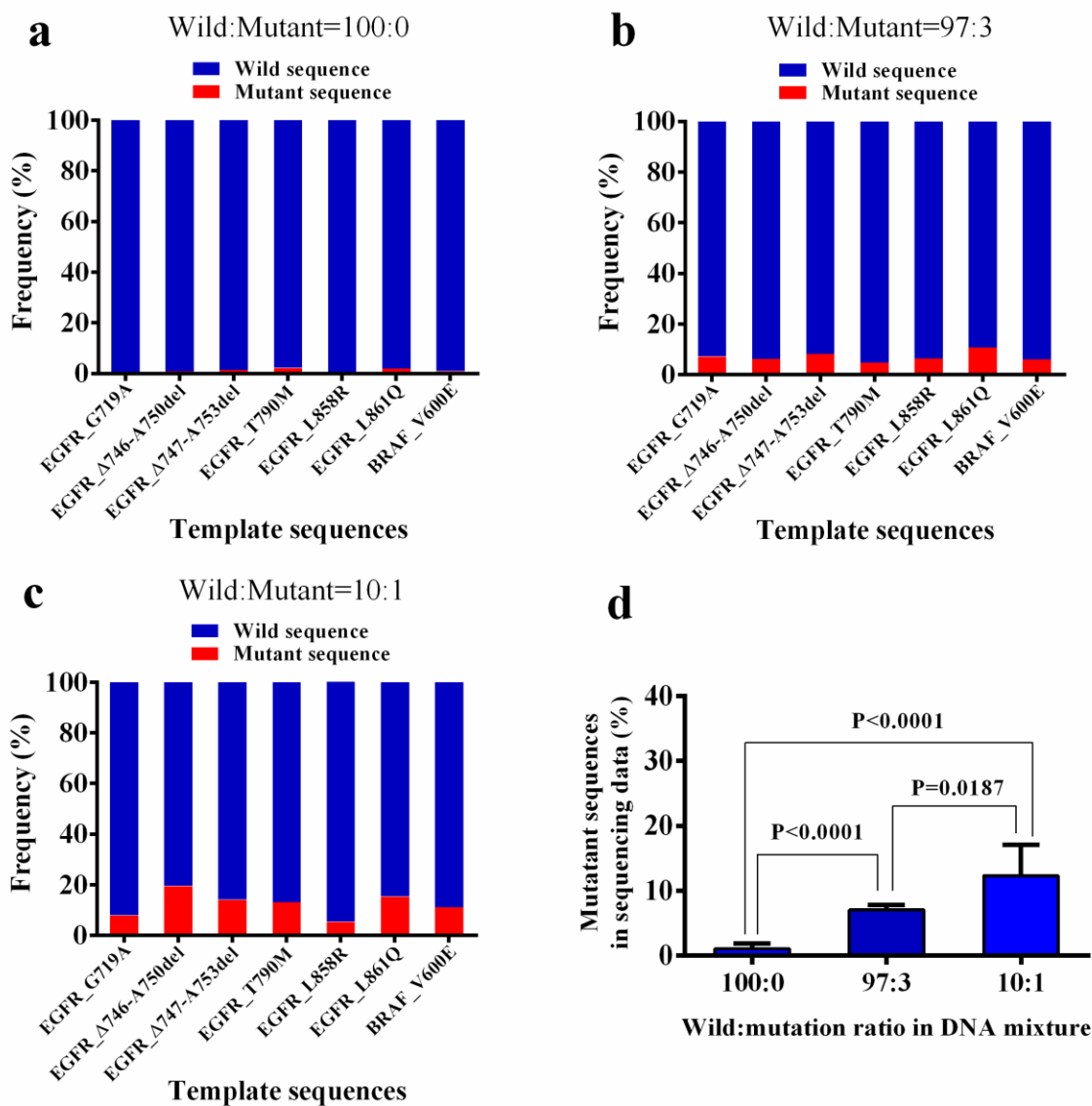386 were not shown in the alignment.

15

**Figure 4.** (a), the coverage per base. The frequently mutant position is in red color. Y-axis is the number reads mapped to each position. (b), The dominant and minor base at the each position. (c), The consensus accuracy increased with coverage depth. Sampling-subsampling was performed to simulate low coverage situation.

394

**Figure 5. Detecting mutant sequences in a mixture**. The wild type and mutant DNA were mixed at 100:0, 10:1 and 97:3 ratios. The mixed DNA was subjected to sequencing. Each sequence reads was aligned to the wild type and mutant type reference sequences and alignment scores were calculated. If the alignment score of wild type reference sequence was higher than that of mutant type reference sequence, the original sequence read was classified as wild type. Otherwise, it was classified as mutant type. The frequency of wild type and mutant type sequence reads are calculated for each reference. (a-c) The frequency of wild type and mutant type sequences calculated from the sequencing data. (d), The average of mutant sequences in sequencing data over all template sequences. P value is calculated by two-tailed Student T test.

404