



A Suite Of Programs For Pre- And Post-imputation Data Checking

N.W. Rayner^{1,2,3}, N. Robertson^{1,2}, Anubha Mahajan¹ M.I. McCarthy^{1,2,4}

1) WTCHG, University Oxford, Oxford, United Kingdom; 2) OCDEM, University Oxford, United Kingdom; 3) Wellcome Trust Sanger Institute, Hinxton, Cambridge, United Kingdom ;4) NIHR Oxford Biomedical Research Centre, Churchill Hospital, Oxford, United Kingdom

Contact:

wrayner@well.ox.ac.uk

Introduction

- The availability of the public imputation servers at the Wellcome Trust Sanger Institute and the University of Michigan has greatly reduced the complexity associated with performing an imputation run.
- Despite this, thorough quality control prior to imputation is still vital to ensure the data are correctly aligned to the reference genome.
- Post-imputation, verifying the run has completed successfully can be tricky to do from log files.
- To simplify these processes we have developed a suite of programs in Perl and Java that check and summarise the data at both the pre- and post-imputation stages.

Pre-imputation Checking

- The pre-imputation checking program compares a Plink format marker (.bim) and frequency (.freq) file to the selected reference panel.
- Currently 1000G phase 3 and HRC r1 and r1.1 reference panels are supported.
- The program checks all variants in the .bim file that match the reference for a variety of potential issues (Table 1) using both location information as well as SNP names.

Issue Checked	Removed or updated?
Chromosome and position	Updated
Alleles	Removed
Allele frequency (AF)	Removed if AF difference > 0.2*
Strand	Updated
Palindromic SNPs	Removed if MAF > 0.4
Reference allele	Updated
Variant naming	Updated**

Table 1: Potential issues checked for in .bim file. *AF difference used can be set by the user, or removed entirely. **Variant names are not automatically updated although a list of naming differences is produced.

- Summary statistics for the run are written to a log file (Figure 1) and for each potentially updatable issue a list of variants affected is produced.
- Variants not matching the reference or failing one of the above cut offs are listed for removal.

```
Options Set:
Reference Panel: HRC
Bim filename: P1.qcPlus.bim
Reference filename: HRC.r1-1.0RCH37.wgs.mac5.alleles.tab.gz
Allele frequencies filename: P1.qcPlus.frequency.freq
Allele frequency threshold: 0.2

Matching to HRC

Position Matches
ID matches HRC 480362
ID Doesn't match HRC 48139
Total Position Matches 724481
ID Match
Different position to HRC 1953
No Match to HRC 1661
Skipped (X, Y, MZ) 0
Total in bim file 732139
Total processed 732095

Indels (ignored in r1) 1
SNPs not changed 210061
SNPs to change ref alt 314530
Strand ok 726545
Total Strand ok 726581

Strand to change 56
Total checked 730434
Total checked Strand 726601
Total removed for allele frequency diff > 0.2 605
Palindromic SNPs with Freq > 0.4 2076

Non Matching alleles 1757
ID and allele mismatching 1005; where HRC is . 935
Duplicates removed 23
```

Figure 1: Log file from the pre-imputation checking.

- A shell script of Plink commands is created to automate the use of these variant lists within Plink to update or remove variants.

Post-Imputation Checking

- The post-imputation checking program takes information from the output of the current set of imputation programs or servers and produces a range of charts and tables based on information score, alternate AF, MAF and position.
- The resulting plots allow a quick visual assessment of the quality of an imputation run (Figures 2-7).

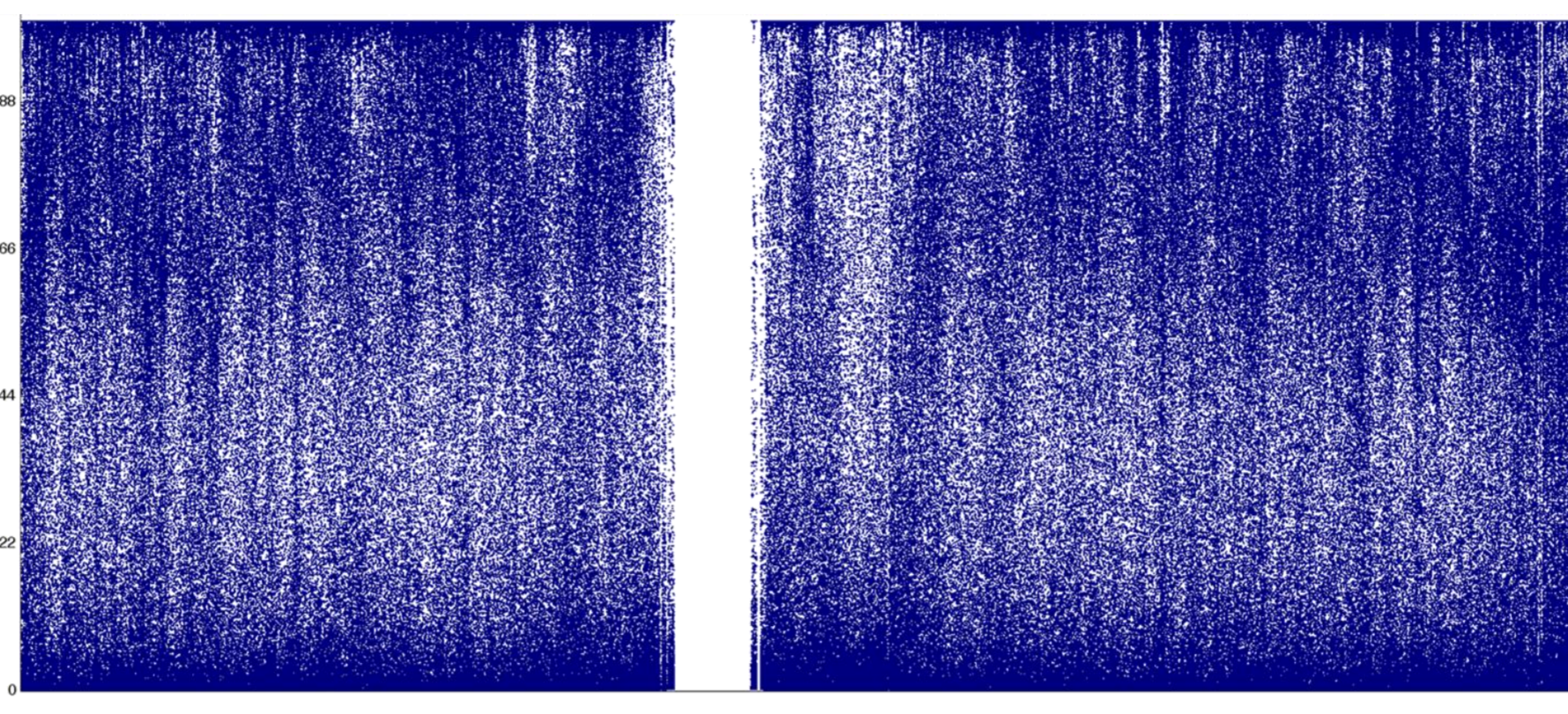


Figure 2: Information score plotted per chromosome. Plots are coloured red if there are 2 sections of >1MB without variants, indicating possible imputation failures.

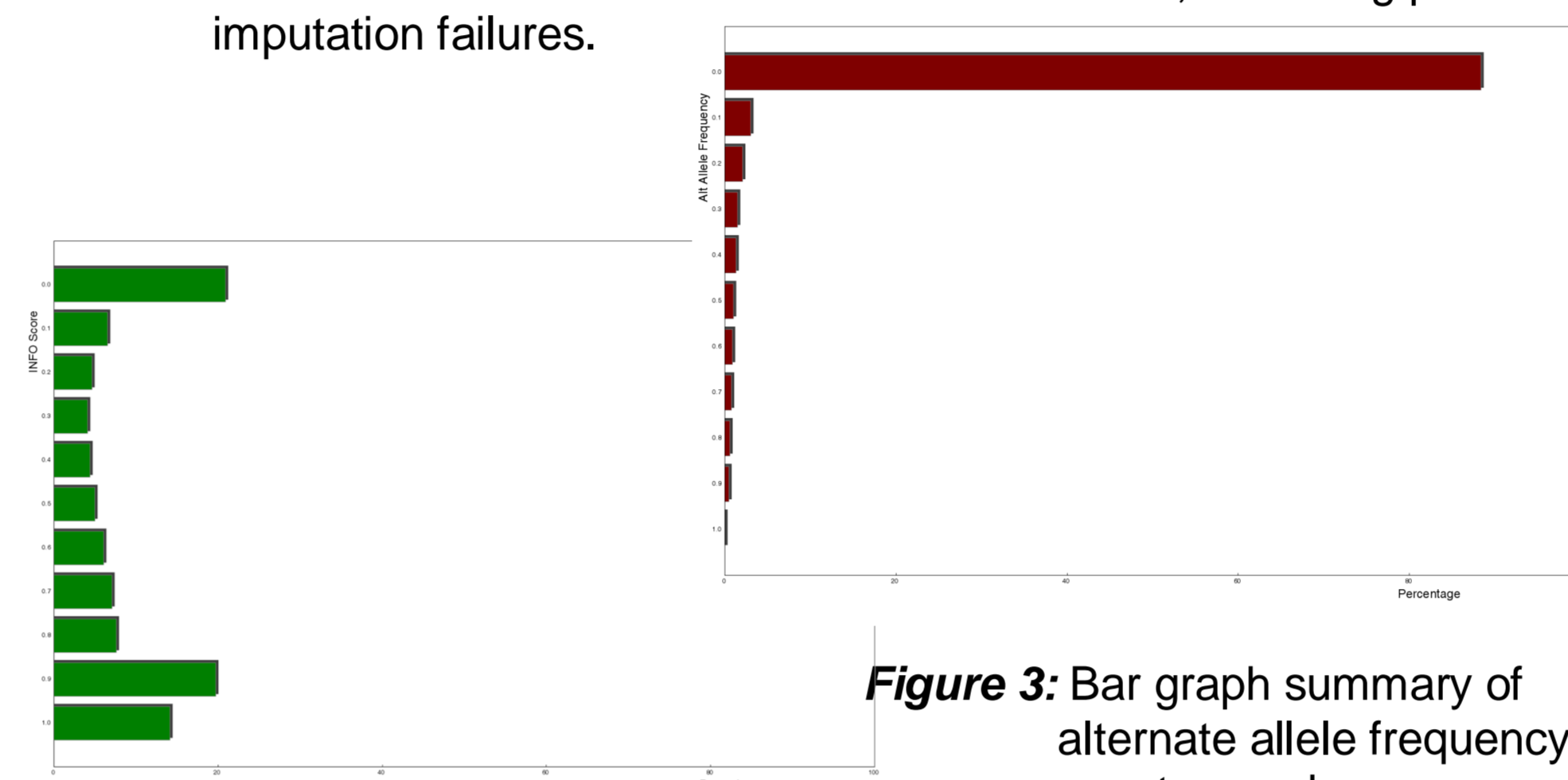


Figure 4: Bar graph summary of Information score counts per chromosome.

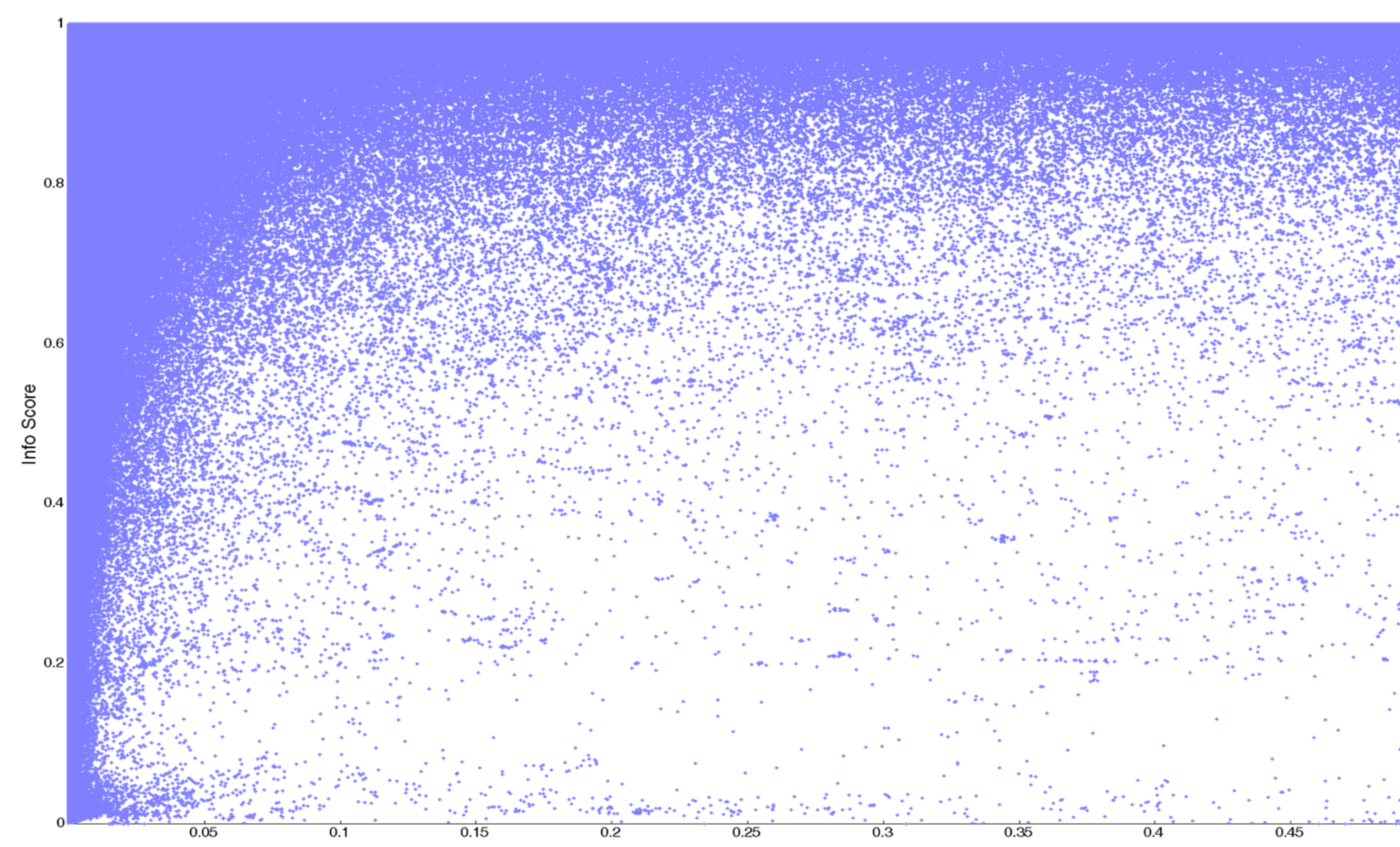


Figure 5: Plot of information score vs MAF across the chromosome.

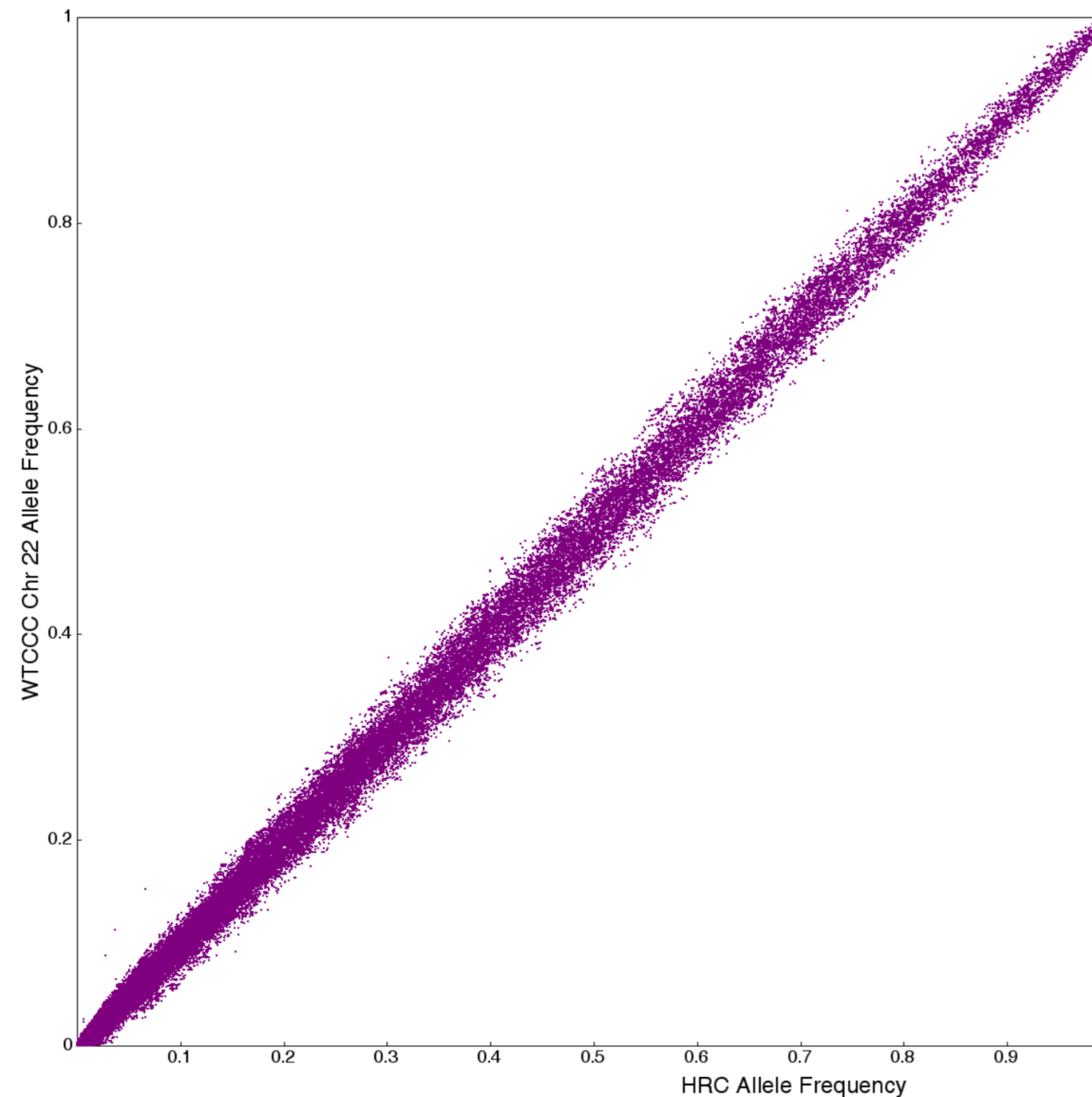


Figure 6: Plot of the allele frequency comparison to the reference (HRC in this case).

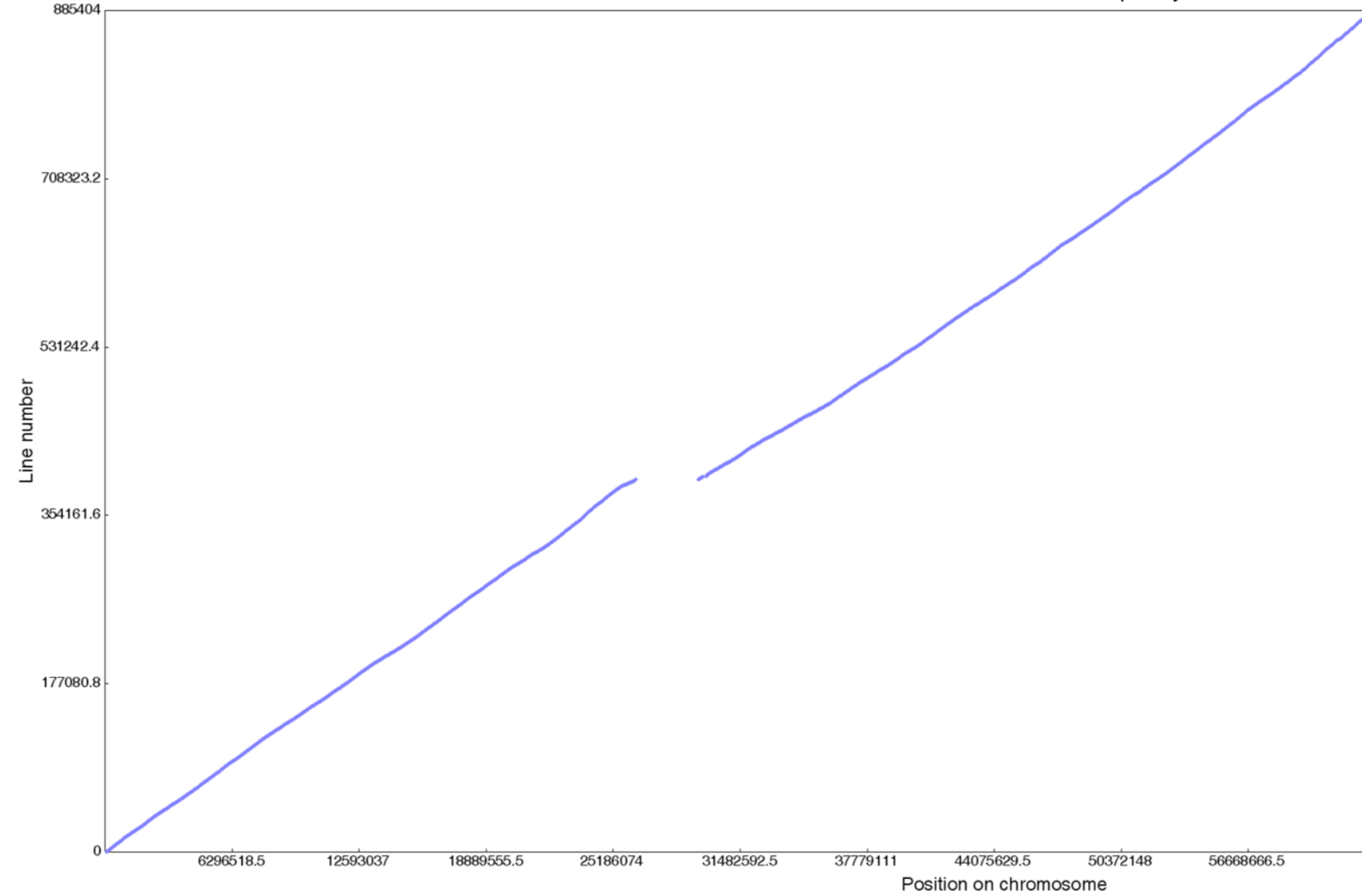
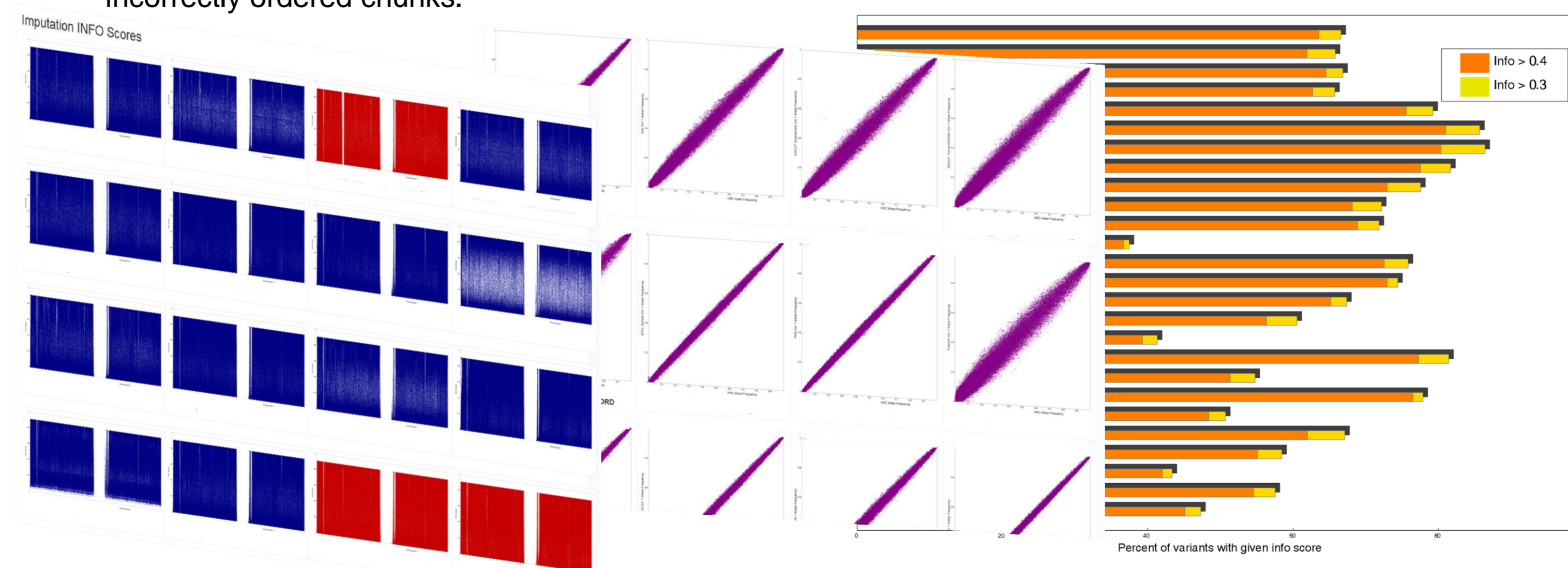


Figure 7: Plot of chromosome position vs. position in file to highlight failed or incorrectly ordered chunks.



Visualisation

- As the output for a single imputation run consists of 132 images and 70 tables for ease of visualisation these are collated into a single html file (Figure 8) viewable in any current web browser.
- The program can also be run to summarise and compare multiple imputed data sets (Figure 9).

Availability

- All the programs and instructions are available at: www.well.ox.ac.uk/~wrayner/tools

Future Work

- The post-imputation program is reliant on external libraries for chart plotting, we plan to offer the program via a web server to avoid installation issues.

Summary

- These programs simplify the checking of large data sets and have been used extensively in single cohort studies, as well as for large consortia such as DIAMANTE where >35 genome-wide data sets were prepared using our pre-imputation program and >25 were checked post-imputation.
- Using the post-imputation program we were able to easily identify a wide range of imputation issues such as AF differences and missing or duplicated regions.

Figure 8: (right) Single html page containing tables, embedded images and the navigation bar.

Figure 9: (below) Plots of multiple cohort summaries to enable quick comparison.

Imputation Report For Cohort: WTCCC

Chromosome index, 22 in total. Tap to the section for a particular chromosome using the buttons below.

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

Cohort Metadata

Information captured during the processing of this cohort

Property	Value
FORMAT	OXFORD
UPLOADED	17/02/2016

Genome-Wide

A: Alternate Allele Frequency: Genome-Wide

Alt Allele Freq	Count	%
0.0	3450271	88.131
0.1	1183640	3.019
0.2	736565	2.001
0.3	587659	1.525
0.4	481965	1.239
0.5	392059	1.003
0.6	340346	0.888
0.7	282086	0.722
0.8	230072	0.611
0.9	135466	0.356
1.0	11580	0.030
Total	3903319	100.000

B: Imputation INFO Score: Genome-Wide

INFO Score	Count	%
0.0	8189206	20.890
0.1	2519907	6.428
0.2	1771425	4.519
0.3	1572599	4.011
0.4	1682161	4.291
0.5	1929613	4.922
0.6	2360593	6.024
0.7	2720880	6.979
0.8	2922346	7.482
0.9	2679712	68.160
1.0	5623307	14.344
Total	3903319	100.000

Chromosome 1

A: Alternate Allele Frequency vs HRC

B: Alternate Allele Frequency

C: Alternate Allele Frequency

Alt Allele Freq	Count	%
0.0	2710694	88.3215
0.1	88676	2.8284
0.2	57673	1.8791
0.3	42529	1.4700
0.4	27526	1.0166
0.5	30222	0.9847
0.6	27396	0.8306
0.7	22110	0.7501
0.8	18058	0.6300
0.9	28994	0.9317
1.0	1165	0.0380
Total	3069122	

D: Imputation INFO Score

E: Imputation INFO Score

INFO Score	Count	%
0.0	442258	20.584
0.1	139462	6.514
0.2	141613	6.514
0.3	125990	4.9920
0.4	134933	4.3716
0.5	152256	4.9609
0.6	186148	6.0822
0.7	216085	7.0456
0.8	232537	7.5767
0.9	604522	19.6969
1.0	433918	14.1382
Total	3069122	

F: Imputation INFO Score

G: Order of Chr Pos