

1 mInDel: an efficient pipeline for high-throughput InDel marker discovery

2 Yuanda Lv¹, Yuhe Liu², Xiaolin Zhang¹ and Han Zhao^{1,*}

3 ¹ Institute of Agricultural Biotechnology, Jiangsu Academy of Agricultural Sciences,
4 Nanjing 210014, China;

5 ² Department Of Crop Sciences, University of Illinois, Urbana-Champaign, 1201 West
6 Gregory Drive, Urbana, IL 61801, USA.

7 *Corresponding author: Han Zhao, zhaohan@jaas.ac.cn

8 **Abstract**

9 **Background**

10 Next-Generation Sequencing (NGS) technologies have emerged as a powerful tool to
11 reveal nucleotide polymorphisms in a high-throughput and cost-effective manner.
12 However, it remains a daunting task to proficiently analyze the enormous volume of data
13 generated from NGS and to identify length polymorphisms for molecular marker discovery.
14 The development of insertion-deletion polymorphism (InDel) markers is in particular
15 computationally intensive, calling for integrated high performance methods to identify
16 InDels with high sensitivity and specificity, which would directly benefit areas from
17 genomic studies to molecular breeding.

18 **Results**

19 We present here a NGS-based tool for InDel marker discovery (mInDel), a
20 high-performance computing pipeline for the development of InDel markers between any
21 two genotypes. The mInDel pipeline proficiently develops InDel markers by comparing
22 shared region size using sliding alignments between assembled contigs or reference
23 genomes. mInDel has successfully designed thousands of InDel markers from maize
24 NGS data locally and genome-wide. The program needs less than 2 hours to run when
25 using 20 threads on a high-performance computing server to implement 40G data.

26 **Conclusions**

27 mInDel is an efficient, integrated pipeline for a high-throughput design of InDel markers
28 between genotypes. It will be particularly applicable to the crop species which require a
29 sufficient amount of DNA markers for molecular breeding selection. mInDel is freely
30 available for downloading at www.github.com/lyd0527/mInDel website.

31 **Background**

32 For most plant species, genetic mapping of a trait of interest relies largely upon the
33 trait-marker association revealed from segregating populations [1-3]. The quantity and

34 quality of molecular markers are therefore instrumental to the resolution and accuracy of
35 successful linkage construction, gene mapping, and cloning studies [4-7]. The
36 development of a proficient and cost-effective marker system is also critical in other
37 studies including population genetics, molecular evolution and marker assisted breeding,
38 etc [8-10].

39 In the past decades, the Southern blot-based marker systems such as restriction fragment
40 length polymorphism (RFLP) have been progressively replaced by PCR-based markers
41 such as random amplified polymorphic DNA (RAPD), simple sequence repeat (SSR), and
42 amplified fragment length polymorphism (AFLP). However, the development and
43 large-scale screening of polymorphic molecular makers at the time usually involved
44 extensive library construction and Sanger sequencing which are time-consuming and
45 labor-intensive. The recent advent of next-generation sequencing (NGS) technologies has
46 revolutionized the pace and throughput of sequence information generation, and greatly
47 accelerated the discovery processes for genetic markers [11,12]. A large amount of single
48 nucleotide polymorphisms (SNP) were developed based on NGS data and showed a
49 strong vitality in several studies [11,13].

50 Compared to other PCR-based marker types and the more recent SNP system,
51 Insertion-Deletion markers (InDels) have shown several advantages. Firstly, InDels are
52 much more abundant than other PCR-based markers such as SSRs in genomes [15-18].
53 Secondly, InDels are potentially multi-allelic and codominant, offering more genomic
54 information than the bi-allelic SNPs [18-20]. Lastly, existing InDels can be readily used to
55 genotype new materials via simple gel electrophoresis, whereas SNP detection requires
56 specialized and usually more expensive equipment (e.g. DHPLC) or assays (e.g.
57 TaqMan). All these characteristics make InDel markers highly suitable for applications
58 from basic genomic research to applied plant molecular breeding.

59 Faced with the great wealth of sequence information generated from NGS however,
60 existing programs or pipelines for Indel marker discovery are struggling to keep in step.
61 They are usually reliant on a reference genome and can only discover small indels (≤ 10
62 nt) due to the limited alignment algorithm; the primer design efficiency is often
63 compromised by the large nucleotide variation flanking the InDels, all of which inhibited
64 their usage in plant molecular breeding. To overcome the aforementioned problems, here
65 we present and validate an efficient and high-throughput pipeline for InDel markers
66 development based on NGS data independent of a reference genome. The mInDel
67 pipeline has been primarily designed to work with diploid plant genomes, but can also be
68 applied to polyploid plant species such as cotton, wheat and potato.

69 **Implementation**

70 Overview of the pipeline

71 mInDel is a comprehensive and high performance marker development program to
72 identify specific InDel primers between genotypes in a high-throughput manner. mInDel
73 mainly implements five modules: the Pre-processing module, the *De novo* assembly
74 module, the overlap PCR primer design in batch module ,the ePCR mapping module and
75 InDel screening and marker development module. In the following, we describe the
76 mInDel pipeline for InDel marker development, including one pre-processing step, three
77 main marker discovery steps and a post-processing step, as shown in Figure 1.

78 Pre-processing

79 This module is intended for quality control of sequence data, filtering against low quality,
80 short reads and low complex regions. First, the Btrim [21] program is employed to remove
81 low quality nucleotides at the 3'-end of a read by Q20 (Phred score of 20). Then, a custom
82 Perl script removes sequence reads shorter than 40 bp (L40). When a sequence from a
83 sequence-pair was removed, the remaining one is put into a separate file and used as a
84 singleton during *De novo* assembly.

85 *De novo* assembly

86 This module is designed for a high quality sequence assembly. We provide two different
87 assembly strategies for two commonly used platforms (Illumina, Sanger and 454). For the
88 Illumina platform with relatively shorter reads, we adopt the de-Bruijn graphs algorithm
89 [22,23], which was specifically designed for processing short sequences and has been
90 shown to be effective for sequence assembly. For Sanger or 454 platform with relatively
91 long reads, the overlap algorithm [24,25] is employed for its efficacy at producing high
92 quality assemblies of genome data. Here, SOAPdenovo2 [26] (Illumina Platform) or
93 Newbler (sanger or 454 Platform) (<http://www.454.com/products/analysis-software/>) are
94 used for sequence assembly using a coverage cutoff of 5 and discarding contigs shorter
95 than 100 bp. By providing two different assembly strategies, mInDel enables users to
96 choose flexibly according to their specific needs.

97 Overlap PCR Probe Design

98 This module is designed to generate high quality primer pairs. It generates large tracks of
99 primer pairs, with estimated PCR products in overlap spanning the entire input region.
100 This module firstly generates overlap fragments from the given sequence using a sliding
101 window method. The overlap fragments are then passed on to the PCR primer selecting
102 program Primer3 [27]. Primer3 then detects primer sets based on user's setting criteria.
103 The output file is parsed and generates a list of probe primers consisting of the primer
104 sequence, the calculated melting temperature, the quality score of the primer set, primer
105 positions, primer lengths, the PCR product length, and the amount of overlap between
106 fragments.

107 ePCR mapping and polymorphic screening

108 This module is designed for ePCR mapping. In silico PCR strategy is used to predict
109 possible InDel differences from orthologous and homologous loci between any two
110 genotypes. Primer sets generated from the previous module are run through in silico PCR
111 analysis, with a threshold less than 3 mismatch bases using Bowtie aligner [28]. PCR
112 products from a given region with different amplicon sizes are used to mine potential
113 InDels among orthologous and homologous loci between the two samples. As a result,
114 InDels larger than certain base pairs are preferentially mapped to the reference genome
115 for single amplified locus screening. Finally, the site-specific InDel primer pairs are
116 chosen.

117 InDel marker development

118 In the final post-processing step, mInDel pipeline performs the post-processing to
119 integrate core results and generate tab-delimited or Excel-compatible files. The final
120 output compiles information for candidate InDel markers including forward and reverse
121 primers, the product size (sample A and sample B), the InDel size (delta) and the
122 chromosome location.

123 Test data

124 A small set of test data is provided with mInDel. It consists of simulated Illumina 100bp *2
125 paired-end reads, and raw reference sequence. These data enable testing of most
126 commonly used features. Users are strongly encouraged to run this test with the
127 command '*sh test.sh*' after installation which usually takes about 5 minutes to run on a
128 8-core 2.5 GHz computer. The test data also serves as the foundation for the example
129 analysis described in the TUTORIAL file. We routinely run this test in the course of adding
130 and refactoring code.

131 Results and discussion

132 InDel markers development between maize inbred lines

133 In addition to the small test data sets included with mInDel, we demonstrate here an
134 example analysis of larger data sets from four maize inbred lines (B73 as the reference
135 genome, Mo17, Qi319 and Zheng58). Mo17vsB73 and Qi319vsZheng58 were two
136 experiments used for testing the mInDel pipeline performance under two different
137 conditions: with and without the reference genome.

138 The Mo17, Qi319 and Zheng58 genome data was acquired from the NCBI's Sequence
139 Read Archive (SRA) database while the latest B73 reference genome data (V3) was from
140 plant ensembl databases (http://plants.ensembl.org/Zea_mays). A total of 24.8 Gb 454
141 sequencing data from Mo17, as well as 30.4Gb and 12.5 Gb Illumina paired-end 100 bp
142 reads for Qi319 and Zheng58 genotypes were retrieved from NCBI's SRA database
143 respectively.

144 The raw FASTQ-formatted reads were then pre-processed by Q20 and L40 filtering using
145 mInDel's quality control module. As a result, 21.2 Gb cleaned 454 reads from Mo17, along
146 with 28.6 Gb and 10.1 Gb high quality Illumina reads from Qi319 and Zheng58 were
147 obtained. The clean reads were then separately assigned into mInDel's *de novo* assembly
148 process. Contig assembly was performed with de bruijn graph method by SOAPdenovo
149 for illumina sequences and with overlap method by Newbler for 454 sequences to ensure
150 assembly quality. 1,351,472 contigs for Mo17, 476,470 contigs for Qi319 and 395,075
151 contigs for Zheng58 (≥ 100 bp) were generated and accounted for approximately 75%,
152 19.5% and 8.8% of the B73 reference genome. The contig N50 of the assemblies was
153 2,012 bp, 1,157 bp and 457 bp, with the longest contig being 33,817 bp, 18,068 bp and
154 7,286 bp for Mo17, Qi319 and Zheng58 respectively.

155 Overlap PCR Primer Design

156 Furthermore, mInDel uses the open source Primer3 software to identify primers and
157 probes with desired thermodynamic properties from candidate sequences in batch mode.
158 At this stage of the pipeline, each contig larger than 300 bp was cut into several segments
159 by sliding window of 300 bp with a step 150 bp, which ensured the overlapping amplified
160 regions to cover the whole sequence. The amount of overlap is calculated from the 5' end
161 of a fragment to the 3' end of the previous (adjacent) fragment and takes into account of
162 the primer lengths. Generally, we have found that when sequencing these PCR products
163 with dye-primer sequencing an overlap 150-200 bp is best. As a result, a total of
164 4,205,672 and 610,091 primer probes from Mo17 and Zheng58 were generated using
165 mInDel.

166 ePCR and polymorphic screening

167 In this study, two groups, Mo17vsB73 and Qi319vsZheng58, were assigned as two
168 distinct data sets for all possible InDels loci based on in silico PCR analysis. Primer
169 probes from Mo17 and Zheng58 contigs were run against the B73 reference genome
170 sequences and Qi319 assembled contigs using in silico PCR analysis with a threshold
171 less than 3 mismatch bases. In total, 1,855,222 primer probes from Mo17 were
172 successfully mapped to B73 genome, while 307,345 from Zheng58 were mapped to
173 Qi319 assembled contigs. Of mapped primers, 756,723 and 226,680 primer pairs were
174 mapped to unique locations in B73 and Qi319 genomes and were considered as
175 site-specific or single-copy primers, whereas primers mapped to multiple positions were
176 discarded. Furthermore, a total of 739,250 and 223,084 InDel loci from group Mo17vsB73
177 and group Qi319vsZheng58 were identified using mInDel.
178 For InDel sizes, 133,918 and 44,389 markers have an InDel of 1-20 bp, while 594,717 and
179 178,695 markers have an InDel greater than or equal to 20 bp. The average length of the
180 InDels was 71 bp and 66 bp, with the maximum InDels up to 387 bp and 371 bp. Finally,
181 single-copy primer pairs with InDel size larger than 20 bp were selected and considered
182 as candidate InDel markers for the ease of PCR detection and electrophoresis screening.

183 Experimental validation of mInDel

184 To verify the accuracy and efficiency of mInDel, primer pairs randomly chosen across
185 genome from each of the two comparison groups were synthesized for polymorphism
186 detection between maize inbred lines Mo17 and B73 and between Qi319 and Zheng58.
187 Of 318 sampling primers from Mo17vsB73, 276 (86.8%) showed polymorphisms with the
188 expected amplicon sizes. In group Qi319vsZheng58, 437 of 602 InDel primer pair (72.6%)
189 showed polymorphisms.

190 Application on QTL fine mapping

191 Numbers of molecular markers are crucial to a high resolution genetic mapping. Marker
192 saturation is one way to achieving fine mapping on recombination-rich regions with a yet
193 saturated genetic map. Here, we successfully narrowed down one grain starch's QTL
194 region from a maize segregating population using InDel markers developed from mInDel
195 pipeline. 90 InDel markers were developed for both regions by mInDel. As a result, 71
196 InDels showed polymorphisms between two parents, and 51 were successfully genetically
197 mapped on regions with the synteny conservation to physical locations.

198 Conclusion

199 mInDel is an efficient and robust high-throughput pipeline in developing InDel markers
200 between any two genomes. While most parameters work well with default settings, mInDel
201 also offers great customization on parameter tuning. All modules are able to run
202 independently, with the input/output features offering ample options and great flexibility.
203 As a freely available, high-throughput and fully integrated software for InDel marker
204 development, mInDel provides a useful resource for various genetics and genomics
205 research and marker assisted breeding.

206 Availability and requirements

207 **Project name:** mInDel

208 **Project home page:** www.github.com/lyd0527/mInDel

209 **Operating system(s):** Linux OS. Tests were performed in Centos and Ubuntu Linux
210 systems. Some minor modifications are needed for other operation systems.

211 **Programming language:** Perl 5.10 or above, Shell.

212 **Other requirements:** Primer 3.2.0 or above and Bowtie any version.

213

214 Abbreviation

215 InDel: insertion and deletion
216 NGS: Next Generation Sequencing

217 **Competing interests**

218 The authors declare that they have no competing interests.

219 **Authors' contributions**

220 HZ and YDL conceived the study and wrote the manuscript. YDL developed the software
221 for the analysis. XLZ and YHL compared the mInDel signatures with the experimentally
222 verified signatures. All authors read and approved the final manuscript.

223 **Acknowledgements**

224 This work was supported by the Natural Science Foundation of China(NSFC) under Grant
225 No. 31271728, Jiangsu Agriculture Science and Technology Innovation Fund (JASTIF)
226 under Grant No. CX(13)5055.

227 **References**

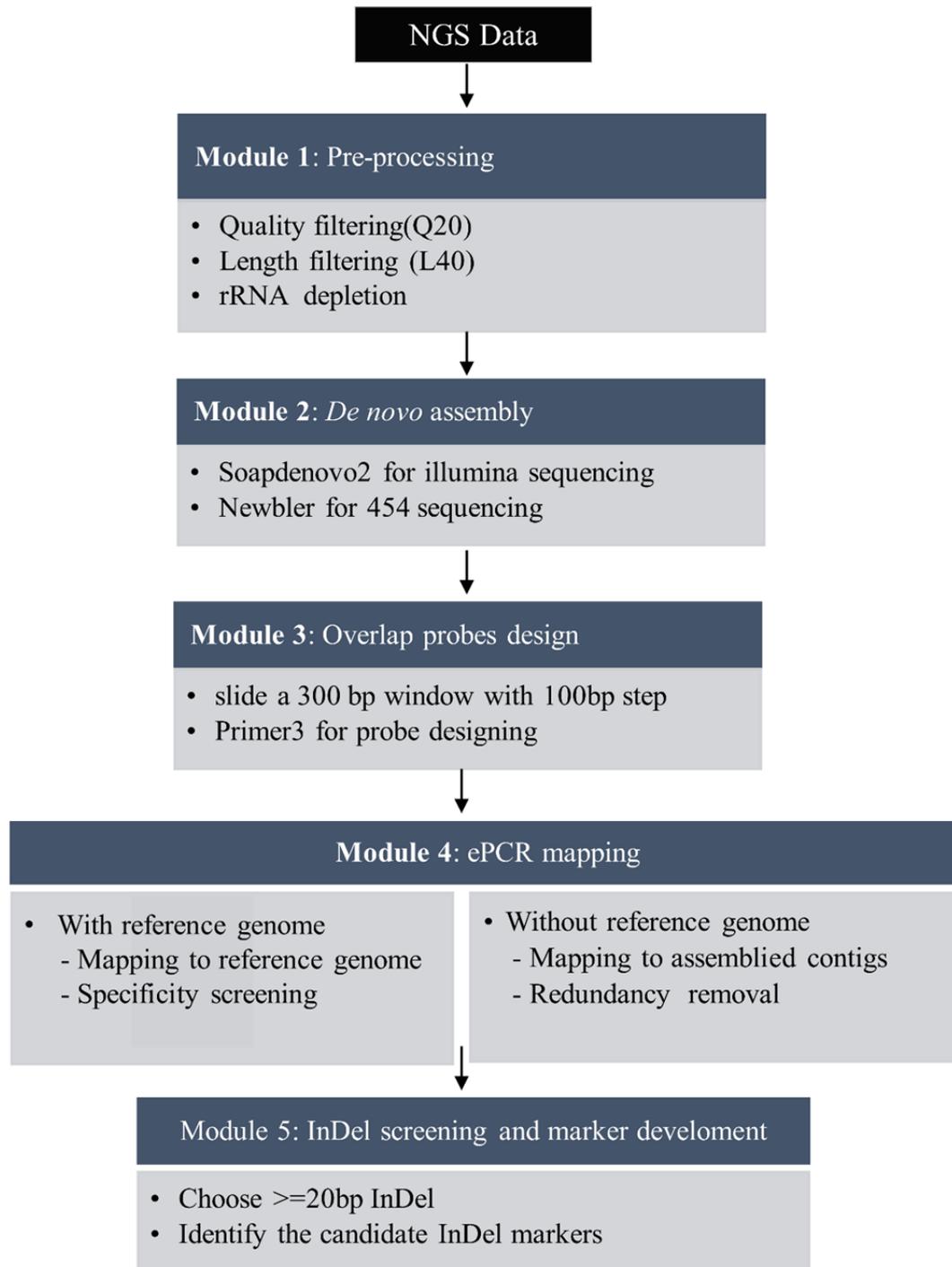
- 228 1. Doerge RW: **Mapping and analysis of quantitative trait loci in experimental populations.** *Nat*
229 *Rev Genet* 2002, **3**(1):43-52.
- 230 2. Pflieger SEP, Lefebvre VER, Causse M: **The candidate gene approach in plant genetics: a**
231 **review.** *Mol Breeding* 2001, **7**(4):275-291.
- 232 3. Tanksley SD: **Mapping polygenes.** *Annu Rev Genet* 1993, **27**:205-233.
- 233 4. Smith J, Chin E, Shu H, Smith OS, Wall SJ, Senior ML, Mitchell SE, Kresovich S, Ziegler J: **An**
234 **evaluation of the utility of SSR loci as molecular markers in maize (*Zea mays* L.):**
235 **comparisons with data from RFLPs and pedigree.** *Theor Appl Genet* 1997, **95**(1-2):163-173.
- 236 5. Remington DL, Ungerer MC, Purugganan MD: **Map-based cloning of quantitative trait loci:**
237 **progress and prospects.** *Genet Res* 2001, **78**(3):213-218.
- 238 6. Young ND: **Potential applications of map-based cloning to plant pathology.** *Physiol Mol Plant*
239 *P* 1990, **37**(2):81-94.
- 240 7. Staub JE, Serquen FC, Gupta M: **Genetic markers, map construction, and their application in**
241 **plant breeding.** *Hortscience* 1996, **31**(5):729-741.
- 242 8. Sunnucks P: **Efficient genetic markers for population biology.** *Trends in Ecology & Evolution*
243 2000, **15**(5):199-203.
- 244 9. Cavalli-Sforza LL, Feldman MW: **The application of molecular genetic approaches to the**
245 **study of human evolution.** *Nat Genet* 2003, **33**:266-275.
- 246 10. Eathington SR, Crosbie TM, Edwards MD, Reiter RS, Bull JK: **Molecular markers in a**
247 **commercial breeding program.** *Crop Sci* 2007, **47**(3):154.
- 248 11. Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML: **Genome-wide genetic**

- 249 **marker discovery and genotyping using next-generation sequencing.** *Nat Rev Genet* 2011,
250 **12(7):499-510.**
- 251 12. Varshney RK, Nayak SN, May GD, Jackson SA: **Next-generation sequencing technologies and**
252 **their implications for crop genetics and breeding.** *Trends Biotechnol* 2009, **27(9):522-530.**
- 253 13. Ramos AM, Crooijmans RP, Affara NA, Amaral AJ, Archibald AL, Beever JE, Bendixen C,
254 Churcher C, Clark R, Dehais P *et al*: **Design of a high density SNP genotyping assay in the pig**
255 **using SNPs identified and characterized by next generation sequencing technology.** *PloS one*
256 2009, **4(8):e6524.**
- 257 14. Atienzar FA, Jha AN: **The random amplified polymorphic DNA (RAPD) assay and related**
258 **techniques applied to genotoxicity and carcinogenesis studies: a critical review.** *Mutation*
259 *Research/Reviews in Mutation Research* 2006, **613(2):76-102.**
- 260 15. O'Hanlon PC, Peakall R, Briese DT, Others: **A review of new PCR-based genetic markers and**
261 **their utility to weed ecology.** *WEED RESEARCH-OXFORD-* 2000, **40(3):239-254.**
- 262 16. Powell W, Morgante M, Andre C, Hanafey M, Vogel J, Tingey S, Rafalski A: **The comparison**
263 **of RFLP, RAPD, AFLP and SSR (microsatellite) markers for germplasm analysis.** *Mol*
264 *Breeding* 1996, **2(3):225-238.**
- 265 17. Blears MJ, De Grandis SA, Lee H, Trevors JT: **Amplified fragment length polymorphism**
266 **(AFLP): a review of the procedure and its applications.** *Journal of Industrial Microbiology*
267 *and Biotechnology* 1998, **21(3):99-114.**
- 268 18. Vignal A, Milan D, SanCristobal M, Eggen AE, Others: **A review on SNP and other types of**
269 **molecular markers and their use in animal genetics.** *Genet Sel Evol* 2002, **34(3):275-306.**
- 270 19. Rafalski A: **Applications of single nucleotide polymorphisms in crop genetics.** *Curr Opin*
271 *Plant Biol* 2002, **5(2):94-100.**
- 272 20. V A Li UL, Brandstr O M M, Johansson M, Ellegren H: **Insertion-deletion polymorphisms**
273 **(indels) as genetic markers in natural populations.** *BMC genetics* 2008, **9(1):8.**
- 274 21. Kong Y: **Btrim: a fast, lightweight adapter and quality trimming program for**
275 **next-generation sequencing technologies.** *Genomics* 2011, **98(2):152-153.**
- 276 22. Iqbal Z, Caccamo M, Turner I, Flicek P, McVean G: **De novo assembly and genotyping of**
277 **variants using colored de Bruijn graphs.** *Nat Genet* 2012, **44(2):226-232.**
- 278 23. Chaisson MJ, Brinza D, Pevzner PA: **De novo fragment assembly with short mate-paired**
279 **reads: Does the read length matter?** *Genome Res* 2009, **19(2):336-346.**
- 280 24. Kaplan PD, Ouyang Q, Thaler DS, Libchaber A: **Parallel overlap assembly for the construction**
281 **of computational DNA libraries.** *J Theor Biol* 1997, **188(3):333-341.**
- 282 25. DiGuistini S, Liao NY, Platt D, Robertson G, Seidel M, Chan SK, Docking TR, Birol I, Holt RA,
283 Hirst M *et al*: **De novo genome sequence assembly of a filamentous fungus using Sanger, 454**
284 **and Illumina sequence data.** *Genome Biol* 2009, **10(9):R94.**
- 285 26. Koressaar T, Remm M: **Enhancements and modifications of primer design program Primer3.**
286 *Bioinformatics* 2007, **23(10):1289-1291.**
- 287 27. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y *et al*: **SOAPdenovo2:**
288 **an empirically improved memory-efficient short-read de novo assembler.** *Gigascience* 2012,
289 **1(1):18.**
- 290 28. Langmead B, Trapnell C, Pop M, Salzberg SL, Others: **Ultrafast and memory-efficient**
291 **alignment of short DNA sequences to the human genome.** *Genome Biol* 2009, **10(3):R25.**

293

294 **Figure 1: mInDel workflow.**

295 Graphic summary of mInDel workflow: the pipeline accepts NGS data as input and then
296 proceeds automatically to perform several independent analyses, most of which can be
297 selected or excluded according to the user's needs. Module 1: Pre-process analysis.
298 Module 2: De novo assembly (Illumina or 454 sequencing platform). Module 3: Overlap
299 primer design. Module 4: ePCR mapping and specificity screening. Module 5: InDel
300 screening and marker development.



301

302 Table 1: Summary of three genotypes genome data.

Types	B73(V3)	Mo17	Qi319	Zheng58
Raw data	/	24.8 Gb	30.4 Gb	12.5 Gb
Sequencing platform	/	454 sequencing	Illumina paired-end	Illumina 100 bp*2
Cleaned data	/	21.2 Gb	28.6 Gb	10.1 Gb
Assembled size	2.1G	1.5G	408.9 M	183.6 M
No. of Contigs	/	1,351,472	476,470	395,075

Longest contig	/	33,817 bp	18,068 bp	7,286 bp
N50 size	/	2,012 bp	1,157 bp	457 bp

303

304 Table 2: Summary of InDel markers developed from two groups.

Types	Mo17vsB73 (with reference genome)	Qi319vsZheng58 (without reference genome)
Total InDels	739,250	223,084
Mean size	71	66
1~10bp	76,376	25,144
11~20bp	63,714	21,291
21~100bp	374,847	123,788
>100bp	213,698	52,861
sampling polymorphic rate	86.8% (276 in 318)	72.6% (437 in 602)

305