

# Visualizing the Geography of Genetic Variants

Joseph H. Marcus<sup>\*,1</sup> and John Novembre<sup>†,1</sup>

<sup>1</sup>Department of Human Genetics, University of Chicago, Chicago, IL

## Abstract

One of the key characteristics of any genetic variant, beyond its potential phenotypic effects or its frequency, is its geographic distribution. The geographic distribution of a genetic variant can shed light on where the variant first arose, in what populations spread within, and in turn can help one learn about historical patterns of migration, genetic drift, and natural selection. Collectively the geographic distribution of genetic variants can help to explain how populations have been related through time (e.g. levels of gene flow and divergence). For these reasons, visual inspection of geographic maps for genetic variants is common practice in genetic studies. Here we develop an interactive visualization tool for illuminating the geographic distribution of genetic variants. Through an efficient RESTful API and dynamic front-end the Geography of Genetic Variants (GGV) browser rapidly provides maps of allele frequencies in populations distributed across the globe.

\*josephhmarcus@gmail.com

<sup>†</sup>jnovembre@uchicago.edu

## 16 Introduction

17 Genetics researchers often face the following problem: the researcher has identified one or more  
18 genetic variants of interest (e.g. from a genome-wide association, eQTL, or pharmacogenomic study)  
19 and would like to know the geographic distribution of the variant. For example, the researcher may  
20 hope to address: 1) implications for genomic medicine (e.g. is a risk allele geographically localized  
21 to a certain patient population?) (e.g. ???); 2) design of follow-up studies (e.g. what population  
22 should be studied to observe variant carriers?) (e.g. ???); or 3) the biology of the variant in question  
23 (e.g. does the variant correlate with a known environmental factor in a manner suggestive of some  
24 geographically localized selection pressure?) (Novembre and Di Rienzo 2009). A simple geographic  
25 map of the distribution of a genetic variant can be incredibly insightful for these questions.

26 Contemporary population genetics researchers are also faced with the challenge of large, high  
27 dimensional datasets. For example, it is not uncommon for a researcher in human genetics to have a  
28 dataset comprised of thousands of individuals measured at hundreds of thousands or even millions of  
29 single nucleotide variants (SNVs). One common approach to visualizing such high-dimensional data  
30 is to compress the SNV dimensions down to a small number of latent factors, using a method such  
31 as principal components analysis (???, Patterson et al 2009), a model-based clustering method such  
32 as STRUCTURE (???) or ADMIXTURE (???). While these approaches are extremely valuable,  
33 researchers can use them to often without inspecting the underlying variant data in more detail. A  
34 natural approach to gaining more insight to the overall structure of a population genetic dataset is  
35 to visually inspect what geographic patterns arise in allele frequency maps.

36 Unfortunately, generating geographic allele frequency maps is time-consuming for the average  
37 researcher as it requires a combination of data wrangling methods ((??)) and geographic plotting  
38 techniques that are unfamiliar to most. Our aim here is to produce a tailored system for rapidly  
39 constructing informative geographic maps of allele frequency variation.

40 Our work is inspired by past tools such as the ALFRED database (xx) and the maps available  
41 on the HGDP Selection browser ((??)). One of us (J.N.) developed the scripts for the HGDP  
42 Selection Browser maps using GMTtools ((??)), a powerful system of geographic plotting scripts  
43 with a substantial learning curve. The plots from the HGDP Selection Browser have proved useful,  
44 and have appeared in research articles (e.g exExample), books (e.g. (??)), and have been ported  
45 and made available on the UCSC Genome Browser (available under the XX track of the browser).

46 Reference datasets for population genetic variation has greatly expanded since the release of the  
47 HGDP Illumina 650Y dataset ((??)) that formed the basis of the HGDP Selection Browser maps.  
48 The most notable advance is the publication of the 1000 Genomes phase 3 data (@1000GenomesCon-  
49 sortium) though additional datasets are coming online ((??), (??)).

50 In addition, novel approaches for data visualization have become more widely available. In  
51 particular, web-based visualization tools, such as Data Driven Documents (d3.js), offer powerful  
52 approaches to realize these aims through useful methods for interactivity, advantages of software  
53 development in modern web-browsers, a large open-source development community, and ease of  
54 share-ability (Bostock, Ogievetsky, and Heer 2011).

55 Taking advantage of these recent advances, we aim to address significant visualization chal-  
56 lenges that are inherit in the production of geographic allele frequency maps, including dynamic  
57 interaction, display of rare genetic variation, and representation of uncertainty in estimated allele  
58 frequencies due to variable sample sizes.

59 **Fundamental Approach**

60 The Geography of Genetic Variants browser (GGV) uses the SVG and mapping utilities of d3.js  
61 to generate interactive frequency maps, allowing for quick and dynamic displays of the geographic  
62 distribution of a genetic variant.

63 In order to allow for easy access to large commonly used public genomic datasets, such as the  
64 1000genomes project, Human Genome Diversity project and POPRES project, we have developed  
65 an SQL database and RESTful API for querying allele frequencies by chromosome and position or  
66 reference SNP identifier [Figure 1].

67 Users can query single variants by chromosome and position identifiers, by rsids ((??)), or  
68 users can simply choose a random variant from within a dataset. While many applications require  
69 inspection of the distribution of a specific variant, from our experience, it can be very helpful to  
70 quickly view the geographic distribution of several randomly chosen variants to quickly gain a sense  
71 of structure in a dataset. For instance, in data with deep population subdivision, it is obvious in  
72 the consistent patterns of differentiation observed across markers. We also expect this will be useful  
73 in teaching contexts – as it provides a highly visual way for learners to understand human genetic  
74 variation.

75 After a query, the GGV displays the frequency of a variant in a given population as a pie chart  
76 where each slice represents minor and major allele frequencies. Pie charts are displayed as points  
77 positioned at the population’s region of origin where the projection/map-view is chosen based off  
78 of the geographic proximity of populations present in a given dataset [Figure 2].

79 **Representing variable certainty in frequency data**

80 One under-appreciated problem with allele frequency maps is that not all data points have equal  
81 levels of certainty. For some locations, sample sizes are small, and the reported allele frequency may  
82 be quite far from the true population frequency due to sampling error. To address this issue, we  
83 use varying transparency in a population’s pie chart: sample allele frequencies with higher levels of  
84 sampling error will be made more transparent, and hence less visible on the map [Figure 3].

85 **Representing rare variants in frequency data**

86 An additional challenge is that allele frequencies between variants often differ greatly, sometimes  
87 by orders of magnitude in a single dataset. This has not been an pervasive problem until recently,  
88 as most population genetic samples were based on genotype arrays biased towards variants that  
89 are common in human populations (5-50% in minor allele frequency). With the advent of next  
90 generation sequencing technologies and large samples of thousands of individuals, it is common for  
91 datasets to contain rare variants ((@1000genomes2012integrated, Nelson et al. 2012, Tennessen et  
92 al. (2012)).

93 In current visualization schemes, such as the HGDP Selection Browser, rare variants would  
94 be represented as narrow slivers in a pie chart, nearly invisible to the naked eye. To address  
95 this challenge we re-scale frequencies for rare variants, so that small frequencies become visible.  
96 Specifically, we will use a “frequency scale” that is indicated below the map, and redundantly using  
97 color, that will indicate a constant scale for the frequencies displayed [Figure 4]. Much like scientific  
98 notation, this allows a wide range of frequencies to be displayed].

<sup>99</sup> **Interface features**

<sup>100</sup> In many datasets where populations are sampled densely in geographic space, one problem is that  
<sup>101</sup> allele frequency plots begin to overlap each other and obscure information. To address this issue,  
<sup>102</sup> we use force-directed layouts of the populations such that no two points are overlapping each other,  
<sup>103</sup> and yet the points will be pulled towards their true origins. Lines anchoring the points visibly to  
<sup>104</sup> their original sampling locations will be used to make sure true sampling locations are indicated  
<sup>105</sup> [Figure 5].

<sup>106</sup> **Underlying Frequency API**

<sup>107</sup> [Joe - worth saying more about the API structure and giving example calls and structure of JSON?]

<sup>108</sup> **Caveats**

<sup>109</sup> A major challenge of using a geographic representation of genetic variation in humans is that the  
<sup>110</sup> samples must be associated with a geographic location. While doing so is generally immensely help-  
<sup>111</sup> ful, it has inherent complexity. For example, practitioners must make choices regarding representing  
<sup>112</sup> where an individual was sampled for the study (e.g. the city of a major research center) or choos-  
<sup>113</sup> ing a location that is more representative of an individual's ancestral origins (e.g. as determined  
<sup>114</sup> in practice by the birthplaces of recent ancestors, such as grandparents). We do not proscribe a  
<sup>115</sup> general solution to this problem and instead use a set of locations for each dataset that aligns with  
<sup>116</sup> those used by the initial analysts of the data. A future feature will be to allow alternative location  
<sup>117</sup> schema to be used for the populations in a dataset.

<sup>118</sup> **Conclusion**

<sup>119</sup> By allowing rapid generation of allele frequency maps, we hope to facilitate the interpretation  
<sup>120</sup> of variant function and history by practicing geneticists. We also hope the ability to query ran-  
<sup>121</sup> dom variants from major human population genetic samples will allow students to appreciate the  
<sup>122</sup> structure of human genetic diversity in a more approachable and intuitive form than alternative  
<sup>123</sup> visualizations.

<sup>124</sup> **Acknowledgements**

<sup>125</sup> Support for this work was provided by the NIH-BD2K initiative (1U01 CA198933-0). The authors  
<sup>126</sup> would also like to thank XX for supportive conversations.

<sup>127</sup> **To Add**

- <sup>128</sup> • Provide information about links to ggv, freq\_api, and github for contributions
- <sup>129</sup> • Fill up citation and fix citation formatting .bst?
- <sup>130</sup> • Fill in figures

<sup>131</sup> **Acknowledgements**

<sup>132</sup> **References**

- <sup>133</sup> Bostock, Michael, Vadim Ogievetsky, and Jeffrey Heer. 2011. “D<sup>3</sup> Data-Driven Documents.” *Visualization and Computer Graphics, IEEE Transactions on* 17 (12). IEEE: 2301–9.
- <sup>134</sup> Nelson, Matthew R, Daniel Wegmann, Margaret G Ehm, Darren Kessner, Pamela St Jean, Claudio Verzilli, Judong Shen, et al. 2012. “An Abundance of Rare Functional Variants in 202 Drug Target Genes Sequenced in 14,002 People.” *Science* 337 (6090). American Association for the Advancement of Science: 100–104.
- <sup>135</sup> Novembre, John, and Anna Di Rienzo. 2009. “Spatial Patterns of Variation Due to Natural Selection in Humans.” *Nature Reviews Genetics* 10 (11). Nature Publishing Group: 745–55.
- <sup>136</sup> Tennessen, Jacob A, Abigail W Bigham, Timothy D O’Connor, Wenqing Fu, Eimear E Kenny, Simon Gravel, Sean McGee, et al. 2012. “Evolution and Functional Impact of Rare Coding Variation from Deep Sequencing of Human Exomes.” *Science* 337 (6090). American Association for the Advancement of Science: 64–69.