

Chapter 3

SNP Databases

Christopher Phillips

Abstract

Researchers interested in obtaining detailed information on SNPs now work in a golden age of online database availability: never has so much data and such a wealth of information been freely accessible for such a substantial proportion of the 18 million single nucleotide polymorphism (SNP) loci currently characterized in the human genome. This chapter describes the major SNP databases available for human genetics studies. Tools and strategies are outlined that can help researchers properly formulate a database query to be able to access the most appropriate information needed for their research aims, including medical or population genetics analysis – an approach that is getting increased attention given the expanding scale of online SNP data.

Key words: Single nucleotide polymorphism, database, search, query, National Center for Biotechnology Information, dbSNP Entrez, HapMap.

1. Introduction

In silico research as a part of the preparation for an experimental genetics study is now an essential preamble to the choice of genomic regions to analyze and markers to use, the design of genotyping approaches, and the listing of appropriate samples to characterize. This chapter provides a simple guide to the structure and use of the major online SNP databases, adapted to **Sections 2 and 3**, by linking each database to a particular research planning task: finding sets of single nucleotide polymorphisms (SNPs) that share common characteristics (NCBI Entrez); obtaining detailed information on a SNP locus and collating other genetically relevant data (dbSNP); exploring SNPs in coding regions (SNPper and PupaSuite); performing simple scrutiny of linkage

Phillips

disequilibrium (LD) block structure and choosing SNP markers to tag chromosome regions (HapMap); and assessing population genetics parameters from online SNP data (Haplotter and SPSmart).

Some straightforward, common sense advice is given about Internet browsing (*see* **Notes 1** and **2**), processing of SNP data, once obtained, and direct use of generic search engines such as Google – to look across the Web space before focusing on known SNP databases. The latter approach can yield interesting results, but otherwise this chapter assumes the user will go directly to a particular SNP database gateway (*see* **Table 3.1**) to initiate a directed search of online data.

Table 3.1
The major online single nucleotide polymorphism (SNP) databases

Database	Host organization	Gateway URL for initiating SNP data searches
dbSNP	NCBI	http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp
HapMap	The HapMap Consortium	http://www.hapmap.org/cgi-perl/gbrowse/
Ensembl	EMBL-EBI/Sanger Center	http://www.ensembl.org/Homo_sapiens/index.html
Santa Cruz	University of California, Santa Cruz	http://genome.ucsc.edu/cgi-bin/hgGateway
Perlegen	Perlegen Sciences	http://genome.perlegen.com/browser/index_v2.html
Assays-on-Demand	Applera (Applied Biosystems)	https://products.appliedbiosystems.com/ab/en/US/adirect/ab?cmd=ABGTKeywordSearch&catID=600769
SeattleSNPs	US NHLBI (PGA)	http://gvs.gs.washington.edu/GVS/

NCBI National Center for Biotechnology Information, NHLBI, National Heart, Lung, and Blood Institute, PGA Program for Genomic Applications

2. Materials

2.1. The Major SNP Databases

Suggesting a definitive list of the major online SNP databases runs the risk of becoming out of date once done, but *dbSNP*, the SNP database of the National Center for Biotechnology Information (NCBI), and *HapMap* head the list in **Table 3.1**, which is otherwise not intended to indicate an order of size or usefulness. NCBI continues to be by far the most important and comprehensive set of genomic databases available, while the HapMap project has an

97 ever-closer relationship with dbSNP in collating human SNP data.
 98 To summarize a complex and far-reaching project, HapMap was
 99 intended to concentrate global resources on the characterization
 100 of the *variant* part of the genome as a natural extension of the
 101 work of the original Human Genome Mapping Project in estab-
 102 lishing the *invariant* sequence common to everyone and held in
 103 NCBI (1, 2). An important part of the initial work of HapMap was
 104 to check the efficiency of dbSNP, i.e., how well did the dbSNP
 105 catalogue represent the true extent of SNP variability in humans?
 106 This was achieved by resequencing ten ENCODE regions
 107 (detailed in **Table 3.2** of (1) and at [http://www.hapmap.org/
 108 downloads/encode1.html.en](http://www.hapmap.org/downloads/encode1.html.en)) and extrapolating the SNP variabil-
 109 ity found to the genome as a whole. Two findings emerged from
 110 this comparison: firstly the false-negative rate of dbSNP (i.e., how
 111 often SNPs were present but not detected) although very low was
 112 significant for rare SNPs – loci with allele frequencies around 1%
 113 (0.01) or less; secondly the overriding majority of common varia-
 114 tion SNPs had been captured by dbSNP or if absent had proxies in
 115 the same region in tight correlation and listed by dbSNP. It is
 116
 117
 118

Table 3.2
HapMap study populations (1–4): phase I/phase II (5–11): added to phase III. Many
published studies, including those of HapMap (1), merge CHB and JPT to a
“population panel” abbreviated to ASN

	Abbreviation	Samples	Full description	Group
1	YRI	180	Yoruba in Ibadan, Nigeria	African
2	CEU	180	Utah residents with northern and western European ancestry	European
3	CHB	90	Han Chinese in Beijing, China	East Asian (ASN)
4	JPT	90	Japanese in Tokyo, Japan	East Asian (ASN)
5	ASW	90	African ancestry in southwest USA	African
6	CHD	100	Chinese in metropolitan Denver, Colorado, USA	East Asian
7	GIH	100	Gujarati Indians in Houston, Texas, USA	South Asian
8	LWK	100	Luhya in Webuye, Kenya	African
9	MEX	90	Mexican ancestry in Los Angeles, California, USA	Native American
10	MKK	180	Maasai in Kinyawa, Kenya	African
11	TSI	100	Tuscans in Italy	European

AQ2

145 interesting that estimates of false-positive rates in dbSNP (i.e.,
146 incorrectly listing a nucleotide position as a SNP) were not detailed
147 by HapMap, indicating that these were negligible and therefore
148 dbSNP had developed very efficient systems for confirming that
149 SNPs were real (*see Note 3*). In summary, dbSNP has proved to be
150 both a comprehensive and a reliable catalogue of human SNP
151 variability with an efficient system to cross-reference multiple
152 submissions of the same SNPs from centers outside NCBI (*see*
153 **Note 4**). Since 2003, HapMap has been the major contributor
154 of SNP data to dbSNP. The other databases listed in **Table 3.1**
155 both parallel and feed data into dbSNP, so they either provide an
156 alternative system of browsing and searching the core human
157 genome SNP data (Ensemble and Santa Cruz), or list the SNPs
158 generated by their own independent genotyping initiatives with
159 stand-alone browser systems dedicated to the data they have gen-
160 erated (Perlegen, Assays-on-Demand, and Seattle SNPs).

161 1. dbSNP. The strength of NCBI lies in the breadth of genomic
162 databases held under the single umbrella. This means that
163 queries to any of the NCBI databases can tap into the relation-
164 ships that exist between the subject of interest and each of
165 some twenty or more major databases within NCBI. So
166 genetics research involving SNPs is easily set in the context
167 of supporting information that details published studies of
168 the SNP, context sequence of the SNP, gene structure and
169 function (if this is where the SNP is sited), and how the SNP
170 variation is expressed as a phenotype. These data are handled
171 in NCBI by *PubMed*, *GenBank*, *Gene*, and Online Mendelian
172 Inheritance in Man (*OMIM*) databases, respectively (*see*
173 **Note 5**). In addition, NCBI benefits from a unified approach
174 to constructing database queries, so once the user is familiar
175 with the way to query one NCBI database, the same rules will
176 apply to all other queries made. When accessing the most
177 extensive NCBI data, comprising SNP, gene, protein, pub-
178 lications, phenotype, and sequence, one can execute data
179 queries directly from a menu of choices in a global system
180 termed “*Entrez*” (outlined in detail in **Section 3.1**). The SNP
181 Entrez system *EntrezSNP* has a homepage menu listing the
182 principal SNP criteria ([http://www.ncbi.nlm.nih.gov/sites/
183 entrez?db=snp](http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp)) that help to define a search. This is the main
184 starting point of EntrezSNP and this SNP-focused menu
185 differs from those of other Entrez databases such as Entrez-
186 Gene and EntrezProtein. Therefore, dbSNP can be accessed
187 in three ways: by using EntrezSNP; by following hyperlinks
188 embedded in other NCBI databases, and by direct access to
189 SNP summary pages, termed “*Cluster Reports*” – forming the
190 core data page for each locus in dbSNP. A Cluster Report can
191 be thought of as the SNP “homepage” listing a full set of the
192 key parameters in a standardized format.

193 Reference to a SNP within NCBI, within all other SNP
194 databases, and now, almost universally, in the scientific litera-
195 ture is made using a unique identifier: the *rs-number*, consist-
196 ing of a number prefixed with “rs.” As an open database,
197 dbSNP receives submissions from genotyping centers and
198 collates the data into a merged reference set (*see Note 4*).
199 Since different centers routinely report identical SNPs to
200 dbSNP, the submissions are clustered into reference SNPs
201 (termed “*refSNPs*”) based on genome-wide comparisons of
202 the context sequence submitted. For these reasons the dis-
203 tinction between reference SNPs and submitted SNPs is made
204 by rs and ss, respectively: prefixing a unique SNP identifica-
205 tion number with rs, while creating a number for each sub-
206 mission prefixed with ss. All rs-numbers are displayed
207 throughout NCBI as hyperlinks returning the Cluster
208 Report.

209 2. SNP-related databases in NCBI: PubMed, GenBank, Gene,
210 and OMIM.

211 *PubMed* is the NCBI bibliographic database that provides the
212 starting point for researchers to assess the current published
213 state of the art in their chosen area of study. Data comprise ten
214 million published articles from about 5,000 peer-review jour-
215 nals. PubMed is by default predominantly text-oriented so it
216 works by matching text recognized in the query to the text in
217 the data records, including key words in the article body text
218 itself. Therefore, to work efficiently the system needs to care-
219 fully regulate vocabulary, which is done by a separate database
220 of words used to index PubMed known as *MeSH* (i.e., medical
221 subheadings) – searchable itself using the search menu at the
222 top left of each NCBI homepage. It can be an important
223 check to clarify the vocabulary relating to a trait or disease of
224 interest before performing PubMed searches by subject.
225 Searches using rs-numbers can be an efficient way to find
226 studies related to a research aim, but note that the habit
227 persists in many publications of identifying SNPs in genes
228 by the amino acid substitution they create (e.g., shorthand
229 such as MC1R V60L), so these will not be returned from
230 queries (*see Note 6*).

231 *GenBank* is the nucleotide sequence database of NCBI.
232 This simple description belies the scale of the information
233 held – a collection of sequences comprising 60 gigabases of
234 data from more than 130,000 species updated daily. Despite
235 the complexity of GenBank, most users interested in SNP
236 analysis will simply require a specific context sequence seg-
237 ment of about 120–200 nucleotides to design a genotyping
238 assay for the SNPs of interest. As explained later the
239
240

241 application of RepeatMasker and a neighbor SNP scan of
242 ± 100 bases of context sequence makes it more advantageous
243 to collect it directly from the SNP Cluster Report.

244 **Gene** is the gene catalogue of NCBI and like dbSNP pre-
245 sents a single summary page format of relevant information
246 for coding regions, including function summaries, transcrip-
247 tion structure, genome maps, bibliography, protein data,
248 sequences, and related links to supporting data. NCBI uses
249 the near-standard gene identifiers that take the form of the
250 letter/number combination standardized by the Human
251 Genome Organization (HUGO) ([http://www.gene.ucl.ac.
252 uk/nomenclature/](http://www.gene.ucl.ac.uk/nomenclature/)) or throughout NCBI by a GeneID num-
253 ber (*see Note 7*). Working with data that include a large
254 proportion of text-based information can be difficult, so to
255 view the context of a SNP or list of SNPs sited in genes it may
256 be preferable to use a purely graphical display of SNP posi-
257 tions aligned with intronic, exonic, and 5'/3'-flanking region
258 sequences such as that given by SNPper (*see Section 3.3.1*).
259 Taking text-based data even further towards an article format,
260 the OMIM database has summary pages written as articles
261 describing a phenotype, trait, or disorder with a known or
262 suspected genetic basis. As such, both Gene and OMIM are
263 best consulted along with PubMed during the initial stages of
264 a study design to gain an overview of the current understand-
265 ing of a disease process. Luckily, OMIM is highly readable
266 and can be described as an online textbook expanded and
267 updated as knowledge of a trait or condition is consolidated.
268 Searches of OMIM just provide the descriptive text and a list
269 of articles without the benefit of the search items highlighted
270 within the text; publications must then be read to gather the
271 links to the area of interest. Similarly, rs-numbers are not
272 regularly listed in the OMIM article body text.

- 273 3. HapMap. The original stated aim of the HapMap Project – to
274 determine the haplotype structure of the human genome – has
275 expanded to encompass the characterization of all common
276 human sequence variation. The inclusion of copy number
277 variation and the broadening of ENCODE resequencing
278 efforts to capture rare variation will extend this even further,
279 but HapMap remains dominated by common SNPs and their
280 haplotypes: the correlated arrangement of loci in segments
281 defined by highly variable recombination rates. HapMap data
282 have been structured into study phases I–III with different
283 ranges of SNPs, SNP characteristics, and study populations.
284 It is not always easy to find how each phase was defined but, in
285 short, phase I encompassed about one million SNPs in four
286 populations to give one common SNP per 5,000 bases, phase
287 II consolidated SNP coverage with a further 2.5 million
288

289 markers, and phase III has added another seven study popu-
290 lations. Current study population details are outlined in
291 **Table 3.2** and at the time of writing phase III data have
292 become publicly available.

293 The HapMap Web site provides a wide range of data, but
294 of principal interest will be the genome browser, the SNP
295 summary pages, and the HapMap data mart. HapMap has
296 taken the view that the vast majority of users will start with a
297 graphical overview of a chromosome segment and work out-
298 wards from there, so HapMap presents perhaps the best
299 graphical genome browser for SNP variability currently avail-
300 able. Although the default map details (tracks) are relatively
301 sparse, this provides clarity, while numerous other tracks can
302 be added and kept as the user's default arrangement for
303 future browsing. The chromosome coordinates and SNP
304 positions stay as fixed tracks throughout. This representation
305 usefully complements dbSNP since any SNP not character-
306 ized by HapMap has a hyperlink rs-number in position to
307 gain the Cluster Report. SNPs characterized by HapMap are
308 linked to their own summary pages, which are briefer in
309 content, so again linking out to dbSNP can be the best
310 approach here too. The HapMap graphical browser really
311 becomes informative when used to study the haplotype struc-
312 ture around the sites of interest (*see Section 3.4*) – originally
313 mainly coding regions, but increasingly including intergenic
314 regions identified by genome-wide association studies. The
315 methods of graphical representation of haplotype structure
316 can be a challenge to the first-time visitor to HapMap and it is
317 recommended that users familiarize themselves with
318 approaches used by HapMap and in key papers to display
319 LD and that they understand the characteristics of the prin-
320 cipal SNP association metrics of r^2 and D' (3, 4).

321 4. Ensembl, Santa Cruz, Perlegen, and Assays-on-Demand.

322 *Ensembl* and *Santa Cruz* genome databases largely provide
323 alternatives to NCBI to access most of the same SNP and
324 genome data. Ensemble specializes in the analysis of gene-
325 me features and sequence to best identify and annotate
326 genes and has a large range of species under study. This
327 provides the most informative approach for users inter-
328 ested in comparative genomic approaches: where com-
329 monality of nucleotide or protein sequence can be
330 identified by comparing different species. Ensembl has
331 had a pivotal role in the complex task of gene identifica-
332 tion and characterization, pioneering automated gene anno-
333 tation techniques. Hosted in Ensembl, the Vertebrate
334 Genome Annotation (*VEGA*) database provides a range of
335 genome browsers (5). The main aim of VEGA is in providing
336

Phillips

337 the high-quality manual annotation of vertebrate genome
338 sequence. Lastly, Ensembl provides close integration with
339 the high-quality protein sequence database of *Swiss-Prot/Uni-*
340 *Prot* (<http://www.ebi.ac.uk/swissprot/>). This comprises
341 manually annotated protein sequences with content that is
342 fully linked with the Ensembl gene annotation pipeline.
343 *Santa Cruz* has several features that can provide easier ways
344 than NCBI to obtain information for SNP analysis – for
345 example, the simple process of collecting extended context
346 sequence for a SNP is more straightforward in Santa Cruz
347 than from within dbSNP (*see Note 7*). Therefore, on occa-
348 sions, working with two Web pages with different genome
349 data browsers (essentially accessing the same underlying infor-
350 mation) can be the optimum approach. The guide to Santa
351 Cruz queries is at *Perlegen* and Applied Biosystems’s *Assays-*
352 *on-Demand* are private databases of SNP variability infor-
353 mation that has been submitted to dbSNP and is publicly avail-
354 able, but can also be accessed from each company’s Web site
355 with dedicated filtered search pages. Filters parallel the query
356 process of Entrez by offering a choice of criteria that reduce
357 the data set returned to a small, more manageable group of
358 items meeting the criteria. Both databases elected to study
359 US European, US African-American, and US Chinese popu-
360 lation panels that to a large extent mirror those of HapMap’s
361 CEU, YRI, and CHB, so data obtained can be combined
362 from different sources to allow meaningful comparisons of
363 population variability or less often directly between different
364 samples but originating from the same population group
365 (although comparing YRI Africans with African-Americans
366 highlights the about 20–30% European admixture in the
367 latter). The easiest way to compare SNP data from similar
368 populations in different databases is to use SPSmart (*see Sec-*
369 *tion 3.5.2*). Note that Perlegen uses an internal SNP identi-
370 fier with the format “PS+8 digit no.” (e.g., PS04631975)
371 but accepts rs-number queries, while SPSmart provides a list
372 of these numbers in its returned data.

373 *Assays-on-Demand* SNP data are in large part based on the
374 Celera SNP database generated during the private genome
375 annotation performed by Celera after the human sequence
376 had been completed in parallel to the completion of the
377 public Human Genome Mapping Project in 2000. Celera
378 genome data were available on a subscription basis (as Celera
379 Discovery System, or *CDS*) between 2002 and 2006, but now
380 all Celera’s SNP data have been incorporated into dbSNP and
381 can be individually filtered in a search in Entrez with the
382 inclusive term “AND Celera” or the exclusive term “NOT
383
384

385 Celera” (*see Section 3.3.1*). Accessing Celera SNP data is also
386 possible through Assays-on-Demand; users in the latter case
387 can utilize a stand-alone tool called “SNPbrowser” compris-
388 ing five million SNPs from public and CDS sources. This
389 allows access to some of the original CDS SNP and gene
390 annotation but is of most use as an alternative to HapMap
391 for the definition of haplotype structure in a particular chro-
392 mosome segment (*see Section 2.2.2*). Particularly in the
393 population genetics field, Assays-on-Demand allows a simple
394 system to review a large data set of SNP allele frequency
395 variability from three major population groups and so has
396 provided a core search step for many studies seeking to isolate
397 and develop ancestry informative marker SNPs (5).

- 398 5. SeattleSNPs. The SeattleSNPs initiative is funded as part of
399 the US National Heart, Lung, and Blood Institute (NHLBI)
400 Program for Genomic Applications (PGA) – the latter abbrevi-
401 ation is used by dbSNP to reference SeattleSNPs SNP
402 submissions. The project has undertaken the resequencing
403 of more than 300 genes identified as primarily important
404 in the inflammatory response, but also including cardiovas-
405 cular disease and the immunity (a full list of completed genes
406 is at http://pga.gs.washington.edu/finished_genes.html).
407 Although it is important to stress that the gene list mentioned
408 above is not prescriptive – users are encouraged to nominate
409 candidates for consideration. Therefore, SeattleSNPs provid-
410 es a key opportunity to capture and characterize low-fre-
411 quency SNPs from whole sequence data that would otherwise
412 escape detection or be subject to acquisition bias (*see Note 9*).
413 As sequencing technology has recently undergone one of the
414 periodic quantum leaps in throughput, the chance to properly
415 discover and catalogue new low-frequency SNPs by resequen-
416 cing sufficiently large sample groups or individuals with a
417 particular disorder will form the next major phase of SNP
418 databasing. The extended ENCODE studies and the Seat-
419 tleSNPs initiative stand at the vanguard of this work, with the
420 1,000 Genomes Project poised at the time of publication to
421 take resequencing to the next level of resolution: that of full
422 individual genomes. The evident drawback of SeattleSNPs
423 comes from a focus on targeting a subset of genes or the
424 pathways they occupy with the bias this might represent in
425 attempting to understand the disease process. This is mainly
426 due to the need to direct resources to the best areas for
427 detailed SNP genotyping, and the fact that SeattleSNPs is
428 actively engaged in association studies allows it to combine
429 the knowledge this generates with new targets for resequen-
430 cing in the genome. The SNP data from resequencing is fed
431
432

Phillips

433 into a database known as the *Genome Variation Server* (GVS)
 434 and users are encouraged to access the tutorials that explain
 435 optimum use of SeattleSNPs and GVS at [http://www.open](http://www.openhelix.com/downloads/seattlesnps/seattlesnps_home.shtml)
 436 [helix.com/downloads/seattlesnps/seattlesnps_home.shtml](http://www.openhelix.com/downloads/seattlesnps/seattlesnps_home.shtml).

437
 438 **2.2. A Selection of**
 439 **Tools To Aid Analysis**
 440 **of SNP Data**

The following tools are available to use as Web-based search systems or stand-alone programs that can help to make directed searches of the databases outlined previously.

1. NCBI tools: dbSNP-announce, MyNCBI, MapViewer, and Genome Workbench.

441 Although not strictly online tools, *dbSNP-announce* (http://www.ncbi.nlm.nih.gov/About/news/announce_submit.html) and *MyNCBI* (<http://www.ncbi.nlm.nih.gov/entrez/login.fcgi>) are important subscription-based adjuncts to any use of dbSNP. Subscribing to dbSNP-announce provides automatic reports to the user's e-mail address of dbSNP updates. As well as reporting the release of each new build, announcing newly added features, and outlining corrections or discovered problems with past or present builds, there is an archive for referencing possible problems with, or qualifications to, previously obtained search data. MyNCBI, requires a single subscription step to provide a search workspace for the user that provides a clipboard permitting combined searches from stored results obtained at different times (*see Section 3.1.7*).

442 *MapViewer* integrates the bulk of the NCBI databases into a customizable genome map of aligned components termed "map elements." The SNP data map element, termed "*Variation*," can be included with any other genome feature in a custom map. A simple, clean icon set against each SNP marker positioned on the map showing a chromosome segment provides a clear summary description of the locus. Map browsing offers an intuitive way to review large numbers of SNPs in one session. Exploring a chromosome segment as a map is the best way to scrutinize the position and characteristics of nearby genome features of importance such as neighbor SNPs, genes, and their transcripts. Furthermore, the features around each SNP can be scrutinized easily through a series of hyperlinks embedded in many of the key map elements such as Genes.

443 NCBI *Genome Workbench* (<http://www.ncbi.nlm.nih.gov/projects/gbench/>) is a stand-alone program that works locally, i.e., independently of individual online access to NCBI. Once installed, it can access and display genomic data from NCBI and combine this with the user's own data in a series of graphical representations. The program is available for download and installation in any operating system format,

AQ5

AQ6

481 and offers considerable flexibility in how the user chooses
482 to align and compare genomic data. This extends to a
483 range of alignment views, phylogenetic tree views, and
484 tabular views of data. It can also align user's data to
485 those of public databases, and retrieve BLAST results. A
486 full guide is beyond the scope of this chapter, so users are
487 encouraged to explore this tool and the five tutorials
488 ([http://www.ncbi.nlm.nih.gov/projects/gbench/tutorial.](http://www.ncbi.nlm.nih.gov/projects/gbench/tutorial.html)
489 [html](http://www.ncbi.nlm.nih.gov/projects/gbench/tutorial.html)) for themselves.

490 2. Checking SNP assay primer designs: BLAST and Santa Cruz
491 In Silico PCR.

492 *BLAST* is a tool for assessing/calculating sequence similarity
493 between a query sequence and the target sequence(s) avail-
494 able in the NCBI GenBank nucleotide databases. Users inter-
495 ested in developing SNP assay designs will query Nucleotide
496 BLAST in two ways: (1) finding the location of a submitted
497 sequence that includes the SNP, as the query “does the sub-
498 mitted sequence exist in a GenBank database?”, and (2)
499 checking for coincidental similarity in a sequence, normally
500 a PCR primer, the query being “what is the degree of speci-
501 ficity of the submitted sequence?” These sequence compar-
502 isons can be made by choosing the standard BLAST (termed
503 “*blastn*”) and *Search for short and near exact matches* options,
504 respectively. As a simple and quick alternative to BLAST, the
505 Santa Cruz *In Silico PCR* tool ([http://genome.ucsc.edu/](http://genome.ucsc.edu/cgi-bin/hgPcr?command=start)
506 [cgi-bin/hgPcr?command=start](http://genome.ucsc.edu/cgi-bin/hgPcr?command=start)) offers a straightforward sys-
507 tem that indicates the expected PCR product sequence from
508 primer designs submitted by the user from comparisons to the
509 current human reference nucleotide sequence. This tool is
510 highly recommended since it provides a simple check before
511 committing to primer purchases.

512 3. Exploring haplotype block structure maps: Haploview and
513 SNPBrowserTM.

514 *Haploview* (<http://www.broad.mit.edu/mpg/haploview/>) is
515 an essential adjunct to HapMap browsing comprising a Java
516 applet tool that permits the analysis and visualization
517 of haplotype block patterns in HapMap data, choosing
518 tagSNPs (7, 8), and estimating haplotype frequencies (*see*
519 **Section 3.4**).

520 The Applied Biosystems *SNPBrowser*TM tool ([http://](http://marketing.appliedbiosystems.com/mk/get/snpb_landing)
521 marketing.appliedbiosystems.com/mk/get/snpb_landing)
522 provides a stand-alone database of five million Celera SNPs
523 that is downloaded to the user's PC and can therefore be
524 accessed offline or in parallel to online searches. The SNP
525 data are presented as a chromosome segment map showing
526 haplotype block distributions defined by Celera's own pair-
527 wise analysis of 160,000 SNPs (termed “*backbone validated*”
528

SNPs”), so it provides an alternative to HapMap in the annotation of human haplotype blocks, although it can display both HapMap and Celera haplotype maps. Additionally, SNPBrowserTM is easily configured to tailor haplotype block annotation displayed, SNP type, population studied, and size of the region shown. SNPBrowserTM works along the same lines as Assays-on-Demand in providing a shopping list of SNPs based on user’s criteria that can then be ordered as commercial singleplex (Applied Biosystems TaqManTM) or multiplex (Applied Biosystems SNPlexTM) SNP genotyping assays.

4. Mapping SNPs and mutations in genes: SNPper.

SNPper provides a tool for the extraction and re-presentation of SNP data from public databases focused on coding regions, offering the clearest system for scrutinizing SNP positions in and around genes (9). Once the user has provided the gene identifier, SNPper will list exonic, intronic, and 5’/3’-regions, plus embedded SNP positions within these, either as a plain nucleotide sequence or as triplet code groups with their amino acid codes. Although the same output can be achieved with GenBank and Santa Cruz nucleotide browsers, SNPper is a much quicker and simpler system for listing SNPs in a gene of interest with a clean and intuitive graphical summary of the gene. This particularly suits the cataloguing of mutation sites in coding regions since these are usually defined by the amino acid changes they produce and SNPper allows their identification in relation to the SNP landscape that surrounds them, providing a straightforward way to develop genotyping assays.

5. Exploring the effect of SNPs on gene action: PupaSuite, Polyphen, and ESEfinder.

PupaSuite (“Pupa” stands for putative phenotype alterations) encompasses two tools – PupaSNP and SNPeffect – that aid the identification of SNPs effecting the processing of genes (10, 11), namely, sites of intron/exon boundaries or exonic splicing enhancers (ESEs), predicted transcription factor binding sites, and amino acid sequence changes. PupaSuite works with the Ensembl gene annotation and SNP database and can process an uploaded SNP list, but the user can also provide individually identified SNP sites with the aim of exploring their effect on gene action. The utility of PupaSuite is the ability to explore the effect of SNPs on transcriptional activity and splicing as well as protein sequence – an increasingly important step when analyzing coding regions.

PolyPhen (<http://genetics.bwh.harvard.edu/pph/data/index.html>) is a tool that usefully predicts the possible impact of an amino acid sequence change on the properties of a

577 protein (12). Although it will not accept nucleotide input
578 directly as it holds a nonsynonymous SNP database compris-
579 ing about 50,000 SNPs from dbSNP, PolyPhen can check
580 whether a SNP is nonsynonymous or not, using the site tool
581 SNP2Prot. Effects on proteins are tentatively defined as
582 unknown, benign, possibly damaging, and probably dama-
583 ging. Users can input rs-numbers directly for comparison
584 against the PolyPhen data, but they are advised to go directly
585 to PupaSuite for novel coding SNPs discovered in their study.
586 As with SNPper, this tool is particularly applicable to the
587 characterization of mutations that are, by definition, SNPs
588 at very low frequency.

589 Of the three tools that help define SNP effects, the most
590 specialized is *ESEfinder* (<http://exon.cshl.edu/ESE/>), a
591 tool dedicated to identifying precursor RNA splice site
592 changes from SNPs sited at *exonic splicing enhancers* (ESEs)
593 (13). As such, SNPs at the ESE positions of proteins that
594 routinely undergo alternative splicing can profoundly affect
595 the final protein structure. ESEfinder makes use of databases
596 of different ESE sequence motifs to help identify putative
597 SNPs influencing splice patterns.

598 6. Using SNP haplotypes to detect signatures of selection: Hap-
599 lotter and SWEEP™.

600 Compared with the tools available for studying gene and
601 genome structure described above, population genetics
602 tools are latecomers to SNP database analysis. Data of gen-
603 ome-wide patterns of polymorphic marker variation provide a
604 powerful tool for studying the history of migration, bottle-
605 necks/expansions, and adaptation in human populations. For
606 those interested in analyzing such events, a major advantage
607 in using SNP data is the distribution of SNPs at much higher
608 densities compared with microsatellite or insertion–deletion
609 variation and in the advanced characterization of SNP-based
610 haplotypes. Therefore, SNPs are obvious candidate markers
611 for the analysis of patterns of haplotype structure that can
612 indicate signatures of past natural selection. Positive selection
613 will amplify the frequency of a particular haplotype surround-
614 ing a favorable, novel gene variant because the haplotype
615 accompanying the variant on the same chromosome strand
616 also rises rapidly in frequency throughout the population.
617 Before recombination disrupts this association, much higher
618 SNP homozygosity is seen, as identical haplotypes are more
619 likely to be found on each chromosome. Therefore, homo-
620 zygosity is raised in the immediate vicinity of the selected gene
621 variant and diminishes with distance, as recombination
622 increasingly breaks up associations. This is the basis of the
623 extended haplotype homozygosity (EHH) test that aims to
624

625 detect signatures of recent selection by analyzing irregularly
626 long haplotype homozygosity patterns (14). Two tools are
627 available for EHH analysis: *Haplotter* uses HapMap data and
628 is accessed online, while *SWEEP* is a stand-alone program that
629 can use data from any source that has been phased (i.e., allele
630 combinations assigned to one of two strands).

631 Haplotter (<http://hg-wen.uchicago.edu/selection/haplotter.htm>) measures a value iHS (15) that expresses the
632 contrast between haplotypes with changed frequencies and
633 the surrounding genome landscape, so it can reveal frequency
634 rises in ancestral alleles (positive contrasts as the allele
635 increases in frequency) as well as in variant alleles (negative
636 contrasts). Haplotter can work from gene identifiers or a
637 single SNP landmark (slower and varied in coverage). The
638 program returns plots of iHS , plus standard selection signature
639 or population diversity measures H , D , and F_{st} , followed
640 by a table of adjacent genes, colored light blue when showing
641 significant evidence of selection. The major advantage of
642 Haplotter is it allows an unbiased approach to finding regions
643 with indications of recent selection, so in use it is likely to
644 reveal interesting new candidates for more detailed study.
645 This can enable studies to focus on the phenotypes such loci
646 exhibit as a way to explore differences in susceptibility to
647 disease between populations. An advantage of using HapMap
648 data is that the study populations will be extended to allow
649 examination of more widely distributed patterns of local
650 selection.

651 The stand-alone program *SWEEP*TM (<http://www.broad.mit.edu/mpg/sweep/>) acts like Haplotter to measure
652 the rate of decay of homozygosity with distance from putative
653 regions subject to selection (14). Although it requires time
654 and care to become familiar with use of the program, the
655 graphical output, particularly diagrams termed “*bifurcation
656 plots*,” provides very good representations of results summar-
657 izing extended homozygosity versus genomic distance to the
658 core haplotypes.

659 *SPSmart* (<http://spsmart.cesga.es/>) is a tool that performs
660 the simple task of re-presenting SNP allele frequency data
661 from multiple sources as pie charts identical to those of the
662 HapMap browser. So *SPSmart* allows the user to review SNP
663 variability across a wider range of populations than is feasible
664 from single databases accessed one by one. This appears to
665 offer little extra value if, for example, the study populations of
666 HapMap phase II and Perlegen are considered, with only a
667 comparison of YRI Africans and African-Americans of poten-
668 tial interest. However, *SPSmart* also processes data from the
669 Stanford and Michigan University initiatives that have
670
671
672

673 genotyped some 650,000 SNPs in the CEPH human genome
 674 diversity panel (HGDP) comprising over 1,000 samples from
 675 51 global populations. Incorporation of HapMap phase III
 676 populations has also boosted the scope of global variability
 677 that can be accessed with SPSmart.
 678
 679
 680

681 3. Methods

682 3.1. Finding Sets of 683 SNPs That Share 684 Particular 685 Characteristics: NCBI 686 Entrez and Boolean 687 Rules of Database 688 Searching

- 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719
 720
1. The NCBI Entrez system uses *Boolean* terms or *operators* to define searches. These include the principal operators: *AND*, *OR*, and *NOT*, summarized in **Fig. 3.1**. Operators are the key parameters that define the relationship between criteria that describe database entries. In Entrez these descriptive details or criteria are put in groups termed “*fields*” that are defined by *tags* (alternatively qualifiers). Field details can be written in lowercase letters (but following an appropriate format, or *syntax*) ahead of their tags, which are always given in capital letters with fixed syntax within square brackets, for example, to define search criteria “SNPs on chromosome 22” the field would be 22 denoted by the tag [CHR] written as 22[CHR], the chromosome field syntax being a number or X or Y. Users can either manually construct their own search with any combination of fields/tags and operators or simply choose tags from a menu on the EntrezSNP homepage (<http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp>) and provide fields to make a query using a default AND operator. For users unfamiliar with searches in NCBI, the latter option of choice from a menu can be easier to start with. The principal fields and their tags provided in EntrezSNP are given in **Table 3.3**. Fields separated by spaces alone also default to AND, e.g., query “HERC2[GENE] coding non-synon[FUNC]” finds the nonsynonymous SNPs in HERC2 exons.

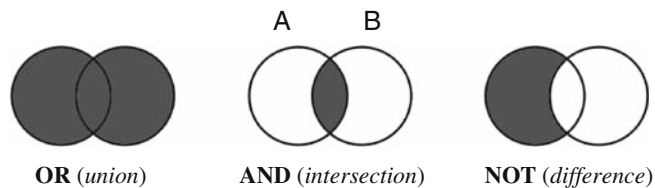


Fig. 3.1. Boolean operators. *OR* applies to all items in A or B, *AND* applies to items found in both A and B, and *NOT* applies to items in A not found in B.

Phillips

Table 3.3
Key EntrezSNP fields and their tags

Description	Tag	Search field used	Example query
Observed alleles	[ALLELE]	IUPAC allele code (see Table 3.4)	R[ALLELE] find SNPs with A/G substitutions
Chromosome	[CHR]	Number/X, Y	21[CHR] OR 22[CHR] find SNPs on chromosomes 21 & 22
Base position	[BPOS]	Ranged number & AND & [CHR]	18000:28000[BPOS] AND Y[CHR] find SNPs in 10 kb segment of Y chromosome
Heterozygosity	[HET]	Ranged number	30:50[HET] find SNPs with heterozygosity value in range 30–50%
Function Class	[FUNC]	Locus region, intron, etc. (8 in total)	Coding nonsynon[FUNC]
Build	[CBID]	Number	125[CBID] search build 125
Gene location	[GENE]	Gene symbol	DARC[GENE] search for SNPs in Duffy blood group, chemokine receptor
Genotyping method	[METHOD]	Description as listed at URL below	Hybridize[METHOD] search for SNPs found by chip hybridization
http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=snp – METHOD)			
Map weight	[MPWT]	Number: 1=once, 2=twice, 3=3–9 times	NOT (2[HIT] OR 3[HIT]) exclude SNPs mapping twice or more in genome

SNP single nucleotide polymorphism.

- As a worked example, a search could be written longhand: “find unique SNPs in dbSNP on human chromosome 22 that are AC substitutions and show heterozygosity of more than 45%.” To perform a manual EntrezSNP search, place the following search description in the search box: 1[MPWT] AND human[ORGN] AND 22[CHR] AND M[ALLELE] AND 45:50[HET]. Note the field/tag items follow the same order as the longhand query, but this is not essential. The *IUPAC allele codes* applicable to the [ALLELE] field/tag are listed in **Table 3.4**. To perform the same search using the Entrez menu system, go to the limits menu – <http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp&TabCmd=Limits> – and choose tick boxes in (respectively) *Map Weight*, *Organism*, *Chromosome*, *Variation Allele*, and *Heterozygosity*.
- Note that the heterozygosity tag in the above example search uses *ranging*: a range of values to define the field, described by a colon (:) in the middle of range limits. The menu-based

769 **Table 3.4**

770 **IUPAC SNP substitution codes used in Entrez as the field with the [ALLELE] tag. In**
 771 **addition A, C, G, and T can also be used with [ALLELE] to select all SNPs showing**
 772 **that base as an allele**
 773

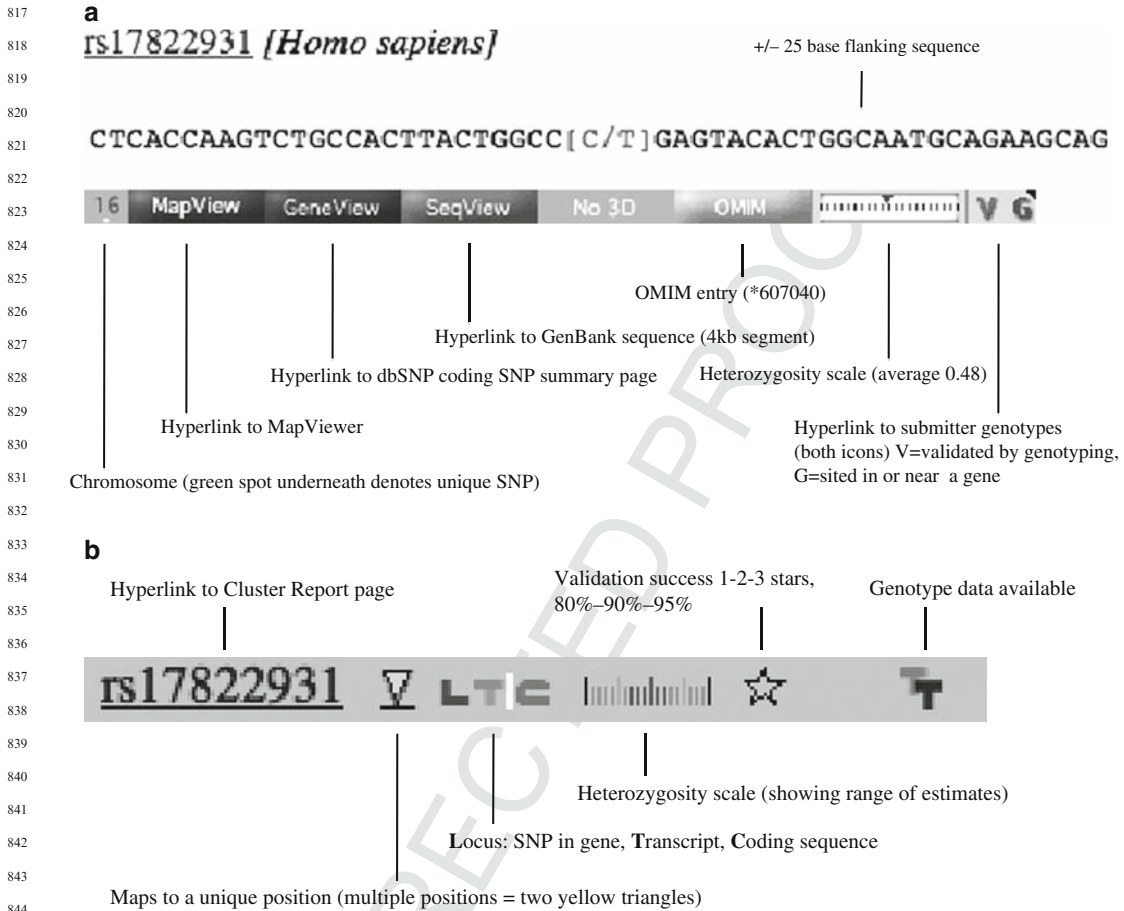
774 Code	Substitution	Code	Substitution
776 M	A or C	V	A or C or G
777 R	A or G	H	A or C or T
778 W	A or T	D	A or G or T
779 S	C or G	B	C or G or T
780 Y	C or T	N	A or C or G or T
781 K	G or T		

782 *N* can also denote an indeterminate base.
 783
 784

785
 786
 787
 788 system only allows heterozygosity ranges of 10%, so fine-
 789 tuned searches such as 49–50% heterozygosity require man-
 790 ual construction. Two other modifiers of operator function
 791 exist for manual searches: *parentheses* and the *wild-card aster-*
 792 *isk* (*). Parentheses group search terms into logical sets to
 793 obtain items that further operations can search. To a large
 794 extent, the logic follows that used in a normal sentence, for
 795 example, in a PubMed search “find articles on the effects of
 796 heat and humidity on multiple sclerosis” is (heat OR humid-
 797 ity) AND multiple sclerosis, while “find articles on the effects
 798 of heat as well as the effects of humidity on multiple sclerosis”
 799 is heat OR (humidity AND multiple sclerosis). The wild-card
 800 asterisk in place of missing text allows a partial entry to be used
 801 as a query term, e.g., using BRC*[GENE] will find both
 802 BRCA1 and BRCA2 genes.
 803

- 804 4. Each SNP in the EntrezSNP list that is returned from a query
 805 defaults to a summary graphic with components that describe
 806 the key parameters of the SNP, outlined in **Fig. 3.2a**. For the
 807 above example, query EntrezSNP lists 911 SNPs in order of
 808 chromosome position. If multiple chromosomes are listed,
 809 these are in order: Y, X, 22, 21, etc. A useful feature is the
 810 ability to change the default listing order amongst six options,
 811 including SNP ID (ascending rs-number) and heterozygos-
 812 ity. When assessing the role of particular SNPs in a disease
 813 process or by association with a candidate region, a particu-
 814 larly useful feature is the “Cited in PubMed” tab. Click the
 815 “Links/Pubmed (SNP Cited)” hyperlinks in this list to obtain
 816 each publication abstract.

Phillips



845 Fig. 3.2. Key to summary graphics for single nucleotide polymorphisms (SNPs). (a) Key to summary graphics for
 846 EntrezSNP search return of example SNP rs17822931. (b) Key to summary graphics for SNPs shown in the chromosome
 847 view in NCBI MapViewer of example SNP rs17822931.

- 848
 849
 850 5. It is important to note that not all possible search fields can be
 851 accessed from the EntrezSNP limits menu: some 14 from a
 852 total of 24 are available (the full list and details are at [http://](http://www.ncbi.nlm.nih.gov/sites/entrez?db=snp)
 853 www.ncbi.nlm.nih.gov/sites/entrez?db=snp). By far the
 854 most useful search field for medical genetics studies omitted
 855 from the limits menu is Gene. This allows the listing of SNPs
 856 found within and close to a gene (or multiple genes using the
 857 OR operator) described by the query. SNPs at the 5'-end and
 858 the 3'-end (that order) of the gene are also listed, but it is
 859 prudent to extend the search using chromosome/base posi-
 860 tion tags – [CHR] AND [BPOS] – to capture potential
 861 promoter SNPs further afield.
 862 6. An Entrez search can be initiated from the NCBI Entrez
 863 homepage by selecting “all databases” from the Search
 864 drop-down menu or specifically from Entrez SNP (“SNP”

865 from the same menu). While searches from the EntrezSNP
 866 homepage return just a list of SNPs meeting the criteria, using
 867 the *all databases* option creates an NCBI-wide search with
 868 hyperlinked numbers of entries from 35 different databases
 869 that can be explored individually: making a good starting
 870 point in the early stages of a genetic study. The cross-database
 871 returns page groups six text-based databases, 26 genomic
 872 databases, and three catalogues (books, journals, and MeSH
 873 vocabulary) into three separate boxes. A cross-database search
 874 can be made using specific or general terms to obtain, respec-
 875 tively, a focused query of the broadest possible coverage
 876 within NCBI or a more open ended survey. For instance,
 877 “rs2293855” as a query returns a single PubMed reference
 878 to a possible role of this SNP in obesity, however with no
 879 reference to the gene MTMR9, where it resides, while “obe-
 880 sity” as a query lists no specific SNPs, but more than 111,000
 881 publications and 681 genes, including MTMR9.

882 7. EntrezSNP gives the most efficient system for progressive
 883 searches as the lists generated can be stored in a clipboard
 884 and then sent to MyNCBI (avoiding an 8 h inactivity delete
 885 step for the clipboard), exported as a text file, combined with
 886 new searches, or re-searched itself. This uses the clipboard and
 887 history tabs at the top of the EntrezSNP page. The clipboard
 888 is a workspace for holding up to 500 items, while history lists
 889 the database search activity as numbers prefixed by a hash (#).
 890 Making Entrez searches at different times by exporting to
 891 MyNCBI allows the user to monitor the number of returns
 892 obtained with different search term combinations. Previous
 893 searches can be combined as hash fields with operators (e.g.,
 894 #1 AND #2 gives items common to both searches). It is also
 895 possible to use hash fields together with normal fields, helping
 896 to build a stepwise record of the search process as it is mod-
 897 ified in incremental stages.

898 8. To automatically reduce SNP numbers returned by a search,
 899 certain fields are best used as filters with fixed values including
 900 the organism and map weight. Therefore use of the huma-
 901 n[ORGN] field/tag ensures only human SNPs are listed and
 902 1[MPWT] ensures all SNPs are unique (i.e., single map
 903 weight). The SNP validation tag also provides a system to
 904 filter out SNPs detected by sequence comparisons alone,
 905 using by frequency[VALIDATION].
 906 .

907
 908 **3.2. Obtaining Detailed**
 909 **Information on a SNP:**
 910 **dbSNP Cluster Reports**
 911

912 1. The Cluster Report page of dbSNP provides most, if not all,
 the information needed to assess the characteristics of a SNP
 and design a genotyping assay if required. Each page is broken
 down into seven sections with largely self-explanatory

913 headings: Submitter records; Fasta Sequence; GeneView;
914 Integrated Maps; NCBI Resource Links; Population Diver-
915 sity; and Validation Summary. This section of the chapter
916 outlines the steps required to (1) obtain sufficient context
917 sequence to design a genotyping assay, (2) scrutinize the map
918 position of the SNP with accompanying genomic features,
919 and (3) begin analysis of the population variation shown by
920 the SNP.

- 921 2. The Fasta section is named after the fast-all sequence similar-
922 ity program used by dbSNP to detect identical SNP submis-
923 sions given in the Submitter section (*see Note 4*). Fasta lists
924 the flanking sequence surrounding the SNP position – the
925 quantity of nucleotides listed is variable in extent but always
926 arranged as groups of ten bases, six groups per line, with the
927 SNP positioned on a separate line as an IUPAC code base.
928 The single Fasta header line contains summary locus details
929 explained in the “Legend” hyperlink in the title above. Reli-
930 ance on dbSNP Fasta sequence alone for primer designs can
931 cause problems (*see Note 10*); however, one clear advantage
932 of dbSNP Fasta is the inclusion of information from the initial
933 submitter and from RepeatMasker analysis (*see Note 11*)
934 which predicts the likely genomic uniqueness of the SNP
935 context sequence. This can help avoid certain sequence seg-
936 ments that may occur in multiple genomic locations and
937 therefore reduce the specificity of any genotyping assay
938 designs.
- 939 3. An essential additional aid to the primer design process is the
940 neighbor SNP detection tool found in the Integrated Maps
941 section. Clicking on the “View” hyperlink under the Neigh-
942 bor SNP heading of each assembly (for *ref_assembly*: i.e., the
943 reference sequence is recommended) gives all SNPs within \pm
944 100 bp of the reported SNP, including at 0 bp, the target
945 substitution itself. This permits easy location and masking of
946 variable sites that can interfere with the predictable binding of
947 primers designed for the assay.
- 948 4. GeneView gives a graphical overview of the SNP if it is located
949 in a gene, giving a position mark using nine color codes for
950 one of 16 predicted functions (hyperlink guide: *color legend*).
951 The cSNP radio button directs one to the coding SNP listing
952 page for each gene (*see Note 6*).
- 953 5. The Integrated Maps section provides a “snapshot” genome
954 view of the SNP position as a red mark-point by clicking the
955 chromosome number hyperlink. This whole genome view can
956 be used to map a series of SNPs by using the OR operator
957 between each rs-number in the query string in the search box.
958 Including “NOT Celera” eliminates the double mark-points
959 and position listing (*see Note 7*). From this overall map any
960

961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008

AQ9

chromosome can be viewed in detail in NCBI MapViewer, centered on the SNP position by clicking the number hyperlink of each chromosome showing a mark-point. Once in MapViewer, configure the view by clicking “maps and options” to show the “Variation” map as the master to ensure a summary icon set accompanies each SNP position (outlined in **Fig. 3.2b**). Any number of other tracks (map components) can be selected as part of the genome landscape around the SNP, although the single most useful of these is Gene.

6. The Population Diversity section summarizes SNP allele frequency distributions in different populations using gold (reference allele) and blue (alternative allele) bars. Clicking on “Genotype Detail” provides the Genotype and Allele Frequency Report (still in beta status at the time of publication). It is possible to obtain the individual SNP genotypes submitted from each contributing laboratory – data that can be particularly useful when using standard DNA samples, such as Coriell panels (<http://ccr.coriell.org/>), as genotyping controls in an assay. The easiest way to achieve this directly from HapMap, the major source of SNP population variability data in dbSNP, is to go to the individual SNP information page in HapMap (http://www.hapmap.org/cgi-perl/snp_details?name=rsnumber) and click the “retrieve genotypes” hyperlinks. Genotypes are always listed in the same sample-ID order and so can be downloaded directly to Excel and correctly ordered as rows per population per SNP using the text to columns option (*see Note 1*).

3.3. Exploring SNPs in Coding Regions: SNPper and PupaSuite

1. SNPper requires a subscription before the user can explore a gene of interest which can be found using “Gene Finder” by inputting standard identifiers (*see Note 8*). Once the gene has been obtained, click on the “Annotated” sequence hyperlink to obtain the nucleotide listing marked as follows: green, 5’/3’; black lowercase, exonic noncoding; black uppercase, coding; gray, intronic; blue underlined, SNPs. A useful approach for the reliable detection of mutations or scrutiny of coding SNPs is to click “View amino acid sequence” to obtain the coding nucleotides as triplet codes above their accompanying amino acids. To locate a novel mutation from the standard “*wild-type amino acid/codon/variant amino acid*” format as normally reported in the literature (e.g., V60L in MC1R) it is necessary to carefully count the relevant nucleotide and codon numbers from the leftmost reference numbers (pencil annotations of a printout are recommended).

Phillips

2. PupaSuite can accept a list of genes using Ensembl or GeneID identifiers or can review a defined chromosome segment to search for SNPs and suggest an effect. PupaSuite is of particular interest if novel or uncharacterized SNPs are being studied as there is the opportunity to apply the same predictive tools to these loci. To explore the above three options, upload the relevant data to “Upload/paste file of genes,” “Search a region,” and “Have you got new SNPs?”, respectively. There is an option to define gene flanking regions as numbers of nucleotides upstream of the translation start site to find SNPs that may affect transcription factor binding sites. Therefore, PupaSuite is a particularly useful tool for the identification of SNP sites associated with changes to intron/exon boundaries or transcription factor binding. Lastly, additional functional annotations are provided to help assess the impact of the uploaded SNPs, including gene ontology, homology data, and OMIM references.

3.4. Simple Reviewing of SNP Haplotype Block Structure: HapMap

1. Users new to SNP analysis may hesitate before undertaking the process of analyzing human haplotype block structure in regions of interest. The accurate mapping of haplotype blocks, interpretation of D' and r^2 values, selecting tag SNPs to track blocks (3, 4, 7, 8), and assessment of genome-wide patterns of association are all specialist tasks needing care and experience (16). However, all current genetic analysis approaches require an understanding of the likely patterns of association between a set of SNPs and correlating genes or regions of interest; therefore, using HaploView within the HapMap database browser can provide a simple overview to start this process. Once HaploView has been installed on the user's own PC as a Java applet, it is possible to work directly on data from HapMap or Perlegen, but it is easier to start by configuring and viewing LD maps in the HapMap genome browser.
2. Add a gene (or region) of interest to the “Landmark or Region” search box and tick the three “Analysis” tracks: *Phased Haplotype Display*, *LD Plot*, and *tag SNP Picker*. Clearer graphics can be obtained by initially selecting one population at a time by selecting each of “Annotate LD Plot/Phased Haplotype Display” and clicking “Configure...,” then choosing a single population radio button.
3. The phased haplotype display presents the alternative haplotype blocks as blue and yellow segments matching the chromosome lengths occupied. The ease with which the user can interpret these depends on the number and length of the haplotypes in the region displayed. As an example, a very

AQ10

simple pattern is shown by ATM: a large but highly conserved gene (strong selective constraints apply to ATM, OMIM: 607585). The phased haplotype plot clearly shows that two haplotypes account for almost two equal halves of the CEU sample. No fewer than 27 of the 31 blocks define this division and the pattern is underlined by a series of identical equal-segment CEU pie charts for the genotyped SNPs in ATM. Note that at blocks 7 (left to right) and 13 a third and fourth common haplotype can be discerned and the third haplotype is characterized by different SNP alleles at blocks 17, 23, 25, 29, and 30. Finally, a singleton (literally a single CEU sample) and a minor-frequency haplotype can be seen in blocks 25 and 29, respectively. The pattern shown by ATM is, in fact, relatively common in the human genome and is termed “yin-yang haplotypes” (17).

4. The LD plots represent the extent of LD between SNPs in the region queried shown as inverted pyramid graphics. The default color scale, also in widespread use in the literature, shows maximum LD as dark red blocks and minimum LD as light gray blocks. Two example genes, CAPG and DTNBP1, illustrate how these plots can summarize both simple and complex predicted LD patterns: showing, respectively, a single, simple pyramid and multiple overlapping pyramids with heterogeneous LD values within each pyramid (checkerboards of red and gray blocks). While this partly reflects gene size and therefore SNP density (note the sevenfold difference between each gene), LD plots can provide a summary overview of recombination and SNP association patterns in the region.

5. The Tag SNP display, once configured, updates automatically between genes and many users may wish to rely just on this system to collect tag SNPs to combine with other core loci (nonsynonymous coding SNPs and translation/transcription-modifying SNPs identified by PupaSuite) to construct simple directed association studies. Although this process has largely been replaced by a standard two-stage approach of whole-genome scans then follow-up directed SNP genotyping, HapMap browsing can give a simple system for assessing the transportability, i.e., the applicability of a tag in multiple populations (8), power, and positioning of the tag SNPs that now form the core battery of markers in whole-genome analyses.

**3.5. Assessing
Population Genetics
Parameters from
Online SNP Data:
Haplotter and SPSmart**

1. Haplotter provides a useful way to begin exploring the population genetics parameters of iHS (outlined in Section 2.2.5), H , D , and F_{st} , in a genomic region. Queries are initiated by chromosome region, gene, or SNP and this will return four graphics which summarize the above-

mentioned parameters in the same order, with plots for each of the three HapMap panels (i.e., CHB and JPT populations are combined as panel ASN). The Fst graphic plots the three population comparisons to give a useful overview of genomic divergence – in particular the outliers plotted as single points can highlight those SNPs that show very strong interpopulation diversity. A table is given of iHS values around the region of interest with levels diagnostic of EHH highlighted in blue. An often-quoted example that users can investigate for themselves is the gene LCT (gene-ID 3938), this shows a very broad peak of elevated iHS in the CEU population extending well beyond the LCT chromosome region, underlined by high CEU-YRI and CEU-ASN Fst values and blue-labeled iHS levels in the accompanying table. Both the original study of selection patterns in LCT (18) and the Haplotter paper (1) ably explain these patterns.

2. In a fashion identical to Haplotter, SPSmart reviews a region, gene, or SNP list with the primary aim of summarizing the population variability found in multiple SNP databases as pie charts and key population metrics: observed H (heterozygosity), expected H , F_s , F_{st} , and *divergence* (In). Usefully SPSmart also pulls from dbSNP chromosome and position, validation status, reference and ancestral allele, and the minor allele frequency, providing alongside the population metrics a succinct one-line summary of each SNP. To explore the population variability of a set of SNPs, choose the SNP databases from HapMap phase II, HapMap phase III, Perlegen, and Stanford/Michigan CEPH-HGDP (4, 4+7, 3, and 51 populations, respectively) and provide the rs-numbers or locations. Clicking “metasearch” permits selection of a population or any combination from each of the five databases (but note the overlap between HapMap phase II and HapMap phase III) prior to uploading the SNPs of interest. For example, to review European frequency variability for the SNP rs12075, click each of the databases, tick the populations of interest, (e.g., CEU, TSI, European American, Italy-Sardinian, France Basque), add the rs12075 query to the “Search by SNPs” box, then (after choosing optional filters) click “search.” Pie charts and population metrics (and their downloadable data) are returned as separate tabs, while missing data are clearly labeled. The evident Basque divergence for rs12075 demonstrates the simplicity but informative value of HapMap style pie charts as an aid to reviewing SNP variation across multiple population-based databases.

4. Notes

1. Several approaches to database searching using a PC can help the user considerably when manipulating the data obtained from a query. Using tabbed Web page holders in the Web browser of choice (Internet Explorer; Firefox; Safari) allows simple switching between pages. While much SNP data is numerical, all information can be uploaded to a simple individual database in Excel, which now also offers sophisticated text-handling capacity, for offline processing. Although it is rarely recommended by specialists, Excel can offer a simple stand-alone database system by adapting cells to use functions such as LOOKUP, COUNTIF, or those specifically geared to database searches prefixed with “D,” such as DGET. Excel compensates for a lack of power by providing a simple and easily mastered set-up of small-scale personal databases suiting many SNP studies. The “Text to Columns” tool in the Excel “Data” menu is a straightforward way to directly process plain text files downloaded from the Web, while preserving the structure of data items separated by spaces, semicolons, or other delimiters. Simple text editors themselves are highly efficient systems for holding and searching data. For example, it is possible to find a single SNP amongst a list of 650,000 in real time using the “Find” function available in all text editors.
2. Google can be used directly to search for specific items such as rs-numbers or mitochondrial substitution sites – the latter being a particularly fruitful approach to finding medical or population studies reporting diagnostic mitochondrial haplotypes (19). For example, entering the search string “human mtDNA G6261A” into Google provides a list of papers reporting this mutation and a supposed role as a cancer risk factor. Care should be taken to ensure full use of the adjacency function of Google searches (known as the Boolean operator *NEAR*), which is not part of most genome database search engines. Therefore, to avoid very long lists of returns, it is advisable to include terms such as *human mtDNA* alongside the standard Cambridge Reference Sequence descriptions. As a compliment to PubMed queries, Google Scholar (http://scholar.google.com/advanced_scholar_search) should also be part of every researcher’s online SNP query bookmarks.
3. HapMap experienced minor problems when collating project data generated in different genotyping centers for the same SNP sites, for example, SNP rs1355497 was amongst 37 SNPs reported as showing fixed-difference allele frequencies (1) but has since been shown to be an invariant, monomorphic SNP (also *see Note 9*).

Phillips

1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248

4. Since a SNP is characterized by the context sequence each side of the nucleotide substitution site, it should be possible to uniquely define a SNP by referencing organism, chromosome, and base-pair position. However, a small but significant proportion of SNPs are nonunique, so the context sequence and its likelihood of repetition in multiple locations become critical in identifying whether a SNP is unique or not. The submission criteria of dbSNP are very effective at detecting nonunique SNPs, with a process that uses the FASTA program to check a minimum 100 bp flanking sequence to assess if the SNP can be positioned uniquely in the genome and can be matched with other submissions of the same SNP. The proportion of nonunique SNPs remains very small in dbSNP at about 5% and is much more common in certain regions, e.g., pericentromeric areas of each chromosome.
5. Very useful and readable guides to the routine use of the NCBI sites are detailed in a PDF handbook that can be downloaded chapter by chapter (<http://www.ncbi.nlm.nih.gov/books/bv.fcgi?rid=handbook.part.1>). Chapters particularly relevant to SNP research include Chapters 2 (PubMed), 5 (dbSNP), 7 (OMIM), 15 (Entrez), 16 (BLAST), and 20 (Map Viewer). Download them by clicking the PDF icon on each chapter summary page.
6. SNP sites are still routinely described by the amino acid substitution they create rather than an rs-number, particularly if they are mutations or rare enough to escape detection by dbSNP. The easiest way to obtain the rs-number (if it exists) is to record the gene identification number from NCBI Gene (e.g., query “MC1R AND human” gives GeneID 4157) then go to the coding SNP part of dbSNP using the following URL finishing with the ID number to obtain the listing of known coding substitutions and their affected amino acid residues: http://www.ncbi.nlm.nih.gov/SNP/snp_ref.cgi?locusId=number. For example, in gene MC1R, R151C is amongst the most commonly described coding SNPs, and was revealed to be rs1805007 using the procedure described above.
7. Human genes are consistently identified across different databases by a *gene symbol* (sometimes termed the “*gene name*”) comprising a series of uppercase letters and numbers (<http://www.genenames.org/>), provided by HUGO. A *gene ID* consists of a number alone and refers to the number codes given to each gene by NCBI Gene. These can be used both within NCBI as the main point of reference for a gene (e.g., when reviewing coding SNPs) and in certain other databases. A useful gene ID converter tool is provided at <http://idconverter.bioinfo.cnio.es/>.

1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296

8. Use of the single map weight filter in EntrezSNP does not lead to the exclusion of all the SNPs in dbSNP with different Celera locations; however, the list of returns is headed by SNPs that carry the warning “*Mapped unambiguously on non-reference assembly only.*” Note, however, that when one uses the NCBI MapViewer (*see Section 3.2.5*) both reference sequence and Celera locations are marked on the genome-wide chromosome map. Therefore, it is advisable to include the term *NOT Celera* at the end of a multiple SNP list in MapViewer queries.
9. Resequencing is the only method of SNP characterization that avoids acquisition bias. This is the phenomenon where the characteristics of the SNP itself affect the chances of its detection by genotyping methods. Examples of SNP features that mean the loci are either not detected or incorrectly genotyped by large-scale projects such as HapMap include triallelic SNPs (the medically important CRP promoter rs3091244 being a notable example), SNPs with very low frequency minor alleles (also missed by resequencing if insufficient samples are typed), and SNPs with very dense arrays of closely neighboring SNPs such as those of the hypervariable major histocompatibility complex. Acquisition bias can also describe the process of selecting SNPs from databases using criteria which bias the SNP lists produced from a query.
10. Often the Fasta section lists less than 100 bp of context sequence each side of a SNP (e.g., rs1805009) – often owing to the fact that a submitting laboratory only provided short segments. The easiest way to obtain ± 100 bp of context sequence for assay primer design purposes is to use the Santa Cruz genome assembly. Add the SNP rs-number to the URL <http://genome.ucsc.edu/cgi-bin/hgTracks?position=rsnumber> (several dbSNP builds available), click the highlighted SNP in the map, click “view DNA for this feature,” then opt for 100 bases upstream/downstream. The SNP base is the reference allele and is not marked so it is best to use 100+0 and 0 +100 in two separate sequence dumps. Another potential problem in the Cluster Report Fasta section is the occasional (and apparently ad hoc) listing of neighbor SNPs as IUPAC codes (*see Table 3.4*). For example, rs1805006 includes no fewer than six other SNPs in ± 100 bp of sequence, given as K, R, R, Y, R, R (in that order), that may cause problems once the sequence is inserted into primer design software. Processing the SNP context sequence directly from Santa Cruz avoids this problem.
11. Fasta section nucleotides are presented in two ways: in uppercase/lowercase letters and in black/green. Uppercase letters denote a normal, unique, genomic sequence, while lowercase letters are used for a sequence identified by RepeatMasker

Phillips

(<http://repeatmasker.genome.washington.edu/cgi/bin/RepeatMasker>) as a low-complexity or repetitive element sequence. Green denotes a sequence identified by the submitter during the SNP assay process (a single green SNP base signifying identification by sequence comparison), while black denotes a flanking sequence used by NCBI from the nucleotide databases as part of the SNP submission checks.

Acknowledgements

The author would like to thank Maviky Lareu, Antonio Salas, and Angel Carracedo, University of Santiago de Compostela, for useful discussions in the preparation of this chapter. The work was in part supported by funding from Xunta de Galicia: PGIDTIT06P-XIB228195PR and the Spanish MEC: BIO2006-06178.

References

1. The International HapMap Consortium. (2005) A haplotype map of the human genome. *Nature* **437**, 1299–1320.
2. The International HapMap Consortium. (2007) A second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**, 851–861.
3. Gabriel, S. B., Schaffner, S. F., Nguyen, H., Moore, J. M., Roy, J., Blumenstiel, B. et al. (2002) The structure of haplotype blocks in the human genome. *Science* **296**, 2225–2229.
4. Wall, J. D. and Pritchard, J. K. (2003) Haplotype blocks and linkage disequilibrium in the human genome. *Nat. Rev. Genet* **4**, 587–597.
5. Ashurst, J. L., Chen, C. K., Gilbert, J. G., Jekosch K., Keenan S., Meidl P. et al. (2005) The Vertebrate Genome Annotation (Vega) database. *Nucleic Acids Res.* **33**, D459–465.
6. Yang, N., Li, H., Criswell, L. A., Gregersen, P. K., Alarcon-Riquelme, M. E., Kittles, R. et al. (2005) Examination of ancestry and ethnic affiliation using highly informative diallelic DNA markers. *Hum. Genet.* **118**, 382–392.
7. de Bakker, P. I., Yelensky, R., Pe'er, I., Gabriel, S. B., Daly, M. J. and Altshuler, D. (2005) Efficiency and power in genetic association studies. *Nat. Genet.* **37**, 1217–1223.
8. de Bakker, P. I. W., Noel, N. P., Burtt, N. P., Graham, R. R., Guiducci, C., Yelensky, R., Drake, J.A. et al. (2006) Transferability of tag SNPs in genetic association studies in multiple populations. *Nat. Genet.* **38**, 1298–1303.
9. Riva, A. and Kohane, I. S. (2002) SNPper: retrieval and analysis of human SNPs. *Bioinformatics* **18**, 1681–1685.
10. Conde, L., Vaquerizas, J. M., Santoyo, J., Shahrou, F., Ruiz-Llrente, S., Robledo, M. et al. (2004) PupaSNP Finder: a web tool for finding SNPs with putative effect at transcriptional level. *Nucleic Acids Res.* **32**, W242–248.
11. Conde L., Vaquerizas J.M., Dopazo H., Arbiza L., Reumers J., Rousseau F. et al. (2006) PupaSuite: finding functional single nucleotide polymorphisms for large-scale genotyping purposes. *Nucleic Acids Res.* **34**, W621–625.
12. Ramensky, V., Bork, P., and Sunyaev, S. (2002) Human non-synonymous SNPs: server and survey. *Nucleic Acids Res* **30**, 3894–3900.
13. Cartegni, L., Wang, J., Zhu, Z., Zhang, M. Q., and Krainer, A. R. (2003) ESEfinder: A web resource to identify exonic splicing enhancers. *Nucleic Acids Res.* **31**, 3568–3571.
14. Voight, B. F., Kudaravalli, S., Wen, X. and Pritchard, J. K. (2006) A map of recent positive selection in the human genome. *PLoS Biol* **4**, 446–458.
15. Sabeti, P. C., Reich, D. E., Higgins, J. M., Levine, H. Z. P., Richter, D. J., Schaffner, S. F. et al. (2002) Detecting recent positive selection in the human genome from haplotype structure. *Nature* **419**, 832–837.

SNP Databases

- 1345 16. Haiman, C. A. and Stram, D. O. (2008)
1346 Utilizing HapMap and tagging SNPs. *Methods Mol. Med.* **141**, 37–54.
1347
- 1348 17. Zhang, J., Rowe, W.L., Clark, A.G. and
1349 Buetow, K.H. (2003) Genomewide distribution
1350 of high-frequency, completely mismatching SNP
1351 haplotype pairs observed to be common across
1352 human populations. *Am. J. Hum. Genet.* **73**, 1073–1081.
1353
- 1354
- 1355
- 1356
- 1357
- 1358
- 1359
- 1360
- 1361
- 1362
- 1363
- 1364
- 1365
- 1366
- 1367
- 1368
- 1369
- 1370
- 1371
- 1372
- 1373
- 1374
- 1375
- 1376
- 1377
- 1378
- 1379
- 1380
- 1381
- 1382
- 1383
- 1384
- 1385
- 1386
- 1387
- 1388
- 1389
- 1390
- 1391
- 1392
18. Bersaglieri, T., Sabeti, P. C., Patterson, N.,
Vanderploeg, T., Schaffner, S.F., Drake J.A.
et al. (2004) Genetic signatures of strong
recent positive selection at the lactase gene.
Am. J. Hum. Genet. **74**, 1111–1120.
19. Bandelt, H. J., Salas, A. and Bravi, C. M.
(2006) What is a 'novel' mtDNA mutation—
and does 'novelty' really matter? *J. Hum.
Genet.* **51**, 1073–1082.