1   **Genome-wide approaches and technologies to assess human variation**

2   Claudia Gonzaga-Jauregui [1,2]

3   [1]Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, USA.

4   [2]Center for Human Disease Modeling, Duke University, Durham, NC, USA (current affiliation).

5

6   **Abstract**

7   Current genome-wide technologies allow interrogation and exploration of the human genome

8   as never before. Next-generation sequencing (NGS) technologies, along with high resolution

9   Single Nucleotide Polymorphisms (SNP) arrays and array Comparative Genomic Hybrization

10   (aCGH) enable assessment of human genome variation at the finest resolution from base pair

11   changes such as simple nucleotide variants (SNVs) to large copy-number variants (CNVs). The

12   application of these genomic technologies in the clinical setting has also enabled the molecular

13   characterization of genetic disorders and the understanding of the biological functions of more

14   genes in human development, disease, and health. In this review, the current approaches and

15   platforms available for high-throughput human genome analyses, the steps involved in these

16   different methodologies from sample preparation to data analysis, their applications, and

17   limitations are summarized and discussed.

18

19    **Next-generation massively parallel sequencing for high-throughput human genome analysis**

20    Since the beginning of the Human Genome Project, sequencing of the human genome has

21    driven the development of technologies and methods to discover the variation within it. While

22    the sequencing of the original human haploid consensus reference genome cost an estimated

23    $2.7 billion US dollars [1], the subsequent development of better and more efficient sequencing

24    machines and methodologies, and later the development of massively parallel sequencing and

25    next-generation sequencing technologies, dramatically reduced both the cost and time to

26    sequence personal human genomes. In addition, the development and refinement of targeted

27    capture methods and reagents for exome sequencing has resulted in a rapid increase in the

28    number of human exomes analyzed dramatically expanding our knowledge of human genetic

29    variation. Concurrently, bioinformatic algorithms and tools have been developed to manage and

30    analyze the tremendous amount of data generated.

31    The process to sequence a human genome or exome is now relatively straightforward and the

32    methodological differences arise mainly from the preferred capture, amplification, and

33    sequencing platforms used. In summary, the process can be reduced to four different steps

34    (Figure 1): 1) DNA preparation, 2) library construction, 3) sequencing, and 4) analysis.

35

36    1)      **DNA preparation.** Human genomic DNA can be isolated from different sources; generally

37    peripheral blood is preferred as the starting biological material. However, available reagents to

38    stabilize other biological fluids such as saliva have proven to be useful when a blood draw is not

39    possible or insufficient, providing an adequate yield and quality of DNA for sequencing when

40    collection is done properly. Extraction of DNA from tissue biopsies, preferably fresh tissues, can

41    be performed by first digesting the tissue using proteinase K. DNA extraction from formalin-

42    fixed, paraffin-embedded (FFPE) tissues is, although possible, often suboptimal in yield and

43    molecular weight integrity. The DNA yield for tissue biopsies is lower due to the amount of

44    starting available material and there is a risk of DNA degradation during extraction and

45    purification.  Current next-generation technologies allow the preparation of sequencing libraries

46    with as little as 1 ug of genomic DNA.

47

48    **2) Library Preparation.** After extraction and purification, genomic DNA is fragmented by

49    mechanical methods, such as nebulization or sonication, into fragments of ~200-400 bp.

50    Sonication is usually preferred over nebulization because the amount of input DNA is less and

51    the fragment size is more consistent. Fragment ends are enzymatically repaired and adaptors,

52    which can be barcoded, are ligated to the ends.  For whole-genome sequencing, the whole-

53    genome shotgun library is amplified and subjected to next-generation sequencing. For exome

54    sequencing, or other targeted approaches, target capture and enrichment is implemented prior

55    to amplification. Prior to massively parallel sequencing,  human genomic sequencing already

56    used target enrichment approaches, but these were laborious and not highly scalable methods

57    such as PCR for specific segment amplification or cloning of discreet genomic segments using

58    bacterial vectors and including fosmid and BAC library construction. The currently used targeted

59    capture methodologies were originally developed using oligonucleotide probes covalently

60    bound to a solid array glass slide designed to specifically bind target regions or the exons of the

61    target genes [2-5]. Later, solution based capture was developed [6,7] in which target fragments

3

62    specifically hybridize to biotinylated probes that are then pulled down  using streptavidin coated

63    magnetic beads. During exome capture hybridization, it is important to block repetitive DNA,

64    using human Cot-I DNA which is added in excess in the hybridization solution, in order to avoid

65    nonspecific cross-hybridization. The fragmented genomic DNA hybridizes to the complementary

66    probes either on the array or to the biotinylated oligonucleotide probes in solution; any non-

67    target fragments that do not hybridize are later washed away and consequently not captured

68    for subsequent sequencing. The capture efficiency is dependent on the target fragment length,

69    sequence complexity, and GC content of the region. Solution-based capture is cheaper and

70    more scalable than microarray-based capture; thus most commercially available exome capture

71    reagents use solution based capture [8].

72

73     Importantly, different targeted capture designs exist based upon the genome/gene

74    annotation(s) used for design. Additionally, the above mentioned methodologies are well suited

75    for enriching the "whole exome" with capture libraries of ~50 Mb. However if a more targeted

76    approach is desired, such as those for gene panels or specific regions, approaches such as

77    molecular inversion probes (MIPs) and other multiplex or modified PCR-based amplification of

78    targets can be used to enrich for the desired regions on a reduced  scale [9]. Originally MIPs

79    were developed, improved and applied for high-throughput multiplex SNP genotyping [10, 11].

80    The current MIPs technology relies on the specific design of ~70-mer capture probes. The MIP

81    structure is composed of a common linker sequence flanked by homologous targeting arms that

82    hybridize upstream and downstream to the genomic region of interest. A synthesis reaction

83    follows in which a DNA polymerase copies the target sequence using the upstream targeting

4

84    arm as an extension primer. After extension, the 5' end is then ligated to the downstream

85    targeting arm and the probe is circularized. Further post-capture library amplification, barcoding

86    and sequencing adaptor ligation can later be performed using the common MIP linker sequence

87    [12-14]. Current MIPs designs and approaches have proven effective at capturing ~55,000

88    targets or ~6 Mb [13, 14].

89

90    Most of the current sequencing technologies rely on the amplification of the template to be

91    sequenced in order to form clusters of clonally amplified molecules termed "polonies"; which

92    derives from PCR and colony referring to the original bacterial colonies needed to amplify DNA

93    BACs for sequencing. There are currently two predominant approaches for pre-sequencing

94    library amplification. Emulsion PCR amplification is performed in a water-oil emulsion that

95    contains the captured fragment library, dNTPs, polymerase, and beads with oligonucleotide

96    primers complementary to the adaptors initially ligated to the DNA fragments. In the test tube,

97    each of the spheres formed by the water-oil emulsion will perform as an individual isolated PCR

98    reaction in which the template fragments will be clonally amplified. These beads will then be

99    washed and cross-linked or spread into a slide or solid platform in which the sequencing

100   reaction will be performed. The second approach is a solid-phase amplification, in which pairs of

101   oligonucleotide amplification primers are covalently bound to a solid phase and the template

102   amplification takes place by bridge amplification of the target fragment using a pair of primers

103   and generating clusters of clonally amplified target molecules.

104   The efficiency of the targeted capture enrichment step can be easily assessed by qPCR, testing a

105   few target loci in the initial non-amplified input DNA versus the amplified captured DNA and

106    comparing their $C_T$ values. However, this QC step does not provide information on the specificity

107    and sensitivity of the capture method, just the efficiency of the enrichment [5].

108

109    **3) Sequencing.** Sequencing technologies can be divided into two main categories based on the

110    enzyme that they use: i) sequencing by ligation, using a DNA ligase; and ii) sequencing by

111    synthesis, using a DNA polymerase. Sequencing by synthesis is the most commonly used

112    approach and includes Sanger dideoxy sequencing. For sequencing by synthesis next-generation

113    technologies, the distinctions between methods relate to the output signal that is detected

114    when the nucleotide incorporation occurs. We will review Sanger first generation DNA

115    sequencing and the most common and widely used massively parallel next-generation

116    sequencing technologies. Detailed reviews of additional next-generation sequencing

117    technologies are available [15, 16].

118

119    **Sanger dideoxy sequencing.** Sanger dideoxy sequencing [17] remains the gold standard for DNA

120    sequencing due to its high accuracy and read length of ~1 kb; however, the cost, time and

121    scalability of Sanger sequencing make it unfeasible for large-scale sequencing. Originally, Sanger

122    dideoxy sequencing was developed using radiolabeled chain terminating dideoxy nucleotides

123    (ddNTPs) that were individually included in four separate sequencing reactions along with

124    normal unlabeled deoxynucleotides and when incorporated would stop the polymerization

125    reaction; the dideoxy nucleotides competitively inhibit the synthesis reaction of DNA

126    polymerase I. The four separate polymerization reactions were electrophoresed through

127    polyacrylamide gels and the amplified template fragments migrate by molecular weight due to

128    differences in the number of nucleotides the primer was extended. The fragments that stopped

129    first due to the addition of a given ddNTP would be shorter and migrate faster, while longer

130    fragments would migrate more slowly; in this way the sequence of the DNA template could be

131    deduced. Modifications of Sanger dideoxy sequencing came with the utilization of four-color

132    fluorescently-labeled dideoxy nucleotides instead of radioactive ones. These allowed for all four

133    chemistries and capillary electrophoreses to be run simultaneously in the same lane using laser

134    detection to determine the interrogated base; when the truncated fragments pass through the

135    sequencer, the laser excites the fluorophore and the signal of each of the four fluorophores is

136    detected and recorded in a chromatogram.

137

138    **Sequencing by Oligonucleotide Ligation and Detection.** Library amplified fragments are bound

139    through adaptors to a sequencing flow cell slide.  Sequencing by ligation is initiated by the

140    hybridization of a first of five universal primers and then by adding a pool of fluorescently

141    labeled 8-mer oligonucleotides that are labeled depending on their two last base pairs. This

142    produces sixteen different dinucleotide combinations labeled by four different fluorophores on

143    their 5' end. During the sequencing reaction, only the oligonucleotide that is complementary to

144    the template strand will hybridize, bind and be ligated to the nascent strand by a DNA ligase.

145    Four-color imaging is performed by exciting each of the fluorophores and detecting the

146    fluorescent signal across all the spots in the flow cell slide. Silver ions are flushed in order to

147    cleave the recently ligated oligonucleotides releasing the fluorophore and leaving the 5'-PO end

148    of the oligonucleotide free to bind the next one. The cycle is repeated nine times, after which

149    the universal primer with the extended strand is stripped. The next universal primer is used to

150     start the next cycle of sequencing, which again is repeated five times. There are five universal

151     primers used.  The sequencing is said to be performed in color space, which for downstream

152     analysis requires mapping to a color-space reference sequence in order to infer the nucleotide

153     sequence [18, 19].

154

155     **Pyrosequencing.** After emulsion PCR library amplification, the beads are arrayed into a picotiter

156     plate (PTP) that contains millions of micro wells large and deep enough (44 um x 55 um) only to

157     hold a single bead containing a single amplified molecule per well. Smaller beads with

158     sulphurylase and luciferase enzymes attached, necessary for the later pyrosequencing reactions,

159     are flushed and allowed to diffuse into the wells and cover the target beads. Each of the dNTPs

160     is individually flushed one at a time through the PTP and they diffuse through each of the

161     sequencing wells.  When the DNA polymerase incorporates a nucleotide, a pyrophosphate

162     group is released which will be converted by sulphurylase into ATP to phosphorylate luciferase

163     into luciferin. The light produced by luciferin due to the specific incorporation of a nucleotide in

164     that cycle will be recorded by the CCD camera in the machine, producing an output known as a

165     flowgram or pyrogram. The height of the peak is proportional to the bioluminescence signal

166     intensity which in turn is proportional to the number of incorporated nucleotides in that cycle.

167     However, this is both an advantage and disadvantage of the system, as the specificity of the

168     incorporated base is greater but the detector can be saturated by the signal if more than 6

169     nucleotides are incorporated in the same cycle, making it inaccurate for sequencing

170     homopolymer tracts. Between cycles, there is a wash with apyrase, an enzyme that degrades

171     any remaining unincorporated nucleotides and ATP produced from the previous cycle [20-23].

172

173 **Reversible terminators.** Sequencing by synthesis using reversible terminators uses clusters

174 generated by bridge amplification on an eight lane flow cell slide. The sequencing cycle starts

175 with flushing a mixture of four fluorescently labeled 3'-modified nucleotide terminators and an

176 engineered DNA polymerase that is able to incorporate these modified nucleotides. If the

177 nucleotide is complementary to the next base in the primed template, it will be added by the

178 polymerase; the extension will be blocked and the fluorescent signal derived by laser excitation

179 of each of the fluorophores will be detected by a high resolution camera. After imaging, the

180 terminating group of the modified nucleotides is cleaved along with the fluorophore allowing

181 the regeneration of the 3'-OH for the addition of the next specific nucleotide and starting a new

182 cycle. The presence of this terminator is key to this technology's chemistry as it ensures that no

183 additional or nonspecific nucleotides are added in the same cycle, allowing that just one

184 nucleotide per template is imaged per cycle. All the four-color images are processed in order to

185 derive the actual nucleotide sequence [24, 25].

186

187 **Semiconductor Ion Sequencing.** The most recent next-generation sequencing by synthesis

188 technology is based on detecting the hydrogen ion that is released during the DNA synthesis

189 reaction by a very sensitive pH meter – a microchip sensor. After template amplification by

190 emulsion PCR, template bound acrylamide beads are loaded into the semiconductor chip's

191 wells. Nucleotides are allowed to flow through the chip one at a time. When the DNA

192 polymerase incorporates the next complementary dNTP, the reaction produces pyrophosphate

193 and hydrogen due to the hydrolysis of the triphosphate of the incorporating nucleotide. The

194    hydrogen ion released produces a change in pH proportional to the number of dNTPs

195    incorporated in that given nucleotide flow cycle that can be detected by a tantalum oxide

196    coated sensor, which provides increased proton sensitivity. The 0.02 pH change per nucleotide

197    that is incorporated is registered by the sensor, then converted into a voltage signal that is

198    finally digitalized to a sequence output [26, 27].

199

200    **Single molecule sequencing.** The next-next-generation of sequencing technologies involves

201    single-molecule sequencing. The first of these technologies performs real-time single-molecule

202    sequencing, in which individual DNA polymerases are attached to the bottom of nanophotonic

203    platforms (zero-mode waveguide, ZMW detectors) that can sensitively detect the binding of

204    fluorescently phospho-linked dNTPs to the nascent strand in real time. The template DNA is

205    diluted so that only one molecule will be sequenced by one polymerase in each of the wells.

206    When the nucleotide is in the active site of the polymerase, a pulse of fluorescence in the

207    specific wavelength is detected by the ZMW detector. Once the new correct nucleotide is

208    covalently bound by the DNA polymerase, the fluorophore is released, the pulse ends and the

209    recently incorporated nucleotide is left free for the next dNTP incorporation. The processivity of

210    the DNA polymerase allows the sequencing of several hundreds of base pairs using this

211    technology [28, 29].

212

213    2)      **Analysis.** Although none of the next-generation sequencing technologies has reached

214    the accuracy of classic Sanger dideoxy sequencing, they compensate by several fold redundancy

215    of sequencing and essentially oversampling of the same genomic region thereby reducing the

216    noise and error background. However, these massive amounts of data generated by the next-

217    generation sequencing technologies pose different analytical challenges as current algorithms

218    must process information from millions of short sequence reads and deduce variant information

219    contained in these in the context of a complex genome.

220    After the chemistries and 'wet bench' sequencing process, the data generated by the different

221    technologies is exported into sequence files which generally contain the sequence of each of the

222    reads generated plus some encoded quality information for that read. These sequence files are

223    assembled into contiguous genomic sequence and mapped to the reference genome sequence.

224    Through the use of current next-generation sequencing technologies, most of the human

225    genome sequencing projects are in fact re-sequencing projects, meaning that they rely on a

226    haploid reference genome sequence assembly to map and align the sequence reads produced

227    from any individual personal genome and identify variants determined by differences from the

228    haploid human genome reference. Because of the inability to assemble an entire genome from

229    short read sequences without a reference scaffold, all the individual sequencing read data

230    points that do not map to the reference genome used are generally discarded along with

231    duplicate and low quality reads.

232    There are different alignment algorithms which can be used for mapping sequence reads to the

233    reference genome. Mapping algorithms vary in their approaches and how exhaustive their

234    mapping is, which reflects both the accuracy and computational speed with which they can be

235    implemented. Alignment algorithms can be broadly divided into those that build a 'hash' or

236    associative array of either the reference genome or the sequence reads to use as seeds or

237    anchors for the alignment, once the seeds have been aligned to the reference genome, a

238    smaller local Blast–like alignment is performed in order to extend the alignment and ensure

239    more accurate mapping. The second group of algorithms is formed by those that utilize the

240    Burrows-Wheeler transform (BWT) algorithm in which the reference genome is sorted,

241    reordered, and indexed for more efficient access and read alignment, which makes these

242    algorithms faster [30, 31].  The output of these algorithms is generally a sequence

243    alignment/map (SAM) file. This file is generally very large as it contains the mapping of each

244    read to the reference and its quality. Therefore, for better and more efficient handling of this

245    information, SAM files are converted into binary alignment files (BAM) which can be more

246    readily accessed, read and handled by the computer [32]. From the BAM file, information about

247    each base in the genome can be extracted and this is deposited into a Pileup file. This file

248    reports for every position in the genome the base observed by the pilling up of all the reads at

249    that specific position (Figure 2). However, the majority of the bases in the genome will be

250    identical to the reference, therefore of main interest are those positions that are different. After

251    obtaining the pileup for the whole genome, the variable positions are extracted into another file

252    in a process known as variant calling. It is important to evaluate these variant calls for their

253    quality, number of reads across the position and number of reads that called the variant, strand

254    bias, and likelihood of being a true variant, in order to ensure that the majority are true

255    positives. The quality of a variant call can be generally assessed by its pileup information.

256    There are several algorithms that perform variant calling from next-generation sequencing data,

257    they provide genotype data and can be used for quality filtering [33-35]. Once variants are

258    "called', these can be analyzed and filtered through different approaches and using different

259    criteria depending on the purpose of the study. In general, whether the variants detected map

12

260 to genic or intergenic regions, if they are coding or non-coding, as well as whether they

261 represent previously observed polymorphisms or novel variants, as well as the allele frequency

262 spectrum in populations are all important parameters to consider. This is achieved through a

263 process known as variant annotation. When annotating variants, several available databases are

264 queried in order to gather as much information as possible regarding that specific genomic

265 position or coordinate in order to assess the functional impact of the identified variant. Some of

266 the most widely used databases included in several annotation pipelines are:

267

268 *Population Polymorphism Databases:* The database of single nucleotide polymorphisms

269 (**dbSNP**) was established in 1998 by the National Center for Biotechnology Information (NCBI) in

270 order to store and catalogue single nucleotide polymorphisms (SNPs) as the most common form

271 of genetic variation between individuals [36, 37]. Since its creation, the database has been

272 greatly expanded to include simple nucleotide variants (SNVs) both SNPs and Indels. Initially the

273 database was populated by the SNPs discovered during the HGP, later by the additional variants

274 discovered during the HapMap Project, and most recently by variants of the Thousand Genomes

275 Project and the myriad of whole genome and exome sequencing projects of the last lustrum.

276 The **1000 (Thousand) Genomes Project (TGP)** was initiated with the purpose of cataloguing

277 most of the polymorphic and low-frequency variation amongst human genomes, including SNVs

278 and copy-number variants, by sequencing >1000 genomes and exomes of different populations

279 around the world using NGS technologies [38-40]. The most recent data release of the **NHLBI**

280 **Exome Sequencing Project (ESP6500)** contains SNP variants identified through whole exome

281 sequencing of 6503 samples from multiple NHLBI ESP cohorts [41]. In addition, most large-scale

282  sequencing groups use their own internal variant databases to annotate for in-house variant

283  frequencies and genotypes observed, which helps to control for technical replicate errors and

284  specific population variation.

285

286  ***Conservation scores:*** **PhyloP** computes p-values of nucleotide conservation based on a

287  tree model of neutral evolution and multi-species alignments [42, 43]. The likelihood ratio test

288  (**LRT**) considers all possible ancestral sequences to estimate the 'deleteriousness' of a particular

289  substitution based  on a comparative genomic approach across 32 vertebrate species and

290  assuming neutrality from synonymous changes and treating all nucleotide substitutions equally

291  [44]. The Genomic Evolutionary Rate Profiling (**GERP**) approach aims to identify evolutionary

292  constrained regions that have lower nucleotide substitution rates, which may reflect past

293  purifying selection, through the sequence analysis and alignments of 29 mammalian species

294  [45]. **PhastCons** attempts to identify evolutionarily conserved regions through multiple species

295  alignments (46 placental mammals) based on a phylogenetic Hidden Markov Model (phylo-

296  HMM) that uses statistical models for unequal nucleotide substitution rates [46].

297

298  ***Functional prediction algorithms:*** The 'Sorting Tolerant From Intolerant' (**SIFT**) algorithm

299  predicts the functional effect of an amino acid substitution based on the conservation of that

300  amino acid residue in the protein through multiple sequence alignment of closely related

301  proteins from PSI-BLAST [47]. The scores and predictions given by SIFT range from (1 = tolerated

302  to 0 = damaging); the amino acid change is predicted to be damaging if the score <= 0.05, and

303  tolerated if the score > 0.05.

304  The 'Polymorphism Phenotyping' (**PolyPhen**) algorithm predicts the possible functional impact

305  of an amino acid substitution based on the protein structure, phylogenetic conservation and

306  sequence information calculating a naive Bayes posterior probability that the mutation might be

307  damaging [48]. The score reported by PolyPhen reports the probability of the mutation being

308  damaging when it is not damaging over the probability of the mutation being damaging when it

309  is actually damaging; therefore, the scores range from 'problably damaging' (score=0) to

310  'benign' (score=1).

311  **MutationTaster** utilizes a Bayes probabilistic algorithm to predict the functional impact of a

312  given nucleotide change, either SNPs or small indels, based on a training set of known disease

313  causing mutations and common polymorphisms. This algorithm calculates the probability of the

314  change being a polymorphism or a damaging mutation and reports back p-values (not scores) of

315  the prediction being correct [49].

316

317  *Disease/phenotype related databases:* The Online Mendelian Inheritance in Man

318  (**OMIM**) database is a compendium of human diseases and phenotypes of genetic or suspected

319  genetic etiology [50]. To date there are  about 3,700 genetic phenotypes or diseases for which

320  the molecular cause is  known; however there are still at least ~4,000 phenotypes with

321  suspected genetic basis for which the responsible gene(s) remain unknown. The Human Gene

322  Mutation Database (**HGMD**) is a catalogue of known specific mutations reported to be

323  associated with particular genetic diseases or phenotypes [51]. Currently, HGMD includes more

324  than 134,000 single mutations associated to genetic diseases. Unfortunately, some of these

325  mutations are based on single case reports with insufficient support for pathogenicity of the

326  mutation; new exome data in thousands of individuals support the contention that such variants

327  are not damaging or causing the disease they were reported to be associated with as they reach

328  a certain frequency (>2 – 5%) within normal populations.

329  The Pharmacogenomics Knowledge database (**PharmGKB**) documents reported, well

330  characterized variants and polymorphisms related to the metabolism of a wide variety of drugs

331  and compounds. It is a valuable resource to inform analysis of potential medically actionable

332  variants associated with drug metabolism and suggested dosage and guidelines for the

333  utilization of common medications depending on an individual's genotype [52].

334

335  *Other databases and resources:* Other useful information resources to interpret variants

336  in novel genes can be provided by pathway and protein-protein interaction network databases,

337  such as the Kyoto Encyclopedia of Genes and Genomes (**KEGG**) [53, 54] for biological pathways,

338  and the Human Integrated Protein-Protein Interaction rEference (**HIPPIE**) [55] or the Search Tool

339  for the Retrieval of Interacting Genes/Proteins (**STRING**) [56] databases for protein-protein

340  interactions data. These databases can provide information on which gene products directly

341  interact with each other or function upstream or downstream of other known disease

342  associated genes and help elucidate possible functions and roles of novel genes. Additionally,

343  one would like to be informed of expression profiles of the genes of interest across different

344  tissues or throughout development [57, 58]; and phenotypes of mutant model organisms for the

345  genes of interest [59].

346

**Array Comparative Genomic Hybridization (aCGH) as the gold standard for Copy-Number**

**Variant detection**

Next-generation sequencing is allowing the identification of the vast majority of simple

nucleotide variants (SNVs) in personal genomes, from single nucleotide changes to small

insertion/deletion variants. However, one additional, larger in scale and highly important source

of variation in genomes, both polymorphic and rare is structural and copy number variation

(CNV). Structural variation refers to segments of the genome that differ in copy-number,

architecture and/or orientation between genomes; within this, copy-number variants (CNVs) are

regions in the genome that are present in more or less copies than the expected diploid state. It

is now widely recognized that structural variation and CNVs contribute largely to human

genomic variation, both benign polymorphic and pathogenic [60].

Array Comparative Genomic Hybridization (aCGH) remains the gold standard to identify CNVs in

the genome. It is based on the competitive hybridization between a test or proband's DNA and a

control DNA. First, the samples can be digested (or not, depending on the protocol) using an

*Alu*I/*Rsa*I enzyme mix in order to digest the DNA into fragments of ~500 bp in size. Labeling is

performed by priming with random nonamers and then adding fluorescent cyanine dyes, either

Cy3 or Cy5, to the test DNA and the control DNA, respectively, that are incorporated by a nick-

translation reaction using a Klenow fragment polymerase. Competitive hybridization is then

performed on an array platform which contains probes to interrogate the whole genome or

specific regions for several hours depending on the protocol. After extended hybridization, the

array is washed and scanned using lasers that excite each of the fluorophores and the signal

intensity of each probe in the array is measured. The images obtained can then be processed,

369 merged, and analyzed. If the signal intensity for a particular probe is greater when exciting the

370 fluorophore of the test DNA, it indicates that more test DNA was bound and there is a copy-

371 number gain of that probe in the test DNA. Conversely, if the signal intensity for a different

372 probe is greater when exciting the flourophore of the control DNA, then this indicates a copy-

373 number loss of that region in the test DNA. If the signal intensity for both fluorophores is about

374 the same for a given probe, then it is assumed that there is no copy-number change between

375 the test and control DNAs for the region being interrogated by the oligonucleotide probe [61,

376 62].

377 Initially, the grid array platforms used BACs to tile the genome and interrogate for CNVs,

378 however this made the estimation of CNV size and boundaries difficult, inaccurate and

379 overestimated. With the development of the aCGH platforms, and a reference human genome

380 sequence to design probes, came the tiling of the genome using covalently bound

381 oligonucleotide (~60mers) probes that could be chosen to tile the whole genome or specific

382 regions at higher resolution; the major constraint in genome resolution afforded being the

383 number of probes utilized per array design; i.e. the number of 'pixels' used for genome

384 resolution. These currently range from as few as 60,000 probes to as many as 4.2 million probes

385 depending on the platform used. The performance of the oligonucleotide probes can vary

386 depending on their GC content and therefore it is usually necessary to have a change in the

387 same direction in the signal of at least five continuous probes in order to be able to detect a

388 signal representing a true CNV by aCGH.

389 aCGH has proven to be very accurate and successful at identifying large CNVs, however it has its

390 limitations including: i) not being able to resolve CNVs of less than 5kb genome-wide, ii) relying

18

391   on a "control" DNA, which by itself will have CNVs of its own of unknown significance, iii) not

392   being able to detect copy-number changes of more than 4 due to the signal's dynamic range, iv)

393   relying on a reference sequence on which the array design is based and not being able to detect

394   other types of structural variation such as insertions, inversions and balanced translocations.

395   Currently, copy number variation can be deduced from NGS data based on the coverage

396   distribution of single reads across the genome or target regions and the difference between the

397   number of reads across (i.e. read-depth-of-coverage) a given region and the genome-wide

398   average. Alternatively, CNVs and some structural variants can be inferred from NGS data

399   through the library preparation, sequencing and analysis of paired-end reads. Paired-end

400   sequencing (PE) reads are generated from both ends of a fragment. Because the distance

401   between the two ends is known based on the reference genome assembly, it is possible to use

402   this information to estimate the presence of insertion, deletions, or copy number variants in the

403   region [63]. As it is also the case for aCGH, inference of deletions is in general more robust than

404   the identification of copy-number gains. In addition, sequencing data can potentially provide

405   CNV breakpoint data that can help in further understanding the mechanisms of CNV formation

406   in the human genome [64].

407   It is anticipated that eventually genomic sequencing using next-generation technologies will

408   provide accurate information on smaller CNVs, in addition to the larger CNVs readily detected

409   by aCGH, and potentially other structural variants like inversions and novel sequence insertions.

410   However for that to occur, current and novel algorithms for alignment mapping and *de novo*

411   assembly need to be improved and developed. The interpretation of mechanisms for generating

412   complex rearrangements and deducing their structure from short read sequences, and *de novo*

413    assembly algorithms for larger genomic segments and entire genomes, will have to await further

414    developments and refinements.

415

416    **Single Nucleotide Polymorphism (SNP) detection for genome-wide genotyping and CNV**

417    **detection**

418    One of the initial types of genomic variation that was attainable for genome-wide testing were

419    Single Nucleotide Polymorphisms (SNPs), which by definition are, in most cases, bi-allelic

420    positions in the genome that differ in one nucleotide among individuals and are present in at

421    least 1% of the general population. SNP genotyping has been widely used for genome-wide

422    association studies (GWAS) with the intention to find common SNP associations to common

423    complex diseases. SNP detection platforms perform allelic discrimination to interrogate the

424    polymorphic position through different approaches, which can be primer extension by single

425    base incorporation, mismatch hybridization, ligation, and enzymatic cleavage [65].  The primer

426    extension approach can utilize a common primer that can detect either allele or allele-specific

427    primers; primer anneals to the contiguous region next to the interrogated SNP and the

428    nucleotide corresponding to the SNP is incorporated in an allele-specific PCR reaction and

429    identified by either mass spectrometry or fluorescence. In the mismatch hybridization approach,

430    allele-specific oligonucleotide probes are printed and arrayed in a solid support. Genomic DNA is

431    digested, PCR amplified, fragmented, labeled and hybridized to the array. Each SNP is

432    independently tested by several probes that differ just at the SNP position and can discriminate

433    between each of the two alleles by fluorescence detection. Another mismatch based approach

434    utilizes allele-specific primers that have a fluorophore and a quencher and a common primer for

435    PCR amplification. During the primer extension step, the polymerase with 5' exonuclease

436    activity can cleave the perfectly matched primer freeing the reporter fluorophore from the

437    quencher's proximity and genotype is detected by the flourescent signal emitted. Other

438    mismatch based approaches using allele-specific primers tagged differently have also been

439    developed in order to detect SNP genotypes by other methods besides fluorescence, such as

440    mass spectrometry or flow cytometry. For ligation-based methods, two allele-specific

441    oligonucleotide probes are used in addition to a common oligonucleotide that binds adjacently

442    to the SNP site. When one of the allele-specific probes binds to the SNP site perfectly, a DNA

443    ligase will ligate the specific probe with the common oligonucleotide. The ligated allele-specific

444    products can then be detected depending on what was used to tag the different probes, most

445    commonly fluorescent dyes. A variation of this approach is using longer linear oligonucleotide

446    probes whose ends are equivalent to the allele-specific and the common probes. The approach

447    is the same, but when allele-specific binding and ligation occurs the probe gets circularized. This

448    circular probe can then be amplified by rolling circle DNA replication or  PCR; the genotype can

449    be 'called' using fluorescently labeled primers or by fluorescent labeling during the amplification

450    step with subsequent array hybridization. Enzymatic cleavage methods rely on the specificity of

451    DNA endonucleases to recognize specific (which for SNP genotyping are allele-specific)

452    sequences and cleave the DNA evidencing the genotype. Region of interest PCR amplified

453    products can be incubated with specific restriction enzymes and genotypes can be detected by

454    differences in fragment sizes. However, this method is generally low-throughput and limited by

455    the nucleases sequence recognition repertoire. A variation of the enzymatic cleavage approach

456  uses two fluorescently labeled allele-specific probes with an additional common "invader"

457  probe that is complementary to the 3' end of the SNP region. When hybridization to the SNP

458  site occurs, the presence of the invader probe creates an overhang of the allele-specific

459  oligonucleotide that is recognized by a nuclease, which cleaves it and releases the fluorophore

460  for genotype detection.

461  Currently, SNP detection platforms are mostly array based and aimed for high-throughput,

462  genomewide genotyping with varying number of SNP test probes depending on the design of

463  preference. These methodologies also use a combination of two or more of the previously

464  described allele discrimination approaches. Additionally, current SNP arrays designs also include

465  CNV detection probes or inversely, aCGH designs include SNP array probes which allow for the

466  detection of both CNVs and absence of heterozygosity (AOH) in the test sample [66].

467  SNP arrays genotype data can be analyzed based on the B-allele frequency (BAF), which in a

468  diploid state presents three possible states namely AA=0, AB=0.5, and BB= 1. These data can

469  also be visualized in a B-allele frequency plot. When the BAF deviates from the tri-modal

470  expected distribution, this can indicate allelic imbalance that can relate to genomic events such

471  as copy number variants, regions of absence of heterozygosity, or uniparental disomy, in the

472  same assay [67, 68].

473

## Conclusions

474

475    No current genome-wide technology or platform provides information about all types of

476    variation, from SNVs to structural variation including small CNV in the 100bp to 5kb range, at

477    high resolution with high specificity [63]. However, continued developments in sequencing

478    technologies and analysis promise to deliver better, faster and more high-throughput

479    sequencing technologies to assess the complete picture of human genomic variation,

480    spearheading the development of improved and new methods for data analysis, not only to

481    process the peta ($10^{15}$) amounts of data produced, but also to make biological sense of this

482    information. The experimental and technical approaches for human genomic variation discovery

483    will most probably not be the main limiting factor in assessing it, but instead our understanding

484    and our ability to derive biologically relevant lessons and conclusions from such massive data

485    will remain the premier challenge.

486    References:

487        1.  http://www.genome.gov/11006943

488        2.  Albert TJ, Molla MN, Muzny DM, Nazareth L, Wheeler D, Song X, Richmond TA, Middle CM,

489            Rodesch MJ, Packard CJ, Weinstock GM, Gibbs RA; Direct selection of human genomic loci by

490            microarray hybridization; *Nat Methods.* 2007 Nov;4(11):903-5.

491        3.  Okou DT, Steinberg KM, Middle C, Cutler DJ, Albert TJ, Zwick ME; Microarray-based genomic

492            selection for high-throughput resequencing; *Nat Methods.* 2007 Nov;4(11):907-9.

493        4.  Hodges E, Xuan Z, Balija V, Kramer M, Molla MN, Smith SW, Middle CM, Rodesch MJ, Albert TJ,

494            Hannon GJ, McCombie WR; Genome-wide in situ exon capture for selective resequencing; *Nat*

495            *Genet.* 2007 Dec;39(12):1522-7.

496        5.  Hodges E, Rooks M, Xuan Z, Bhattacharjee A, Benjamin Gordon D, Brizuela L, Richard McCombie

497            W, Hannon GJ; Hybrid selection of discrete genomic intervals on custom-designed microarrays

498            for massively parallel sequencing; *Nat Protoc.* 2009;4(6):960-74.

499        6.  Gnirke A, Melnikov A, Maguire J, Rogov P, LeProust EM, Brockman W, Fennell T, Giannoukos G,

500            Fisher S, Russ C, Gabriel S, Jaffe DB, Lander ES, Nusbaum C; Solution hybrid selection with ultra-

501            long oligonucleotides for massively parallel targeted sequencing; *Nat Biotechnol.* 2009

502            Feb;27(2):182-9.

503        7.  Bainbridge MN, Wang M, Burgess DL, Kovar C, Rodesch MJ, D'Ascenzo M, Kitzman J, Wu YQ,

504            Newsham I, Richmond TA, Jeddeloh JA, Muzny D, Albert TJ, Gibbs RA; Whole exome capture in

505            solution with 3 Gbp of data; *Genome Biol.* 2010;11(6):R62.

506    8.   Parla JS, Iossifov I, Grabill I, Spector MS, Kramer M, McCombie WR; A comparative analysis of

507         exome capture; *Genome Biol.* 2011 Sep 29;12(9):R97.

508    9.   Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A, Howard E, Shendure J, Turner

509         DJ; Target-enrichment strategies for next-generation sequencing; Nat Methods. 2010

510         Feb;7(2):111-8.

511    10.  Hardenbol P, Banér J, Jain M, Nilsson M, Namsaraev EA, Karlin-Neumann GA, Fakhrai-Rad H,

512         Ronaghi M, Willis TD, Landegren U, Davis RW; Multiplexed genotyping with sequence-tagged

513         molecular inversion probes; Nat Biotechnol. 2003 Jun;21(6):673-8.

514    11.  Hardenbol P, Yu F, Belmont J, Mackenzie J, Bruckner C, Brundage T, Boudreau A, Chow S, Eberle

515         J, Erbilgin A, Falkowski M, Fitzgerald R, Ghose S, Iartchouk O, Jain M, Karlin-Neumann G, Lu X,

516         Miao X, Moore B, Moorhead M, Namsaraev E, Pasternak S, Prakash E, Tran K, Wang Z, Jones HB,

517         Davis RW, Willis TD, Gibbs RA; Highly multiplexed molecular inversion probe genotyping: over

518         10,000 targeted SNPs genotyped in a single tube assay; Genome Res. 2005 Feb;15(2):269-75.

519    12.  Porreca GJ, Zhang K, Li JB, Xie B, Austin D, Vassallo SL, LeProust EM, Peck BJ, Emig CJ, Dahl F, Gao

520         Y, Church GM, Shendure J; Multiplex amplification of large sets of human exons; *Nat Methods.*

521         2007 Nov;4(11):931-6.

522    13.  Turner EH, Lee C, Ng SB, Nickerson DA, Shendure J; Massively parallel exon capture and library-

523         free resequencing across 16 genomes; *Nat Methods.* 2009 May;6(5):315-6.

524    14.  O'Roak BJ, Vives L, Fu W, Egertson JD, Stanaway IB, Phelps IG, Carvill G, Kumar A, Lee C,

525         Ankenman K, Munson J, Hiatt JB, Turner EH, Levy R, O'Day DR, Krumm N, Coe BP, Martin BK,

526         Borenstein E, Nickerson DA, Mefford HC, Doherty D, Akey JM, Bernier R, Eichler EE, Shendure J;

527     Multiplex targeted sequencing identifies recurrently mutated genes in autism spectrum

528     disorders; *Science* 2012 Dec 21;338(6114):1619-22.

529  15. Mardis ER; Next-generation DNA sequencing methods; Annu Rev Genomics Hum Genet.

530     2008;9:387-402.

531  16. Metzker ML; Sequencing technologies - the next generation; Nat Rev Genet. 2010 Jan;11(1):31-

532     46.

533  17. Sanger F, Nicklen S, Coulson AR; DNA sequencing with chain-terminating inhibitors; PNAS 1977

534     Dec;74(12):5463-7.

535  18. Shendure J, Porreca GJ, Reppas NB, Lin X, McCutcheon JP, Rosenbaum AM, Wang MD, Zhang K,

536     Mitra RD, Church GM; Accurate multiplex polony sequencing of an evolved bacterial genome;

537     Science. 2005 Sep 9;309(5741):1728-32.

538  19. McKernan KJ, Peckham HE, Costa GL, McLaughlin SF, Fu Y, Tsung EF, Clouser CR, Duncan C,

539     Ichikawa JK, Lee CC, et al; Sequence and structural variation in a human genome uncovered by

540     short-read, massively parallel ligation sequencing using two-base encoding; Genome Res. 2009

541     Sep;19(9):1527-41.

542  20. Ronaghi M, Karamohamed S, Pettersson B, Uhlén M, Nyrén P; Real-time DNA sequencing using

543     detection of pyrophosphate release; Anal Biochem. 1996 Nov 1;242(1):84-9.

544  21. Ronaghi M, Uhlén M, Nyrén P; A sequencing method based on real-time pyrophosphate; Science.

545     1998 Jul 17;281(5375):363, 365.

546  22. Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA, Berka J, Braverman MS,

547     Chen YJ, Chen Z, Dewell SB, Du L, Fierro JM, Gomes XV, Godwin BC, He W, Helgesen S, Ho CH,

548    Irzyk GP, Jando SC, Alenquer ML, Jarvie TP, Jirage KB, Kim JB, Knight JR, Lanza JR, Leamon JH,

549    Lefkowitz SM, Lei M, Li J, Lohman KL, Lu H, Makhijani VB, McDade KE, McKenna MP, Myers EW,

550    Nickerson E, Nobile JR, Plant R, Puc BP, Ronan MT, Roth GT, Sarkis GJ, Simons JF, Simpson JW,

551    Srinivasan M, Tartaro KR, Tomasz A, Vogt KA, Volkmer GA, Wang SH, Wang Y, Weiner MP, Yu P,

552    Begley RF, Rothberg JM; Genome sequencing in microfabricated high-density picolitre reactors;

553    Nature. 2005 Sep 15;437(7057):376-80.

554  23. Wheeler DA, Srinivasan M, Egholm M, Shen Y, Chen L, McGuire A, He W, Chen YJ, Makhijani V,

555    Roth GT, Gomes X, Tartaro K, Niazi F, Turcotte CL, Irzyk GP, Lupski JR, Chinault C, Song XZ, Liu Y,

556    Yuan Y, Nazareth L, Qin X, Muzny DM, Margulies M, Weinstock GM, Gibbs RA, Rothberg JM; The

557    complete genome of an individual by massively parallel DNA sequencing; Nature. 2008 Apr

558    17;452(7189):872-6.

559  24. Ju J, Kim DH, Bi L, Meng Q, Bai X, Li Z, Li X, Marma MS, Shi S, Wu J, Edwards JR, Romu A, Turro NJ;

560    Four-color DNA sequencing by synthesis using cleavable fluorescent nucleotide reversible

561    terminators; Proc Natl Acad Sci U S A. 2006 Dec 26;103(52):19635-40.

562  25. Bentley DR, Balasubramanian S, Swerdlow HP, Smith GP, Milton J, Brown CG, Hall KP, Evers DJ,

563    Barnes CL, Bignell HR, et al; Accurate whole human genome sequencing using reversible

564    terminator chemistry; Nature. 2008 Nov 6;456(7218):53-9.

565  26. Rothberg JM, Hinz W, Rearick TM, Schultz J, Mileski W, Davey M, Leamon JH, Johnson K, Milgrew

566    MJ, Edwards M, Hoon J, Simons JF, Marran D, Myers JW, Davidson JF, Branting A, Nobile JR, Puc

567    BP, Light D, Clark TA, Huber M, Branciforte JT, Stoner IB, Cawley SE, Lyons M, Fu Y, Homer N,

568    Sedova M, Miao X, Reed B, Sabina J, Feierstein E, Schorn M, Alanjary M, Dimalanta E, Dressman

569    D, Kasinskas R, Sokolsky T, Fidanza JA, Namsaraev E, McKernan KJ, Williams A, Roth GT, Bustillo J;

570      An integrated semiconductor device enabling non-optical genome sequencing; Nature. 2011 Jul

571      20;475(7356):348-52.

572  27. Merriman B, R D Team IT, Rothberg JM; Progress in Ion Torrent semiconductor chip based

573      sequencing; Electrophoresis. 2012 Dec;33(23):3397-417.

574  28. Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, Peluso P, Rank D, Baybayan P, Bettman B, Bibillo A,

575      Bjornson K, Chaudhuri B, Christians F, Cicero R, Clark S, Dalal R, Dewinter A, Dixon J, Foquet M,

576      Gaertner A, Hardenbol P, Heiner C, Hester K, Holden D, Kearns G, Kong X, Kuse R, Lacroix Y, Lin S,

577      Lundquist P, Ma C, Marks P, Maxham M, Murphy D, Park I, Pham T, Phillips M, Roy J, Sebra R,

578      Shen G, Sorenson J, Tomaney A, Travers K, Trulson M, Vieceli J, Wegener J, Wu D, Yang A,

579      Zaccarin D, Zhao P, Zhong F, Korlach J, Turner S; Real-time DNA sequencing from single

580      polymerase molecules; Science. 2009 Jan 2;323(5910):133-8.

581  29. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, Qin X, Muzny DM, Reid JG, Worley KC, Gibbs

582      RA; Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing

583      technology; PLoS One. 2012;7(11):e47768.

584  30. Flicek P, Birney E; Sense from sequence reads: methods for alignment and assembly; Nat

585      Methods. 2009 Nov;6(11 Suppl):S6-S12.

586  31. Bao S, Jiang R, Kwan W, Wang B, Ma X, Song YQ; Evaluation of next-generation sequencing

587      software in mapping and assembly; J Hum Genet. 2011 Jun;56(6):406-14.

588  32. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R; 1000

589      Genome Project Data Processing Subgroup; The Sequence Alignment/Map format and SAMtools;

590      Bioinformatics. 2009 Aug 15;25(16):2078-9.

591    33. Nielsen R, Paul JS, Albrechtsen A, Song YS; Genotype and SNP calling from next-generation

592        sequencing data; Nat Rev Genet. 2011 Jun;12(6):443-51.

593    34. Shen Y, Wan Z, Coarfa C, Drabek R, Chen L, Ostrowski EA, Liu Y, Weinstock GM, Wheeler DA,

594        Gibbs RA, Yu F; A SNP discovery method to assess variant allele probability from next-generation

595        resequencing data; Genome Res. 2010 Feb;20(2):273-80.

596    35. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D,

597        Gabriel S, Daly M, DePristo MA; The Genome Analysis Toolkit: a MapReduce framework for

598        analyzing next-generation DNA sequencing data; Genome Res. 2010 Sep;20(9):1297-303.

599    36. Sherry ST, Ward M, Sirotkin K; dbSNP-database for single nucleotide polymorphisms and other

600        classes of minor genetic variation; Genome Res. 1999 Aug;9(8):677-9.

601    37. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K; dbSNP: the NCBI

602        database of genetic variation; Nucleic Acids Res. 2001 Jan 1;29(1):308-11.

603    38. Marth GT, Yu F, Indap AR, Garimella K, Gravel S, Leong WF, Tyler-Smith C, Bainbridge M,

604        Blackwell T, Zheng-Bradley X, Chen Y, Challis D, Clarke L, Ball EV, Cibulskis K, Cooper DN, Fulton B,

605        Hartl C, Koboldt D, Muzny D, Smith R, Sougnez C, Stewart C, Ward A, Yu J, Xue Y, Altshuler D,

606        Bustamante CD, Clark AG, Daly M, DePristo M, Flicek P, Gabriel S, Mardis E, Palotie A, Gibbs R;

607        1000 Genomes Project; The functional spectrum of low-frequency coding variation; Genome

608        Biol. 2011 Sep 14;12(9):R84.

609    39. Mills RE, Walter K, Stewart C, Handsaker RE, Chen K, Alkan C, Abyzov A, Yoon SC, Ye K, Cheetham

610        RK, Chinwalla A, Conrad DF, Fu Y, Grubert F, Hajirasouliha I, Hormozdiari F, Iakoucheva LM, Iqbal

611        Z, Kang S, Kidd JM, Konkel MK, Korn J, Khurana E, Kural D, Lam HY, Leng J, Li R, Li Y, Lin CY, Luo R,

612        Mu XJ, Nemesh J, Peckham HE, Rausch T, Scally A, Shi X, Stromberg MP, Stütz AM, Urban AE,

613      Walker JA, Wu J, Zhang Y, Zhang ZD, Batzer MA, Ding L, Marth GT, McVean G, Sebat J, Snyder M,

614      Wang J, Ye K, Eichler EE, Gerstein MB, Hurles ME, Lee C, McCarroll SA, Korbel JO; 1000 Genomes

615      Project; Mapping copy number variation by population-scale genome sequencing; Nature. 2011

616      Feb 3;470(7332):59-65.

617   40. 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM,

618      Handsaker RE, Kang HM, Marth GT, McVean GA; An integrated map of genetic variation from

619      1,092 human genomes; Nature. 2012 Nov 1;491(7422):56-65.

620   41. Exome Variant Server, NHLBI GO Exome Sequencing Project (ESP), Seattle, WA (URL:

621      http://evs.gs.washington.edu/EVS/)

622   42. Siepel A, Pollard KS, Haussler D; New methods for detecting lineage-specific selection.

623      Proceedings of the 10th International Conference on Research in Computational Molecular

624      Biology (RECOMB) 2006: 190–205.

625   43. Lindblad-Toh K, Garber M, Zuk O, Lin MF, Parker BJ, Washietl S, Kheradpour P, Ernst J, Jordan G,

626      Mauceli E, et al; A high-resolution map of human evolutionary constraint using 29 mammals;

627      Nature. 2011 Oct 12;478(7370):476-82.

628   44. Chun S, Fay JC; Identification of deleterious mutations within three human genomes; Genome

629      Res. 2009 Sep;19(9):1553-61.

630   45. Cooper GM, Stone EA, Asimenos G; NISC Comparative Sequencing Program, Green ED, Batzoglou

631      S, Sidow A; Distribution and intensity of constraint in mammalian genomic sequence; Genome

632      Res. 2005 Jul;15(7):901-13.

633    46. Margulies EH, Blanchette M; NISC Comparative Sequencing Program, Haussler D, Green ED;

634        Identification and characterization of multi-species conserved sequences; Genome Res. 2003

635        Dec;13(12):2507-18.

636    47. Kumar P, Henikoff S, Ng PC; Predicting the effects of coding non-synonymous variants on protein

637        function using the SIFT algorithm; Nat Protoc. 2009;4(7):1073-81.

638    48. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev

639        SR; A method and server for predicting damaging missense mutations; Nat Methods. 2010

640        Apr;7(4):248-9.

641    49. Schwarz JM, Rödelsperger C, Schuelke M, Seelow D; MutationTaster evaluates disease-causing

642        potential of sequence alterations; Nat Methods. 2010 Aug;7(8):575-6.

643    50. Online Mendelian Inheritance in Man, OMIM®. McKusick-Nathans Institute of Genetic Medicine,

644        Johns Hopkins University (Baltimore, MD), URL: http://omim.org/

645    51. Cooper DN, Ball EV, Krawczak M; The human gene mutation database; Nucleic Acids Res. 1998

646        Jan 1;26(1):285-7.

647    52. Whirl-Carrillo M, McDonagh EM, Hebert  JM, Gong L, Sangkuhl K, Thorn CF, Altman RB, Klein TE;

648        Pharmacogenomics Knowledge for Personalized Medicine; Clinical Pharmacology & Therapeutics

649        2012; 92(4): 414-417.

650    53. Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M; KEGG: Kyoto Encyclopedia of Genes

651        and Genomes; Nucleic Acids Res. 1999 Jan 1;27(1):29-34.

652    54. Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M; KEGG for integration and interpretation of

653        large-scale molecular data sets; Nucleic Acids Res. 2012 Jan;40(Database issue):D109-14.

654    55. Schaefer MH, Fontaine JF, Vinayagam A, Porras P, Wanker EE, Andrade-Navarro MA; HIPPIE:

655        Integrating protein interaction networks with experiment based quality scores; PLoS One.

656        2012;7(2):e31826.

657    56. Search Tool for the Retrieval of Interacting Genes/Proteins (**STRING**): http://string-db.org/

658    57. Wu C, Orozco C, Boyer J, Leglise M, Goodale J, Batalov S, Hodge CL, Haase J, Janes J, Huss JW 3rd,

659        Su AI; BioGPS: an extensible and customizable portal for querying and organizing gene

660        annotation resources; Genome Biol. 2009;10(11):R130.

661    58. Kapushesky M, Emam I, Holloway E, Kurnosov P, Zorin A, Malone J, Rustici G, Williams E,

662        Parkinson H, Brazma A; Gene expression atlas at the European bioinformatics institute; Nucleic

663        Acids Res. 2010 Jan;38(Database issue):D690-8.

664    59. Eppig JT, Blake JA, Bult CJ, Kadin JA, Richardson JE; the Mouse Genome Database Group; The

665        Mouse Genome Database (MGD): comprehensive resource for genetics and genomics of the

666        laboratory mouse; Nucleic Acids Res 2012; 40(1):D881-86.

667    60. Stankiewicz P, Lupski JR; Structural variation in the human genome and its role in disease; Annu

668        Rev Med. 2010;61:437-55.

669    61. Pinkel D, Albertson DG; Comparative genomic hybridization; Annu Rev Genomics Hum Genet.

670        2005;6:331-54.

671    62. Brady PD, Vermeesch JR; Genomic microarrays: a technology overview; Prenat Diagn. 2012

672        Apr;32(4):336-43.

673    63. Alkan C, Coe BP, Eichler EE; Genome structural variation discovery and genotyping; Nat Rev

674        Genet. 2011; 12(5):363-76.

675    64. Onishi-Seebacher M, Korbel JO; Challenges in studying genomic structural variant formation

676        mechanisms: the short-read dilemma and beyond; Bioessays. 2011; 33(11):840-50.

677    65. Kim S, Misra A; SNP genotyping: technologies and biomedical applications; Annu Rev Biomed

678        Eng. 2007;9:289-320.

679    66. Schaaf CP, Wiszniewska J, Beaudet AL; Copy number and SNP arrays in clinical diagnostics; Annu

680        Rev Genomics Hum Genet. 2011 Sep 22;12:25-51.

681    67. Bruno DL, White SM, Ganesamoorthy D, Burgess T, Butler K, Corrie S, Francis D, Hills L,

682        Prabhakara K, Ngo C, Norris F, Oertel R, Pertile MD, Stark Z, Amor DJ, Slater HR; Pathogenic

683        aberrations revealed exclusively by single nucleotide polymorphism (SNP) genotyping data in

684        5000 samples tested by molecular karyotyping; J Med Genet. 2011 Dec;48(12):831-9.

685    68. Yau C, Holmes CC; CNV discovery using SNP genotyping arrays; Cytogenet Genome Res.
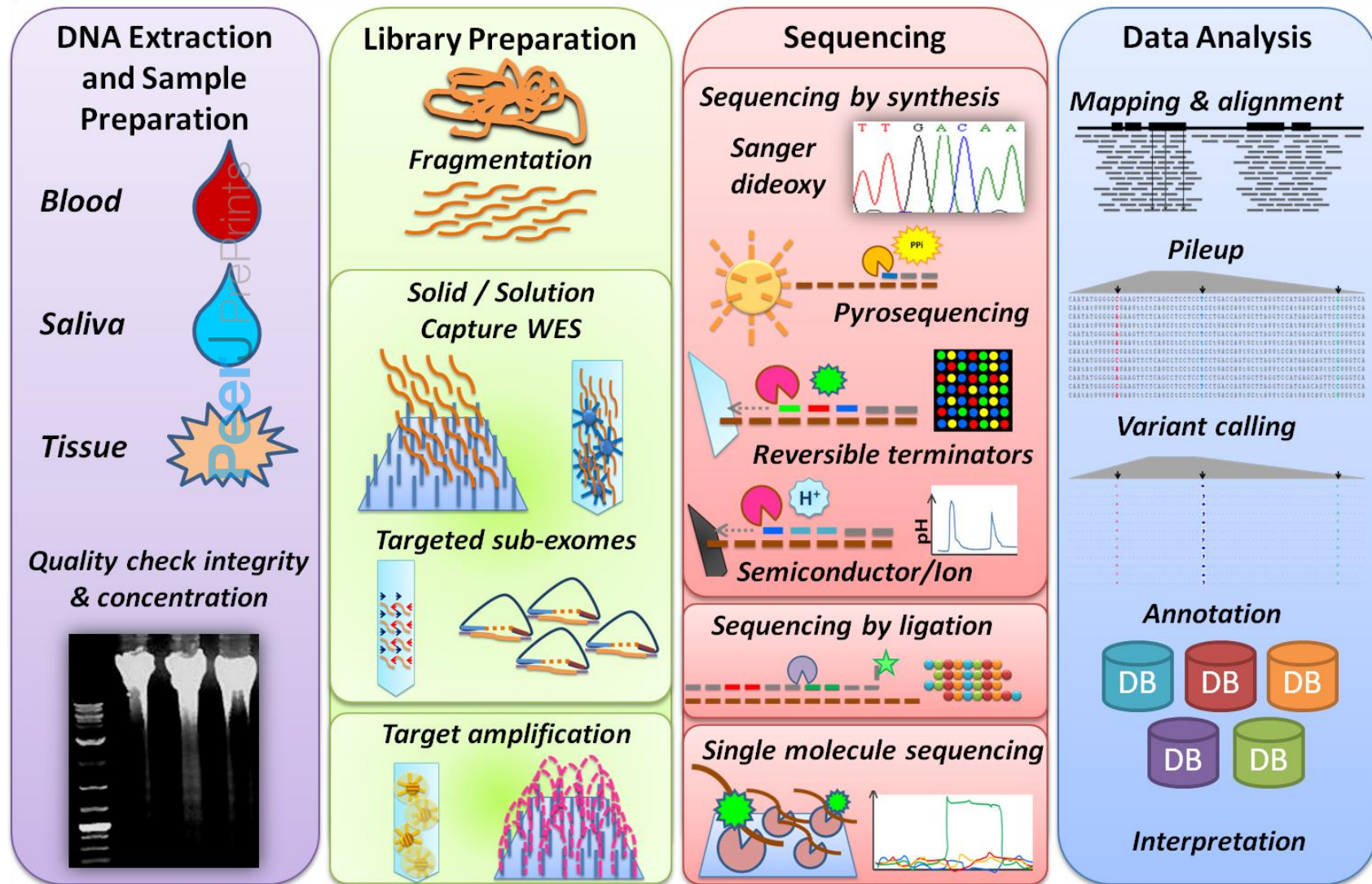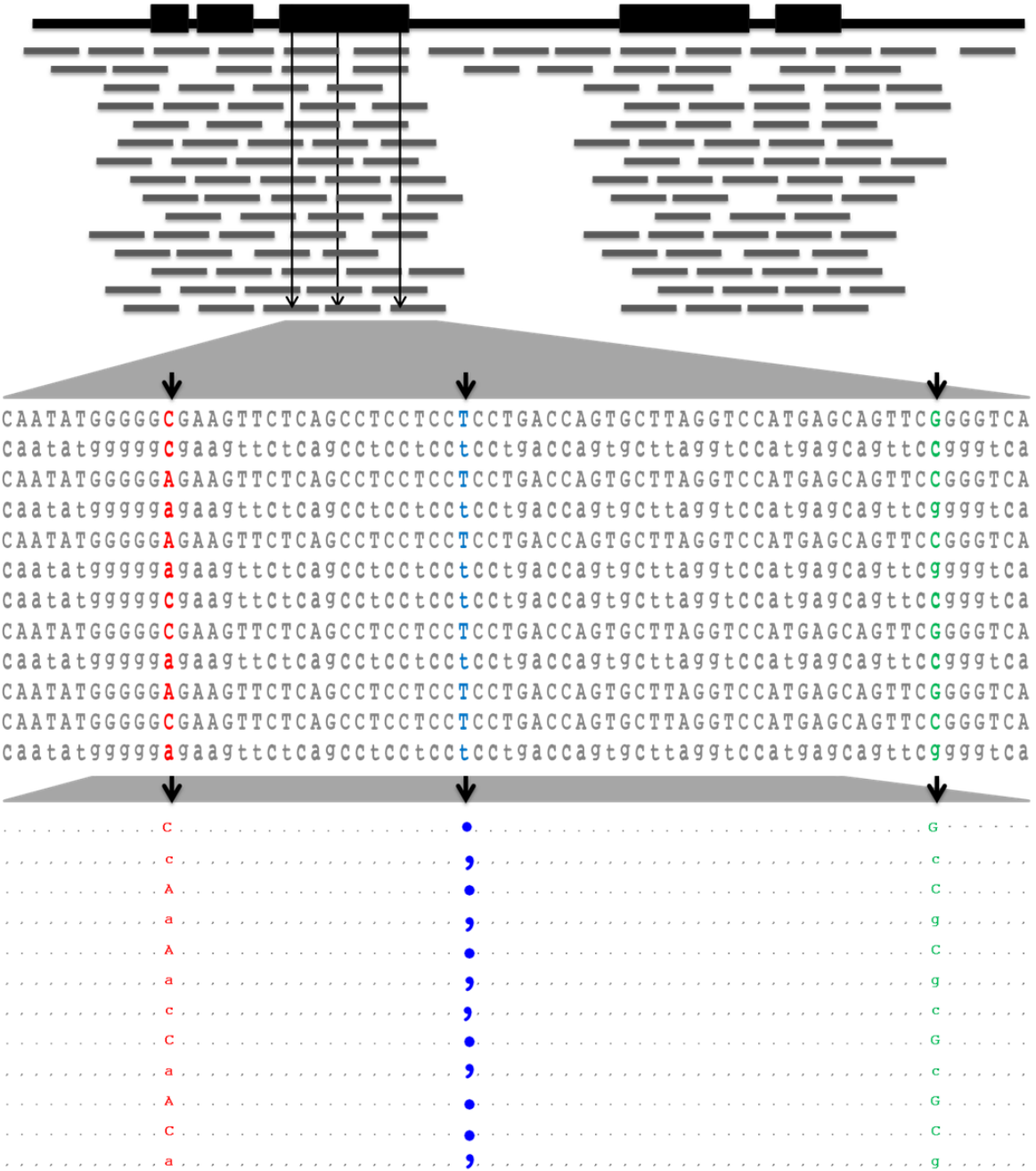
686        2008;123(1-4):307-12

688  **Figure 1**. Overview of the methods used for whole-genome/exome sequencing for human

689  genome variation analyses. 1) DNA extraction from sources such as blood, saliva or tissues is

690  performed. 2) DNA is fragmented and samples are prepared for whole-genome sequencing or

691  target enrichment by a capture method that can be on a solid surface (array) or in solution.

692  Alternatively, for more targeted approaches, PCR or molecular inversion probes (MIPs) can be

693  used for target enrichment. 3) A variety of sequencing technologies are available. These can be

694  subdivided according to the enzyme they use to amplify the target sequences and by the output

695  signal that is detected for sequencing. 4) After sequencing, data generated is mapped and

696  aligned to the human genome reference sequence. A pileup of every base and the nucleotide

697  detected at that position is generated and from the file generated, variants that differ from the

698  reference are extracted and annotated using an extensive variety of databases in order to aid

699  with variant interpretation.

700

35

**Figure 2.**

36

703 **Figure 2.** Schematic representation of the mapping, pileup and variant calling process. Once

704 individual reads are mapped to the reference human genome sequence, a pileup of these reads

705 is generated and every base reported at each aligned position in the genome is reported. In

706 order to facilitate the processing of these data and the files generated, symbols have been

707 assigned to represent a reference base reported in the forward strand (●), a reference base in

708 the reverse strand (ᐧ) and capital and small letters to represent specific variant calls either in the

709 forward or reverse strand respectively..

37