

Qatar Genome: Insights on Genomics from the Middle East

Hamdi Mbarek¹, Geethanjali Devadoss Gandhi^{2,6}, Senthil Selvaraj², Wadha Al-Muftah^{1,8}, Radja Badji¹, Yasser Al-Sarraj^{1,6}, Chadi Saad¹, Dima Darwish¹, Muhammad Alvi¹, Tasnim Fadl¹, Heba Yasin¹, Fatima Alkuwari¹, Rozaimi Razali², Waleed Aamer², Fatemeh Abbaszadeh³, Ikhlaq Ahmed⁴, Younes Mokrab², Karsten Suhre⁵, Omar Albagha^{6,7}, Khalid Fakhro², Ramin Badii³, Said I. Ismail¹ and Asma Althani^{1,9} for the Qatar Genome Program Research Consortium*

Affiliations

1 Qatar Genome Program, Qatar Foundation Research, Development and Innovation, Qatar Foundation, Doha, Qatar.

2 Human Genetics Department, Sidra Medicine, Doha, Qatar.

3 Molecular Genetics Laboratory, Hamad Medical Corporation, Doha, Qatar.

4 Sidra Medicine, Biomedical Informatics - Research Branch, Doha, Qatar.

5 Bioinformatics Core, Weill Cornell Medicine-Qatar, Education City, Doha, Qatar.

6 College of Health & Life Sciences, Hamad Bin Khalifa University, Education City, Doha, Qatar.

7 Center of Genomic and Experimental Medicine, University of Edinburgh, Edinburgh, UK.

8 Department of Genetic Medicine, Weill Cornell Medicine-Qatar, Doha, Qatar.

9 Biomedical Research Center, Qatar University, Doha, Qatar.

*A list of authors appears at the end of the manuscript

Correspondence should be addressed to H.M (hmbarek@qf.org.qa), S.I.I (saismail@qf.org.qa)

Abstract

Despite recent biomedical breakthroughs and large genomic studies growing momentum, the Middle Eastern population, home to over 4000 million people, is under-represented in the human genome variation databases. Here we describe insights from phase 1 of the Qatar Genome Program which whole genome sequenced 6,045 individuals from Qatar. We identified more than 88 million variants of which 24 million are novel and 23 million are singletons. Consistent with the high consanguinity and founder effects in the region, we found that several rare deleterious variants were more common in the Qatari population while others seem to provide protection against diseases and have shaped the genetic architecture of adaptive phenotypes. Insights into the genetic structure of the Qatari population revealed five non-admixed subgroups. Based on sequence data, we obtained and the heritability and genetic marker associations for 45 clinical traits. These results highlight the value of our data as a resource to advance genetic studies in the Arab and neighbouring Middle Eastern populations and will significantly boost the current efforts to improve our understanding of global patterns of human variations, human history and genetic contributions to health and diseases in diverse populations.

Keywords: Population Genomics, Large-Scale Sequencing Project, Diversity, Qatar, Middle East

Introduction

Several countries worldwide have initiated large-scale population genomics projects representing various regions from Africa, Europe, North and South America, South Asia and Australia

(Gudbjartsson, Helgason, et al., 2015; Gudbjartsson, Sulem, et al., 2015; Gurdasani et al., 2019; Manolio et al., 2019; Naslavsky et al., 2020; Stark et al., 2019; Turro et al., 2020; Wu et al., 2019b). In addition to this groundbreaking work, there are also ongoing large collaborative efforts to increase diversity in human genetics, including the All of Us Research Program (Collins & Varmus, 2015), the Human Health and Heredity in Africa (H3Africa) Initiative (C. Rotimi et al., 2014) and the TOPMed Program (Taliun et al., 2021). Such studies provided valuable new insight into human disease, population structure and history of migration (Boomsma et al., 2014; Chiang et al., 2018; Francioli et al., 2014; Gurdasani et al., 2019; Okada et al., 2018; Scott et al., 2016a; Wu et al., 2019a). Despite this notable focus on diversity, there is still considerable effort needed to cover the broad diversity of world ancestries to ensure that discoveries does not conserve historical disparities and to uncover the various diseases etiologies that remain uncharacterized to date (Bentley et al., 2017; Landry et al., 2018; Mills & Rahal, 2019). The Middle-East regions are still underrepresented in the public databases (Abou Tayoun & Rehm, 2020). For instance, the latest version of gnomAD database (3.1) contains data from only 158 Middle-eastern genomes (Karczewski et al., 2020). The Qatar Genome Program (QGP) is a population genome project based in Qatar aiming to sequence the genomes of local population for the purpose of supporting genomic medicine in the country and the region. As part of phase 1, it has sequenced the whole genomes of 6,045 subjects whose specimens were collected and biobanked by the Qatar Biobank (QBB) (Al Thani et al., 2019b) Figure 1a).

Qatar occupies a relatively small surface area of 11,521 km² on the western coast of the Arabian Gulf. Qatar shares its southern border with Saudi Arabia and a maritime border with Bahrain, UAE, and Iran (Figure 1b) and has a population of approximately 2.8 million. The country is located at a historic intersection of ancient and recent migration and admixture (Arauna et al.,

2017; Hellenthal et al., 2014). Similar to other countries in the region, it is known for its unique population structure that is characterized by a high consanguinity rate and increased prevalence of rare genetic diseases (Al-Gazali et al., 2006; Anwar et al., 2014; Hunter-Zinck et al., 2010; Rodriguez-Flores et al., 2014, 2016; Scott et al., 2016a). Recent genetic studies identified indigenous Arabs as the direct descendants of the first Eurasian populations established by early migrations out of Africa (Bentley et al., 2017) (Figure 1c). Moreover, sizable proportions of the population have more recent Persian and African ancestry (Harkness & Khaled, 2014). QBB includes comprehensive phenotyping, providing excellent synergy for discovery when combined with the WGS data, that also enable accurate estimate of allele frequencies for rare and common variants, and well-defined polygenic risk scores for many disease traits. All such features of the local population potentiate discoveries, not only related to millions of people in the immediate neighboring region, but also inform genetic studies in other parts of the world.

Materials and Methods

Qatar Biobank subject recruitment

The Qatar Biobank (QBB) is a longitudinal population-based cohort study examining a population sample of permanent Qatari residents (Qatari nationals, other Arabs and non-Arabs) with follow up every 5 years (Al Thani et al., 2019b) To achieve a representative sample of the permanent population that resides in Qatar, the inclusion criteria of the QBB are: 1. To be Qatari nationals or resident in Qatar for at least 15 years and 2. To be 18 years or older. QBB is inclusive and language specification and tribes name or origin are not part of the inclusion criteria. The participants are recruited from the general public via either social media and the QBB website or through personal recommendations of family and friends.

The study covers extensive baseline sociodemographic data, clinical and behavioral phenotypic data, biological samples (i.e. blood, urine, saliva, DNA, RNA, viable cells and others), as well as clinical biomarkers and Omics data (i.e. genomics, transcriptomics, proteomics, metabolomics etc.) (Al Thani et al., 2019b). Currently the QBB has reached 44.7 % of the target population (60,000) and more than 2 million biological samples. For this study, data from 6,045 Qatari nationals participants were available from QBB population cohort. The percentage female was 56.74% and the mean age was 40 years (SD 12.7 years).

Ethics Statement

All QBB participants signed an Informed Consent Form prior to their participation; QBB study protocol ethical approval was obtained from the Hamad Medical Corporation Ethics Committee in 2011 and continued with QBB Institutional Review Board (IRB) from 2017 onwards and it is renewed on an annual basis (IRB protocol number, QF-QGP-RES-PUB-002).

Qatar Biobank sample collection

Physical and clinical measurements were collected by the QBB, in addition to biological samples (approximately 60ml of blood, 5ml of saliva, and 10ml of urine). Participants were instructed for 8 hours fasting before the visit, but due to different visit shifts samples were mostly spot specimens. Blood samples were analysed to assess 66 different biomarkers associated with disease risk factors. Haematology and blood chemistry biomarkers were analysed at Hamad General Hospital laboratories. EDTA blood samples were separated by centrifugation into plasma, buffy coat (leucocytes) and erythrocytes. All collected samples were aliquoted and stored in 3 different locations (Al Thani et al., 2019a).

DNA isolation and Quality Control

Prior to DNA isolation, each buffy coat sample was registered into the Laboratory Information Management System (LIMS) and assigned with three identifiers: i. the aliquot code, ii. a subject-specific personal number, and iii. a sample-specific serial number. Samples were received in 2D-coded FluidX tubes (Brooks Life Sciences). Upon receiving, samples were scanned on a 2D FluidX Perception Barcode Reader to check for consistency against the sample submission form. The buffy coat samples were processed for DNA isolation using the automated QIASymphony SP instrument according to Qiagen MIDI kit protocol's recommendations. The assessment of DNA quantity and quality was carried out using NanoDrop 8000 (ThermoFisher, Waltham, MA, USA), FlexStation 3 (Molecular Devices, Sunnyvale, CA, USA) and LabChip GX (Perkin Elmer, Waltham, MA, USA). The absorbance at 260 and 280 nm wavelength was measured on Nanodrop 8000 and used to check DNA purity. A fluorescence-based quantification was performed on FlexStation 3 using Quant-iT PicoGreen dsDNA Assay (ThermoFisher). Briefly, an aqueous working solution of the Quant-iT PicoGreen reagent was prepared on the day of the quantification experiment by making a 200-fold dilution of the concentrated DMSO solution in TE. TE buffer was also used for diluting DNA samples and in the assay itself. Sample measurement on FlexStation 3 was performed following the manufacturer's recommendations. DNA integrity was checked on LabChip GX. The Gel-Dye solution, DNA samples and DNA ladder were prepared according to the manufacturer's instructions; the run data was compared to the electropherogram of a typical high molecular weight ladder and assessed for quality. A genomic DNA (gDNA) quality score (GQS) was calculated for each sample. The GQS is derived from the size distribution of the gDNA and it represents the degree of degradation of a given sample, with a score of 5 corresponding to intact gDNA and a score of 0 corresponding to a

highly degraded gDNA. Figure S1 shows the GQS distribution across 50 samples assessed from phase I. The distribution shows $GQS > 3.5$.

Whole genome sequencing

Library construction and sequencing was performed at the Sidra Clinical Genomics Laboratory Sequencing Facility. After extraction of genomic DNA, sample integrity was controlled using the Genomic DNA assay on the Perkin Elmer Caliper Labchip GXII. Concentration was measured using Invitrogen Quant-iT dsDNA Assay on the FlexStation 3. Around 150ng of DNA were used for library construction with the Illumina TruSeq DNA Nano kit. Each library was indexed using the Illumina TruSeq Single Indexes. Library quality and concentration was assessed using the DNA 1k assay on a Perkin Elmer GX2. Libraries were quantified using the KAPA HiFi Library quantification kit on a Roche LightCycler 480. Flow cells were loaded at 1 sample per lane and cluster generation was performed on a cBot 1.0 or 2.0 using the HiSeq X Ten Reagent Kit v2.5. Flow cells were loaded at a cluster density between 1255 and 1412 K/mm² and sequenced on an Illumina HiSeq X instrument to a minimum average coverage of 30x.

Sequencing data processing methods

The Sidra Bioinformatics Core (SBC) developed a pipeline to perform the NGS analysis for QGP and other internal projects (Figure S2). The core also developed a framework to automate the processing of the samples. Data is received from the clinical genomic lab (CGL) in Fastq¹ format. Quality control of Fastq files is performed using FastQC(v0.11.2)², to calculate quality metrics and ensure that raw reads have good quality. Reads are then trimmed and aligned to hs37d5³ reference genome using bwa.kit (v0.7.12)⁴ and a bam⁵ file is generated. Quality control

on mapped reads (BAM files), to evaluate the coverage of each sample, is performed using Picard (v1.117) [CollectWgsMetrics]⁶. The variant calling is performed following GATK 3.4 best practices⁷: Indel realignment and base recalibration (BQSR) is performed on the initial bam then HaplotypeCaller run on each sample to generate an intermediate genomic gVCF (gVCF). Joint Genotyping is performed using all generated gVCF files at once. We first run GenomicsDB⁸ to combine the different samples by regions, then on each region, we run GenotypeGVCFs, apply SNP/Indel recalibration (VQSR) and then merge all regions. Annotation is performed using SnpEff/SnpSift⁹ (v4.3t). The following databases are used within SnpEff/SnpSift for the annotation of the multi-samples VCF file:

- dbSNP build 151
- ClinVar 2019-02-11
- dbNSFP¹⁰ v2.9
- GWAS catalog¹¹
- msigDBdb¹² v5.0

All variants are kept within the VCF file. Copy Number Variation analysis was performed using Canvas¹³ (v1.11.0) and structural variant analysis was performed using Manta¹⁴ (v0.29.6) and Delly¹⁵ (v0.7.8). Both analyses use bam file as input and were performed at the single sample level. Additionally, QGP VCF file was decomposed for multi allelic position and then normalize using vt¹⁶ (v0.5). QGP VCF file was split chromosome wise and this per chromosome VCF file was provided for further analysis as well. All pipeline references are in the supplemental data.

To identify disease-causing variants in HGMD, ClinVar and OMIM, we used VCF file annotated with phenotype/disease information from these databases. To achieve that, we applied successive filtering on the variant list using different criteria (selecting only those located in known

HGMD/OMIM gene, variants with MAF <1% in all databases, except QGP, and the variant should be within or affecting the coding region; missense, nonsense, frameshift, and splice-site variants). Among the final list, we selected those that have been previously reported and flagged as disease-causing “DM/DM?” in HGMD or “Pathogenic/Likely_pathogenic” in ClinVar.

Data Quality Control

QGP phase I study included 6,218 samples. We apply downstream quality control on the multi-sample VCF using the PLINK v2.0 tool (Chang et al., 2015). After quality control, 8 samples were removed for excess heterozygosity, 1 for low-call rates (less than 95%), 65 for gender mismatch, 87 for population outliers (individuals with more than four standard deviation (± 4 SD) away from the mean of the first two multidimensional scaling component), and 10 for identical matching. After these exclusions (N= 171), a final set of 6,045 subjects was obtained for which whole-genome sequencing was performed at a median depth of 32X (Thareja et al. (in press)).

Statistical analyses

We compared the allele counts of QGP samples to allele counts present in gnomAD exome samples for HGMD DM variants. A Fisher’s exact test was used to calculate variations that were significantly overrepresented in the QGP samples (due to founder effect) and corrected for multiple testing using the Bonferroni method.

Hail genomic processing tool

Data preprocessing and analysis was performed using Hail 0.2. allele count, allele number, allele frequency, homozygous count calculation for each subpopulation was performed simultaneously using python scripts written using hail framework. Quality analysis for variant calls and

individual sample were performed using `variant_qc` and `sample_qc` functions respectively.

Sample level statistics for each sample was generated using the Hail.

QGP variant browser

QGP variant browser provides a mechanism for the researchers to be able to search, filter and browse the QGP genomic variants data. This web-based browser supports fast database query response time for searching through more than 88 million records with search and filter functionality on the QGP gene variants and its attributes (e.g. allele frequency, homozygosity etc.).

Results

Genetic variability of the Qatari population

We have identified a total of 88,191,239 variants, which includes 74,991,446 SNVs (74,040,559 bi-allelic SNVs) with 939,405 multi-allelic sites and 13,199,792 INDELS (8,389,562 bi-allelic INDELS) with 2,018,185 multi-allelic sites/microsatellites (Figure 2a-c and Figure S3). Importantly, twenty-eight percent (28%) of the total variants (24,620,313) were novel and not previously reported in dbSNP build 151 or other population databases (gnomAD, 1000 Genomes, and Greater Middle East (GME)) (Figure 2b; Figure S4a-b and 5). Each individual genome presented a median of 3.4 million SNVs and 63,755 novel variants. We estimated the transition to transversion (ti/tv) ratio of 2.05 and heterozygotes to non-ref homozygote (Het/Hom) ratio of 1.85, which is consistent with previous WGS studies (Auton et al., 2015). We found 23 million variants present as singletons which are less when compared to the number of variants falling under the minor allele frequency (MAF) spectrum of <0.1% (2-12 alleles) which

should be around 34 million variants (Figure 2c and Table S1). While considering the novel variants, singletons (45%) being slightly higher than the variants that fall in the category of 2-12 alleles (42%) and only 13% of the novel variants exceed the $MAF > 0.1\%$. Half of the singletons present in QGP were already reported in dbSNP and, each individual carried a median of 1,336 singletons (Figure 2d and Figure S6).

To evaluate the impact and scale of disease-causing variants in our population, we annotated the variant list with disease/phenotype information from HGMD, ClinVar and OMIM databases. In total, we found 4,254 disease-causing mutations (DM), which includes 3,970 SNVs and 284 INDELS (Figure S7a). These variants are located across 1,672 genes that are linked to phenotypes with different modes of inheritance (678 follow autosomal recessive (AR); 315 autosomal dominant (AD); 526 both AR and AD; and 50 X-linked inheritance) (Figure 2e). The vast majority (97%) of these DM variants are rare with $MAF < 1\%$, and among them 30% observed as singletons (Figure S7b). Each individual in the QGP dataset carries a median of 21 DM variants (range of 8-37) (Figure 2f and Figure S7c), slightly less than what have been previously reported (25 DMs/individual in the UK10k (Xue et al., 2012) and 29 DMs/individual in the Uganda genome studies (Gurdasani et al., 2019)). Each individual also carries in the homozygous state a median of 5 DM variants (range of 1-11) compared to 3 homozygous DMs/individual in the Uganda genome and 3–24 homozygous variants in the 1000 Genome project (Auton et al., 2015; Gurdasani et al., 2019). Our data shows that approximately 900 protein-coding genes have at least 1 DM mutation and 26 genes present 15 or more DM mutations (Figure S7d). When QGP data is classified according to ClinVar information (version February 11th 2019), we found that 1,449 variants are classified as “pathogenic” or “likely pathogenic” (Figure S7e). Further classification considering both HGMD and ClinVar, revealed

that 1,011 variants were marked as DM and “pathogenic or likely pathogenic” (Figure 2g), with 160 variants unique to the Qatari population. Interestingly, only a subset of 14 variants, among the 1,011 variants, are shared between the QGP samples and data from Greater Middle East (GME) Variome Project (Scott et al., 2016b) (Table 1). There are also 36 variants which confer protection against several diseases including malaria, obesity, and heart disease (Table S2).

We found some rare pathogenic variants present in Qatari population with high minor allele frequencies due to the founder effect. Some of the examples include variant in the *MPL* gene [MIM: 604498] (rs750046020), previously associated with thrombocytosis, occurs at a MAF of 0.009, and similarly, variants in the genes *CBS* [MIM:236200] (rs398123151) and *KRT5* [MIM: 148040](rs267607448) associated with homocystinuria and Epidermolysis Bullosa, respectively, are observed at a MAF of 0.007”.

Genetic Ancestry and Diversity of the Qatari population

To capture the genetic diversity of the Qatari population and understand its relationship with the world’s populations in both modern and ancient times, we identified five major ancestries: General Arabs (QGP_GAR, 38%), Peninsular Arabs (QGP_PAR, 17%), Arabs of Western Eurasia and Persia (QGP_WEP, 22%), South Asians (QGP_SAS, 1%), Africans (QGP_AFR, 3%) and Admixed (QGP_ADM, 19%) (Razali et al., in press) (Figure S8). We also characterized a group of Peninsular Arabs forming a unique cluster within known descendants originating from the historical homeland of ancient Arab tribes. Analysis of Mitochondrial DNA (mtDNA) and Chromosome Y (Chr Y) in the dataset has enriched the poorly characterized landscape of haplogroups in Arab and Middle East populations in general. Notably, J1a2b a Chr Y haplogroup seen previously in Yemen, has been observed in 1,419 males, which is the largest set of individuals ever sequenced within this haplogroup. We discovered 103 novel Y-Chr SNPs in

these individuals, which aided the expansion of this haplogroup to 29 novel sub-haplogroups. Using this unique dataset, we built a panel for genotype imputation for Arabs and Middle Eastern ancestries which shows an improved imputation score for rare and common allele frequencies variants (Razali et al., in press).

We next characterized the spectrum of genetic variability based on the fine-scale population structure observed in the Qatari population. This analysis highlighted that 70% of the novel variants are cluster-specific, 5% are found in all sub-clusters, and the remaining 25% are shared between one or more sub-clusters (Figure S9a). Similarly, we found that about half (2,139) of the DM variants are cluster-specific and only 68 out of 4,254 DM variants were present in all sub-clusters (Figure S9b). Furthermore, individuals in the QGP_AFR sub-cluster have the highest heterozygotes to non-ref homozygote (Het/Hom) ratio, whereas the ratio was found to be lowest for the QGP_PAR cluster. This reflects the high homozygosity and high consanguinity present within the individuals of this cluster (Figure S9c). Similarly, the median number of singletons is lower for PAR cluster compared to other sub-clusters reflects the closely related individual present in this cluster (Table 2).

Furthermore, runs of homozygosity (ROH) analysis of the QGP done by Razali et al 2021 (Razali et al., in Press), identified per population ROH boundary for short, medium and long ROH. We observed that Peninsular Arabs (PAR) have the lowest median for short ROH after African-based populations. In addition, PAR has the highest median for long ROH, indicating recent consanguinity events. When we analyzed the relationship between genes and the ROH regions, we observed that there are more OMIM genes in ROH regions compared to non-OMIM genes regardless of the ROH classes. PAR was shown to have significantly more OMIM genes compared to the other QGP and 1KG populations.

Burden of Pathogenic Variation

We then focused on the burden of pathogenic variants of recessively inherited disorders in the Qatari population. We found the most common recessive alleles are those linked to structural deformities and developmental disorders, consistent with the fact that such recessive traits prevail in societies where endogamy and consanguinity is practiced (Table S3). However, some of these identified alleles are too common to be classified as pathogenic variants (rs201818754, rs373804633, rs199768740, and rs80358230) as their frequencies in PAR subpopulation exceeding 4%, far more than the associated disease prevalence.

A notable example of an autosomal recessive disorder is Woodhouse-Sakati syndrome [WSS (MIM:241080)], a disease characterized by hypogonadism and hair thinning that often progresses to alopecia totalis. Of the less than 100 individuals reported globally with the disease, 30 are from Middle Eastern families (Bohlega & Alkuraya, 1993). WSS is caused by biallelic pathogenic variants in the [DCAF17 (MIM: 612515)] (previously known as C2orf37) gene. We identified NM_025000.4 (DCAF17): c.436delC (p.Ala147fs) as the sole pathogenic variant of this gene in 88 individuals, in heterozygous state (MAF = 0.007) (Supplementary data). Although all heterozygous individuals were found to be clinically asymptomatic, the alternate allele in these individuals is associated with the decreased levels of Insulin (Pvalue = 2.9E-02; β = -0.225; Figure S10) which could explain diabetes mellitus being one of the characteristic clinical phenotypes in WSS. We also found that c.436delC is enriched (fisher exact test P=7.57E-34; OR=18.45) in one of the founder populations, QGP_PAR subcluster, this is consistent with a previous report that identified DCAF17:c.436delC (rs797045038) as a founder variant in the Qatari population (Ben-Omran et al., 2011). This variant has also been reported in the Kingdom of Saudi Arabia (Alazami et al., 2008), which has a large number of tribes sharing

common and similar carrier frequency with Qatar's native population. Hamad Medical Corporation (HMC) is hosting the national molecular diagnostic laboratories of Qatar, and has identified to date 34 WSS patients and 64 heterozygous carriers. Data from both QGP and HMC laboratories indicates that the carrier frequency for WSS in the Qatari population is approx. 1 in 42 individuals (2.5%) with MAF of 1.25%, which is the highest reported in the world. Remarkably, the carrier frequency of c.436delC (p.Ala147fs) is 7x higher in Qatar than in the same tribe living in neighboring Saudi Arabia and has not yet been reported in population frequency databases, such as gnomAD and 1000 genomes or the 100K Genomes Project that includes patients with rare genetic diseases (Turnbull et al., 2018).

Insights into the genetics of quantitative traits

To gain insights into the genetic architecture of health and disease-related quantitative traits, we performed the first genome wide association studies of a list of 45 quantitative traits in 6,045 individuals from the Qatari population (Thareja et al., 2021). Several important findings of this comprehensive study include replication of multiple associations reported in Caucasian and Asian GWASs; uncovering differences in allele frequencies and LD patterns for replicated loci; and discovery of novel genetic associations mostly with variants common in the QGP but rare in other populations. These findings argue for larger GWAS studies from the region to accurately derive polygenic risk scores optimized for Middle Eastern populations for improved application in precision medicine.

Discussion

Here we characterized a broad spectrum of genetic variation in the Qatari population, in total over 88 million variants (1.86 % of novel variants per individual genome and 24.6 M novel variants in the whole dataset). This large-scale study allowed us to identify five non-admixed

subgroups in QGP (n=6,045) compared to three in the previous study Fakhro et al. 2016 (n=1,005) (Fakhro et al., 2016). We found a larger number of DM variants carried per individual which could be explained by incomplete penetrance, or the individual might carry them in a heterozygous state (Francioli et al., 2014; Xue et al., 2012). We described the distribution of genetic variation across the sub-clusters and found the majority of the novel variants to be cluster-specific. This data support records of high consanguinity and founder effect but also identify a previously unstudied component of the Middle Eastern population. Based on these sequencing results of 6,045 individuals we have recently reported a total of 60 pathogenic and likely pathogenic in 25 ACMG genes in 141 unique individuals (Elfatih et al., 2021) and several other efforts are currently under way to build the catalogues of predicted loss-of-function variants and mendelian disorders mutations and to characterize the pharmacogenomic and the cancer landscapes of the Qatari population. Furthermore, using a combination of whole genomes and exome sequence data and clinical reports, we developed a microarray with Qatari-specific pathogenic variants that could be used to rapidly, accurately and at low cost, screen the Qatari population for pathogenic variants of newborns, premarital couples and patients presenting to the clinic (Rodriguez-Flores, in press).

Previous genetic studies in the Middle East region have assessed the genomic variations linked to health and diseases mostly limited to whole exome sequencing on relatively small sample size (AlSafar et al., 2019; Fattahi et al., 2019; John et al., 2018; Monies et al., 2019; Scott et al., 2016b). Our QGP data have a key advantage over these studies since we are performing large-scale population sequencing using a whole genome approach. Although our work provided various insights into the genomic of the Middle East, we should address one limitation of our

approach is that we are including only Qatari nationals in the first phase. To overcome this limitation, we are including long term residents in our next freezes.

In conclusion, this first phase of the QGP constitutes the largest comprehensive analysis of whole genomes representative of tens of millions of Arabian Peninsula and Middle East inhabitants. Such genetic information is largely lacking in global databases (Easteal et al., 2020). Our next phases will focus on specific diseases relevant to the Qatari population's health burden - e.g. cancer, diabetes and rare diseases - while accelerating the ability to use the genome sequencing data into clinical implementation. We anticipate our data will represent a valuable resource to advance genetic studies in the Arab and neighbouring Middle Eastern populations and will significantly boost the current efforts to improve our understanding of global patterns of human variations, human history and genetic contributions to health and diseases in diverse populations (C. N. Rotimi & Adeyemo, 2021).

Supplemental Information

Supplemental information includes ten figures (S1–S10), three tables (S1–S3), information about *DCAF17* founder mutation, and references for the pipelines.

Declaration of Interests

The authors declare no competing interests.

The Qatar Genome Program Research Consortium

Qatar Genome Project Management: Said I. Ismail¹, Wadha Al-Muftah¹, Radja Badji¹, Hamdi Mbarek¹, Dima Darwish¹, Tasnim Fadl¹, Heba Yasin¹, Maryem Ennaifar¹, Rania Abdellatif¹, Fatima Alkuwari¹, Muhammad Alvi¹, Yasser Al-Sarraj¹, Chadi Saad¹, Asmaa Althani^{1,16}

Biobank and Sample Preparation: Eleni Fethnou², Fatima Qafoud², Eiman Alkhayat², Nahla Afifi²

Sequencing and Genotyping group: Sara Tomei³, Wei Liu³, Stephan Lorenz³

Applied Bioinformatics Core: Najeeb Syed⁴, Hakeem Almabrazi⁴, Fazulur Rehaman Vempalli⁴, Ramzi Temanni⁴

Data Management, Advanced Applications and Computing Infrastructure groups: Tariq Abu Saqri⁵, Mohammedhusen Khatib⁵, Mehshad Hamza⁵, Tariq Abu Zaid⁵, Ahmed El Khouly⁵, Tushar Pathare⁵, Shafeeq Poolat⁵, Shafqat Baig⁵, Anwar Haque⁵, Mohamed Jama⁵, Rashid Al-Ali⁵

Genetic Variability group: Geethanjali Devadoss Gandhi^{6,8}, Senthil Selvaraj⁶, Najeeb Syed⁴, Xavier Estivill⁶, Hamdi Mbarek¹

Population Structure and Genome Reference group: Rozaimi Mohamad Razali⁶, Juan Rodriguez-Flores¹⁷, Elbay Aliyev⁶, Haroon Naeem⁶, Waleed Aamer⁶, Andrew Clark¹⁸, Khalid Fakhro⁶, Younes Mokrab⁶

GWAS group: Gaurav Thareja^{7*}, Yasser Al-Sarraj^{*1,8}, Aziz Belkadi⁷, Maryam Almotawa⁹, Karsten Suhre⁷⁺, Omar Albagha^{+8,15} (* equally contributed, + jointly supervised)

Mendelian Disorders group: Waleed Aamer⁶, Alya Al-Kurbi⁶, Aljazi Al-Maraghi⁶, Geethanjali Devadoss Gandhi^{6,8}, Najeeb Syed⁴, Khalid Fakhro⁶

Loss of Function group: Fatemeh Abbaszadeh^{10*}, Ikhlaq Ahmed^{5*}, Najeeb Syed⁴, Mohammad Abuhaliqa¹⁰, Rashid Al Ali⁵, Khalid Fakhro⁶, Zafar Nawaz¹⁰, Ajayeb Al Nabet Al Marri¹⁰, Xavier Estivill⁶, Puthen V. Jithesh⁸, Ramin Badii¹⁰ (* equally contributed)

Consortium Lead Principal Investigators (in alphabetical order): Omar Albagha^{8,15}, Souhaila Al-Khodor¹¹, Mashael Alshafai¹², Ramin Badii¹⁰, Lotfi Chouchane¹³, Xavier Estivill⁶, Khalid Fakhro⁶, Hamdi Mbarek¹, Younes Mokrab⁶, Puthen V. Jithesh⁸, Karsten Suhre⁷, Zohreh Tatari¹⁴

Affiliations

1. Qatar Genome Program, Qatar Foundation Research, Development and Innovation, Qatar Foundation, Doha, Qatar.
2. Qatar Biobank for Medical Research, Qatar Foundation, Building 317, Hamad Medical City, Doha, Qatar.
3. Sidra Medicine, Integrated Genomics Services, Out-Patient Clinic, Doha, Qatar.
4. Sidra Medicine, Applied Bioinformatics Core - Integrated Genomics Services - Research Branch, Doha, Qatar.
5. Sidra Medicine, Biomedical Informatics – Research Branch, Doha, Qatar.
6. Sidra Medicine, Human Genetics Department, Doha, Qatar.

7. Bioinformatics Core, Weill Cornell Medicine-Qatar, Education City, Doha, Qatar.
8. College of Health and Life Sciences, Hamad Bin Khalifa University, Education City, Doha, Qatar.
9. Qatar Biomedical Research Institute (QBRI), Hamad Bin Khalifa University, Doha, Qatar
10. Molecular Genetics Laboratory, Hamad Medical Corporation, Doha, Qatar.
11. Sidra Medicine, Maternal and Child Health Program, Doha Qatar.
12. College of Health Sciences, Qatar University, Doha, Qatar.
13. Departments of Genetic Medicine, Microbiology and Immunology, Weill Cornell Medicine-Qatar, Doha, Qatar.
14. Sidra Medicine, Clinical Research Center, Doha, Qatar.
15. Center of Genomic and Experimental Medicine, University of Edinburgh, Edinburgh, UK.
16. Biomedical Research Center, Qatar University, Doha, Qatar
17. Department of Genetic Medicine, Weill Cornell Medicine, New York, U.S.A.
18. Department of Molecular Biology and Genetics, Cornell University, New York, U.S.A.

Acknowledgments

The Qatar Genome Program (QGP) and Qatar Biobank (QBB) are both Research and Development entities within Qatar Foundation for Education, Science and Community Development. We are thankful for everyone who contributed to this endeavor including the QGP and QBB team members, in addition to our partners at Hamad Medical Corporation (HMC),

Sidra Medicine and other national stakeholders. We would like to especially thank all participants in this study for their continuous support.

References Pipelines

QGP: <https://qatargenome.org.qa>

QBB: <https://www.qatarbiobank.org.qa>

Sidra Medicine: <https://www.sidra.org>

dbSNP build 151: <http://www.ncbi.nlm.nih.gov/SNP/>

HGMD: <http://www.hgmd.cf.ac.uk/ac/index.php>

ClinVar: <https://www.ncbi.nlm.nih.gov/clinvar/>

gnomAD: <https://gnomad.broadinstitute.org/>

GME Variome: <http://igm.ucsd.edu/gme/> Genomics England:

<https://www.genomicsengland.co.uk/>

Hail-0.2.13: <https://github.com/hail-is/hail/releases/tag/0.2.13>

OMIM: <https://www.omim.org/>

Data and Code Availability

The informed consent given by the study participants does not cover posting of participant level phenotype and genotype data of Qatar Biobank/Qatar Genome Project in public databases. However, access to QBB/QGP data can be obtained through an established ISO-certified process

by submitting a project request at <https://www.qatarbiobank.org.qa/research/how-apply> which is subject to approval by the QBB IRB committee.

References:

- About Tayoun, A. N., & Rehm, H. L. (2020). Genetic variation in the Middle East—an opportunity to advance the human genetics field. *Genome Medicine*, *12*(1).
<https://doi.org/10.1186/s13073-020-00821-7>
- Al-Gazali, L., Hamamy, H., & Al-Arrayad, S. (2006). Genetic disorders in the Arab world. In *British Medical Journal* (Vol. 333, Issue 7573, pp. 831–834). BMJ.
<https://doi.org/10.1136/bmj.38982.704931.AE>
- Al Thani, A., Fthenou, E., Paparrodopoulos, S., Al Marri, A., Shi, Z., Qafoud, F., & Afifi, N. (2019a). Qatar Biobank Cohort Study: Study Design and First Results. *American Journal of Epidemiology*, *188*(8), 1420–1433. <https://doi.org/10.1093/aje/kwz084>
- Al Thani, A., Fthenou, E., Paparrodopoulos, S., Al Marri, A., Shi, Z., Qafoud, F., & Afifi, N. (2019b). Qatar Biobank Cohort Study: Study Design and First Results. *American Journal of Epidemiology*, *188*(8), 1420–1433. <https://doi.org/10.1093/aje/kwz084>
- Alazami, A. M., Al-Saif, A., Al-Semari, A., Bohlega, S., Zlitni, S., Alzahrani, F., Bavi, P., Kaya, N., Colak, D., Khalak, H., Baltus, A., Peterlin, B., Danda, S., Bhatia, K. P., Schneider, S. A., Sakati, N., Walsh, C. A., Al-Mohanna, F., Meyer, B., & Alkuraya, F. S. (2008). Mutations in C2orf37, Encoding a Nucleolar Protein, Cause Hypogonadism, Alopecia, Diabetes Mellitus, Mental Retardation, and Extrapyrmidal Syndrome. *American Journal of Human Genetics*, *83*(6), 684–691. <https://doi.org/10.1016/j.ajhg.2008.10.018>
- AlSafar, H. S., Al-Ali, M., Elbait, G. D., Al-Maini, M. H., Ruta, D., Peramo, B., Henschel, A., &

- Tay, G. K. (2019). Introducing the first whole genomes of nationals from the United Arab Emirates. *Scientific Reports*, 9(1). <https://doi.org/10.1038/s41598-019-50876-9>
- Anwar, W. A., Khyatti, M., & Hemminki, K. (2014). Consanguinity and genetic diseases in North Africa and immigrants to Europe. *European Journal of Public Health*, 24(SUPPL.1), 57–63. <https://doi.org/10.1093/eurpub/cku104>
- Arauna, L. R., Mendoza-Revilla, J., Mas-Sandoval, A., Izaabel, H., Bekada, A., Benhamamouch, S., Fadhlouli-Zid, K., Zalloua, P., Hellenthal, G., & Comas, D. (2017). Recent Historical Migrations Have Shaped the Gene Pool of Arabs and Berbers in North Africa. *Molecular Biology and Evolution*, 34(2), 318–329. <https://doi.org/10.1093/molbev/msw218>
- Auton, A., Abecasis, G. R., Altshuler, D. M., Durbin, R. M., Bentley, D. R., Chakravarti, A., Clark, A. G., Donnelly, P., Eichler, E. E., Flicek, P., Gabriel, S. B., Gibbs, R. A., Green, E. D., Hurles, M. E., Knoppers, B. M., Korbel, J. O., Lander, E. S., Lee, C., Lehrach, H., ... Schloss, J. A. (2015). A global reference for human genetic variation. In *Nature* (Vol. 526, Issue 7571, pp. 68–74). Nature Publishing Group. <https://doi.org/10.1038/nature15393>
- Ben-Omran, T., Ali, R., Almureikhi, M., Alameer, S., Al-Saffar, M., Walsh, C. A., Felie, J. M., & Teebi, A. (2011). Phenotypic heterogeneity in Woodhouse-Sakati syndrome: Two new families with a mutation in the C2orf37 gene. *American Journal of Medical Genetics, Part A*, 155(11), 2647–2653. <https://doi.org/10.1002/ajmg.a.34219>
- Bentley, A. R., Callier, S., & Rotimi, C. N. (2017). Diversity and inclusion in genomic research: why the uneven progress? *Journal of Community Genetics*, 8(4), 255–266. <https://doi.org/10.1007/s12687-017-0316-6>
- Bohlega, S. A., & Alkuraya, F. S. (1993). Woodhouse-Sakati Syndrome. In *GeneReviews*®. <http://www.ncbi.nlm.nih.gov/pubmed/27489925>

- Boomsma, D. I., Wijmenga, C., Slagboom, E. P., Swertz, M. A., Karssen, L. C., Abdellaoui, A., Ye, K., Guryev, V., Vermaat, M., Van Dijk, F., Francioli, L. C., Hottenga, J. J., Laros, J. F. J., Li, Q., Li, Y., Cao, H., Chen, R., Du, Y., Li, N., ... Van Duijn, C. M. (2014). The Genome of the Netherlands: Design, and project goals. *European Journal of Human Genetics*, 22(2), 221–227. <https://doi.org/10.1038/ejhg.2013.118>
- Chang, C. C., Chow, C. C., Tellier, L. C. A. M., Vattikuti, S., Purcell, S. M., & Lee, J. J. (2015). Second-generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience*, 4(1). <https://doi.org/10.1186/s13742-015-0047-8>
- Chiang, C. W. K., Marcus, J. H., Sidore, C., Biddanda, A., Al-Asadi, H., Zoledziewska, M., Pitzalis, M., Busonero, F., Maschio, A., Pistis, G., Steri, M., Angius, A., Lohmueller, K. E., Abecasis, G. R., Schlessinger, D., Cucca, F., & Novembre, J. (2018). Genomic history of the Sardinian population. *Nature Genetics*, 50(10), 1426–1434. <https://doi.org/10.1038/s41588-018-0215-8>
- Collins, F. S., & Varmus, H. (2015). A New Initiative on Precision Medicine. *New England Journal of Medicine*, 372(9), 793–795. <https://doi.org/10.1056/nejmp1500523>
- Easteal, S., Arkell, R. M., Balboa, R. F., Bellingham, S. A., Brown, A. D., Calma, T., Cook, M. C., Davis, M., Dawkins, H. J. S., Dinger, M. E., Dobbie, M. S., Farlow, A., Gwynne, K. G., Hermes, A., Hoy, W. E., Jenkins, M. R., Jiang, S. H., Kaplan, W., Leslie, S., ... Baynam, G. (2020). Equitable Expanded Carrier Screening Needs Indigenous Clinical and Population Genomic Data. *American Journal of Human Genetics*, 107(2), 175–182. <https://doi.org/10.1016/j.ajhg.2020.06.005>
- Elfatih, A., Mifsud, B., Syed, N., Badii, R., Mbarek, H., Abbaszadeh, F., & Estivill, X. (2021). Actionable genomic variants in 6045 participants from the Qatar Genome Program. *Human*

Mutation. <https://doi.org/10.1002/humu.24278>

Fakhro, K. A., Staudt, M. R., Ramstetter, M. D., Robay, A., Malek, J. A., Badii, R., Al-Marri, A. A. N., Khalil, C. A., Al-Shakaki, A., Chidiac, O., Stadler, D., Zirie, M., Jayyousi, A., Salit, J., Mezey, J. G., Crystal, R. G., & Rodriguez-Flores, J. L. (2016). The Qatar genome: A population-specific tool for precision medicine in the Middle East. *Human Genome Variation*, 3. <https://doi.org/10.1038/hgv.2016.16>

Fattahi, Z., Beheshtian, M., Mohseni, M., Poustchi, H., Sellars, E., Nezhadi, S. H., Amini, A., Arzhangi, S., Jalalvand, K., Jamali, P., Mohammadi, Z., Davarnia, B., Nikuei, P., Oladnabi, M., Mohammadzadeh, A., Zohrehvand, E., Nejatizadeh, A., Shekari, M., Bagherzadeh, M., ... Najmabadi, H. (2019). Iranome: A catalog of genomic variations in the Iranian population. *Human Mutation*, 40(11), 1968–1984. <https://doi.org/10.1002/humu.23880>

Francioli, L. C., Menelaou, A., Pulit, S. L., Van Dijk, F., Palamara, P. F., Elbers, C. C., Neerincx, P. B. T., Ye, K., Guryev, V., Kloosterman, W. P., Deelen, P., Abdellaoui, A., Van Leeuwen, E. M., Van Oven, M., Vermaat, M., Li, M., Laros, J. F. J., Karssen, L. C., Kanterakis, A., ... Wijmenga, C. (2014). Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nature Genetics*, 46(8), 818–825. <https://doi.org/10.1038/ng.3021>

Gudbjartsson, D. F., Helgason, H., Gudjonsson, S. A., Zink, F., Oddson, A., Gylfason, A., Besenbacher, S., Magnusson, G., Halldorsson, B. V, Hjartarson, E., Sigurdsson, G. T., Stacey, S. N., Frigge, M. L., Holm, H., Saemundsdottir, J., Helgadóttir, H. T., Johannsdóttir, H., Sigfusson, G., Thorgeirsson, G., ... Stefansson, K. (2015). Large-scale whole-genome sequencing of the Icelandic population. *Nature Genetics*, 47(5), 435–444. <https://doi.org/10.1038/ng.3247>

- Gudbjartsson, D. F., Sulem, P., Helgason, H., Gylfason, A., Gudjonsson, S. A., Zink, F., Oddson, A., Magnusson, G., Halldorsson, B. V, Hjartarson, E., Sigurdsson, G. T., Kong, A., Helgason, A., Masson, G., Magnusson, O. T., Thorsteinsdottir, U., & Stefansson, K. (2015). Sequence variants from whole genome sequencing a large group of Icelanders. *Scientific Data*, 2, 150011. <https://doi.org/10.1038/sdata.2015.11>
- Gurdasani, D., Carstensen, T., Fatumo, S., Chen, G., Franklin, C. S., Prado-Martinez, J., Bouman, H., Abascal, F., Haber, M., Tachmazidou, I., Mathieson, I., Ekoru, K., DeGorter, M. K., Nsubuga, R. N., Finan, C., Wheeler, E., Chen, L., Cooper, D. N., Schiffels, S., ... Sandhu, M. S. (2019). Uganda Genome Resource Enables Insights into Population History and Genomic Discovery in Africa. *Cell*, 179(4), 984-1002.e36. <https://doi.org/10.1016/j.cell.2019.10.004>
- Harkness, G., & Khaled, R. (2014). Modern traditionalism: Consanguineous marriage in Qatar. *Journal of Marriage and Family*, 76(3), 587–603. <https://doi.org/10.1111/jomf.12106>
- Hellenthal, G., Busby, G. B. J., Band, G., Wilson, J. F., Capelli, C., Falush, D., & Myers, S. (2014). A genetic atlas of human admixture history. *Science*, 343(6172), 747–751. <https://doi.org/10.1126/science.1243518>
- Hunter-Zinck, H., Musharoff, S., Salit, J., Al-Ali, K. A., Chouchane, L., Gohar, A., Matthews, R., Butler, M. W., Fuller, J., Hackett, N. R., Crystal, R. G., & Clark, A. G. (2010). Population genetic structure of the people of Qatar. *American Journal of Human Genetics*, 87(1), 17–25. <https://doi.org/10.1016/j.ajhg.2010.05.018>
- John, S. E., Antony, D., Eaaswarkhanth, M., Hebbar, P., Channanath, A. M., Thomas, D., Devarajan, S., Tuomilehto, J., Al-Mulla, F., Alsmadi, O., & Thanaraj, T. A. (2018). Assessment of coding region variants in Kuwaiti population: implications for medical

genetics and population genomics. *Scientific Reports*, 8(1). <https://doi.org/10.1038/s41598-018-34815-8>

Karczewski, K. J., Francioli, L. C., Tiao, G., Cummings, B. B., Alföldi, J., Wang, Q., Collins, R. L., Laricchia, K. M., Ganna, A., Birnbaum, D. P., Gauthier, L. D., Brand, H., Solomonson, M., Watts, N. A., Rhodes, D., Singer-Berk, M., England, E. M., Seaby, E. G., Kosmicki, J. A., ... MacArthur, D. G. (2020). The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*, 581(7809), 434–443. <https://doi.org/10.1038/s41586-020-2308-7>

Landry, L. G., Ali, N., Williams, D. R., Rehm, H. L., & Bonham, V. L. (2018). Lack of diversity in genomic databases is a barrier to translating precision medicine research into practice. *Health Affairs*, 37(5), 780–785. <https://doi.org/10.1377/hlthaff.2017.1595>

Manolio, T. A., Bult, C. J., Chisholm, R. L., Deverka, P. A., Ginsburg, G. S., Jarvik, G. P., McLeod, H. L., Mensah, G. A., Relling, M. V., Roden, D. M., Rowley, R., Tamburro, C., Williams, M. S., & Green, E. D. (2019). Genomic Medicine Year in Review: 2019. *American Journal of Human Genetics*, 105(6), 1072–1075. <https://doi.org/10.1016/j.ajhg.2019.11.006>

Mills, M. C., & Rahal, C. (2019). A scientometric review of genome-wide association studies. In *Communications Biology* (Vol. 2, Issue 1). Nature Research. <https://doi.org/10.1038/s42003-018-0261-x>

Monies, D., Abouelhoda, M., Assoum, M., Moghrabi, N., Rafiullah, R., Almontashiri, N., Alowain, M., Alzaidan, H., Alsayed, M., Subhani, S., Cupler, E., Faden, M., Alhashem, A., Qari, A., Chedrawi, A., Aldhalaan, H., Kurdi, W., Khan, S., Rahbeeni, Z., ... Alkuraya, F. S. (2019). Erratum: Lessons Learned from Large-Scale, First-Tier Clinical Exome

- Sequencing in a Highly Consanguineous Population (The American Journal of Human Genetics (2019) 104(6) (1182–1201), (S0002929719301594), (10.1016/j.ajhg.2019.04.011)). In *American Journal of Human Genetics* (Vol. 105, Issue 4, p. 879). Cell Press. <https://doi.org/10.1016/j.ajhg.2019.09.019>
- Naslavsky, M. S., Scliar, M. O., Yamamoto, G. L., Wang, J. Y. T., Zverinova, S., Karp, T., Nunes, K., Ceroni, J. R. M., de Carvalho, D. L., da Silva Simões, C. E., Bozoklian, D., Nonaka, R., Silva, N. dos S. B., Souza, A. da S., Andrade, H. de S., Passos, M. R. S., Castro, C. F. B., Mendes-Junior, C. T., Mercuri, R. L. V., ... Zatz, M. (2020). Whole-genome sequencing of 1,171 elderly admixed individuals from the largest Latin American metropolis (São Paulo, Brazil). In *bioRxiv* (Vol. 10, p. 24). bioRxiv. <https://doi.org/10.1101/2020.09.15.298026>
- Okada, Y., Momozawa, Y., Sakaue, S., Kanai, M., Ishigaki, K., Akiyama, M., Kishikawa, T., Arai, Y., Sasaki, T., Kosaki, K., Suematsu, M., Matsuda, K., Yamamoto, K., Kubo, M., Hirose, N., & Kamatani, Y. (2018). Deep whole-genome sequencing reveals recent selection signatures linked to evolution and disease risk of Japanese. *Nature Communications*, 9(1), 1631. <https://doi.org/10.1038/s41467-018-03274-0>
- Rodriguez-Flores, J. L., Fakhro, K., Agosto-Perez, F., Ramstetter, M. D., Arbiza, L., Vincent, T. L., Robay, A., Malek, J. A., Suhre, K., Chouchane, L., Badii, R., Al-Marri, A. A. N., Khalil, C. A., Zirie, M., Jayyousi, A., Salit, J., Keinan, A., Clark, A. G., Crystal, R. G., & Mezey, J. G. (2016). Indigenous Arabs are descendants of the earliest split from ancient Eurasian populations. *Genome Research*, 26(2), 151–162. <https://doi.org/10.1101/gr.191478.115>
- Rodriguez-Flores, J. L., Fakhro, K., Hackett, N. R., Salit, J., Fuller, J., Agosto-Perez, F., Gharbiah, M., Malek, J. A., Zirie, M., Jayyousi, A., Badii, R., Al-Nabet Al-Marri, A.,

- Chouchane, L., Stadler, D. J., Mezey, J. G., & Crystal, R. G. (2014). Exome Sequencing Identifies Potential Risk Variants for Mendelian Disorders at High Prevalence in Qatar. *Human Mutation*, 35(1), 105–116. <https://doi.org/10.1002/humu.22460>
- Rotimi, C., Abayomi, A., Abimiku, A., Adabayeri, V. M., Adebamowo, C., Adebisi, E., Ademola, A. D., Adeyemo, A., Adu, D., Affolabi, D., Agongo, G., Ajayi, S., Akarolo-Anthony, S., Akinyemi, R., Akpalu, A., Alberts, M., Alonso Betancourt, O., Alzohairy, A. M., Ameni, G., ... Zar, H. (2014). Research capacity. Enabling the genomic revolution in Africa. In *Science* (Vol. 344, Issue 6190, pp. 1346–1348). American Association for the Advancement of Science. <https://doi.org/10.1126/science.1251546>
- Rotimi, C. N., & Adeyemo, A. A. (2021). From one human genome to a complex tapestry of ancestry. *Nature*, 590(7845), 220–221. <https://doi.org/10.1038/d41586-021-00237-2>
- Scott, E. M., Halees, A., Itan, Y., Spencer, E. G., He, Y., Azab, M. A., Gabriel, S. B., Belkadi, A., Boisson, B., Abel, L., Clark, A. G., Rahim, S. A., Abdel-Hadi, S., Abdel-Salam, G., Abdel-Salam, E., Abdou, M., Abhytankar, A., Adimi, P., Ahmad, J., ... Zhang, S. Y. (2016a). Characterization of greater middle eastern genetic variation for enhanced disease gene discovery. In *Nature Genetics* (Vol. 48, Issue 9, pp. 1071–1079). Nature Research. <https://doi.org/10.1038/ng.3592>
- Scott, E. M., Halees, A., Itan, Y., Spencer, E. G., He, Y., Azab, M. A., Gabriel, S. B., Belkadi, A., Boisson, B., Abel, L., Clark, A. G., Rahim, S. A., Abdel-Hadi, S., Abdel-Salam, G., Abdel-Salam, E., Abdou, M., Abhytankar, A., Adimi, P., Ahmad, J., ... Zhang, S. Y. (2016b). Characterization of greater middle eastern genetic variation for enhanced disease gene discovery. In *Nature Genetics* (Vol. 48, Issue 9, pp. 1071–1079). Nature Research. <https://doi.org/10.1038/ng.3592>

- Stark, Z., Dolman, L., Manolio, T. A., Ozenberger, B., Hill, S. L., Caulfield, M. J., Levy, Y., Glazer, D., Wilson, J., Lawler, M., Boughtwood, T., Braithwaite, J., Goodhand, P., Birney, E., & North, K. N. (2019). Integrating Genomics into Healthcare: A Global Responsibility. In *American Journal of Human Genetics* (Vol. 104, Issue 1, pp. 13–20). Cell Press.
<https://doi.org/10.1016/j.ajhg.2018.11.014>
- Taliun, D., Harris, D. N., Kessler, M. D., Carlson, J., Szpiech, Z. A., Torres, R., Taliun, S. A. G., Corvelo, A., Gogarten, S. M., Kang, H. M., Pitsillides, A. N., LeFaive, J., Lee, S., Tian, X., Browning, B. L., Das, S., Emde, A.-K., Clarke, W. E., Loesch, D. P., ... Abecasis, G. R. (2021). Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*, 590(7845), 290–299. <https://doi.org/10.1038/s41586-021-03205-y>
- Thareja, G., Al-Sarraj, Y., Belkadi, A., Almotawa, M., Ismail, S., Al-Muftah, W., Badji, R., Mbarek, H., Darwish, D., Fadl, T., Yasin, H., Ennaifar, M., Abdellatif, R., Alkuwari, F., Alvi, M., Al-Sarraj, Y., Saad, C., Althani, A., Fethnou, E., ... Albagha, O. M. E. (2021). Whole genome sequencing in the Middle Eastern Qatari population identifies genetic associations with 45 clinically relevant traits. *Nature Communications*, 12(1).
<https://doi.org/10.1038/s41467-021-21381-3>
- Turnbull, C., Scott, R. H., Thomas, E., Jones, L., Murugaesu, N., Pretty, F. B., Halai, D., Baple, E., Craig, C., Hamblin, A., Henderson, S., Patch, C., O'Neill, A., Devereaux, A., Smith, K., Martin, A. R., Sosinsky, A., McDonagh, E. M., Sultana, R., ... Caulfield, M. J. (2018). The 100 000 Genomes Project: Bringing whole genome sequencing to the NHS. *BMJ (Online)*, 361. <https://doi.org/10.1136/bmj.k1687>
- Turro, E., Astle, W. J., Megy, K., Gräf, S., Greene, D., Shamardina, O., Allen, H. L., Sanchis-Juan, A., Frontini, M., Thys, C., Stephens, J., Mapeta, R., Burren, O. S., Downes, K.,

- Haimel, M., Tuna, S., Deevi, S. V. V., Aitman, T. J., Bennett, D. L., ... Raymond, F. L. (2020). Whole-genome sequencing of patients with rare diseases in a national health system. *Nature*, 583(7814), 96–102. <https://doi.org/10.1038/s41586-020-2434-2>
- Wu, D., Dou, J., Chai, X., Bellis, C., Wilm, A., Shih, C. C., Soon, W. W. J., Bertin, N., Lin, C. B., Khor, C. C., DeGiorgio, M., Cheng, S., Bao, L., Karnani, N., Hwang, W. Y. K., Davila, S., Tan, P., Shabbir, A., Moh, A., ... Wang, C. (2019a). Large-Scale Whole-Genome Sequencing of Three Diverse Asian Populations in Singapore. *Cell*, 179(3), 736-749.e15. <https://doi.org/10.1016/j.cell.2019.09.019>
- Wu, D., Dou, J., Chai, X., Bellis, C., Wilm, A., Shih, C. C., Soon, W. W. J., Bertin, N., Lin, C. B., Khor, C. C., DeGiorgio, M., Cheng, S., Bao, L., Karnani, N., Hwang, W. Y. K., Davila, S., Tan, P., Shabbir, A., Moh, A., ... Wang, C. (2019b). Large-Scale Whole-Genome Sequencing of Three Diverse Asian Populations in Singapore. *Cell*, 179(3), 736-749.e15. <https://doi.org/10.1016/j.cell.2019.09.019>
- Xue, Y., Chen, Y., Ayub, Q., Huang, N., Ball, E. V., Mort, M., Phillips, A. D., Shaw, K., Stenson, P. D., Cooper, D. N., & Tyler-Smith, C. (2012). Deleterious- and disease-allele prevalence in healthy individuals: Insights from current predictions, mutation databases, and population-scale resequencing. *American Journal of Human Genetics*, 91(6), 1022–1032. <https://doi.org/10.1016/j.ajhg.2012.10.015>

Figure Titles and Legends

Figure 1. Qatar Genome Program, timelines, and regional context.

a) Three phases project timeline and current status. b) Qatar Geographical map. Qatar is located in the north-eastern coast of the Arabian Peninsula with an area of 11,521 km² sharing borders with Saudi Arabia from the south and maritime borders with Bahrain, UAE, and Iran. c) The Arabian Peninsula is believed to be the first stop in human migration out of Africa, and home for the first ancient Eurasian populations, whom later spread throughout Asia and Europe.

Figure 2. Variants distribution and allele frequency spectrum of QGP data.

a) Number of SNVs and INDELS present within the QGP data. b) Known and novel variants distribution of QGP data. c) QGP variants classification based on minor allele frequency (MAF). d) Proportion of known and novel singletons within the QGP data. e) Classification of DM variants based on pattern of inheritance. Inheritance patterns of genes were derived from OMIM database. f) Distribution of DM variants among individuals in QGP sub clusters. g) QGP variants classified as both DM and pathogenic/likely pathogenic.

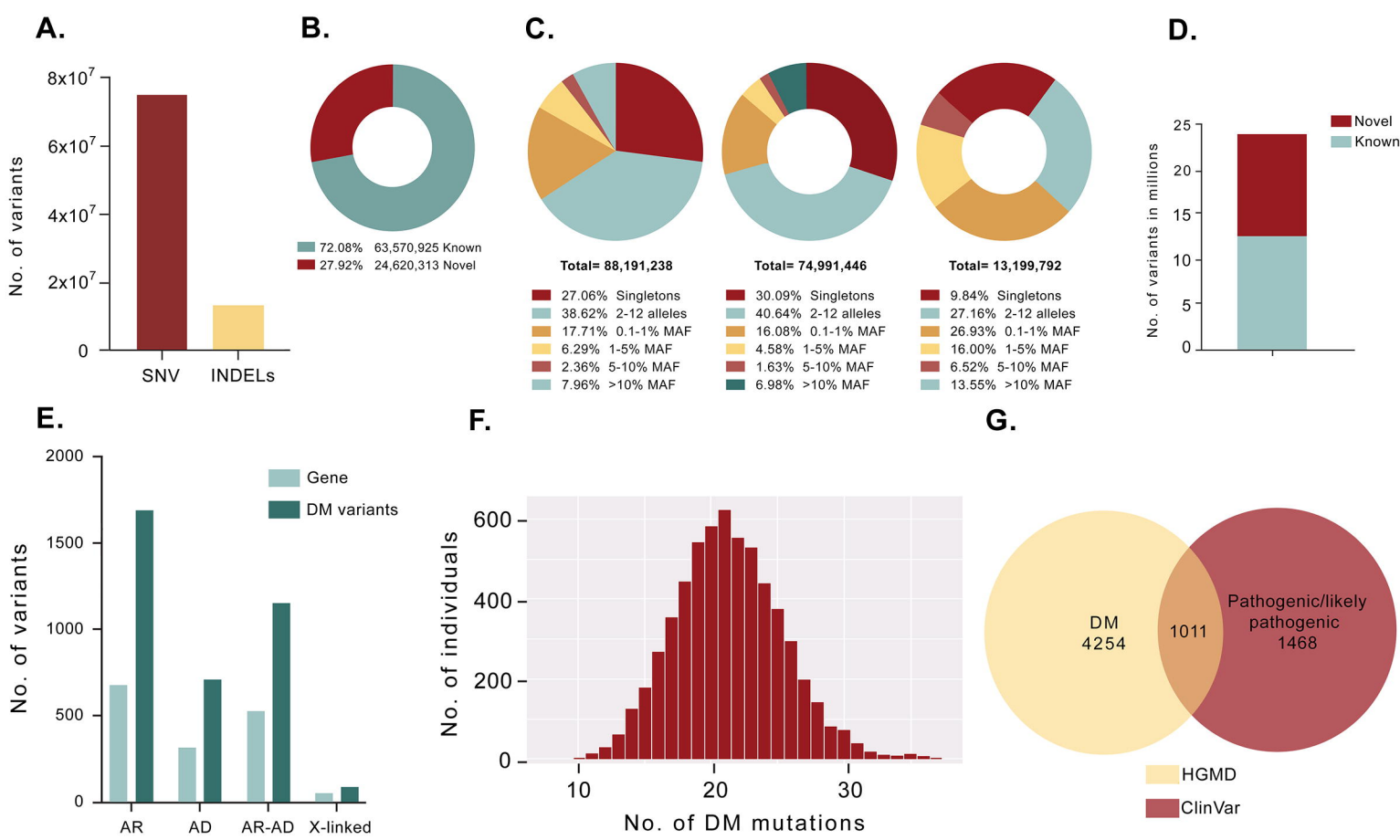
Tables

Table 1. Pathogenic variants unique to the Middle East region. Pathogenic variants exclusively reported in QGP and GME (Greater Middle East) variome project. QGP_MAF - Minor allele frequency in QGP data; GME_MAF – Minor allele frequency in GME variome project.

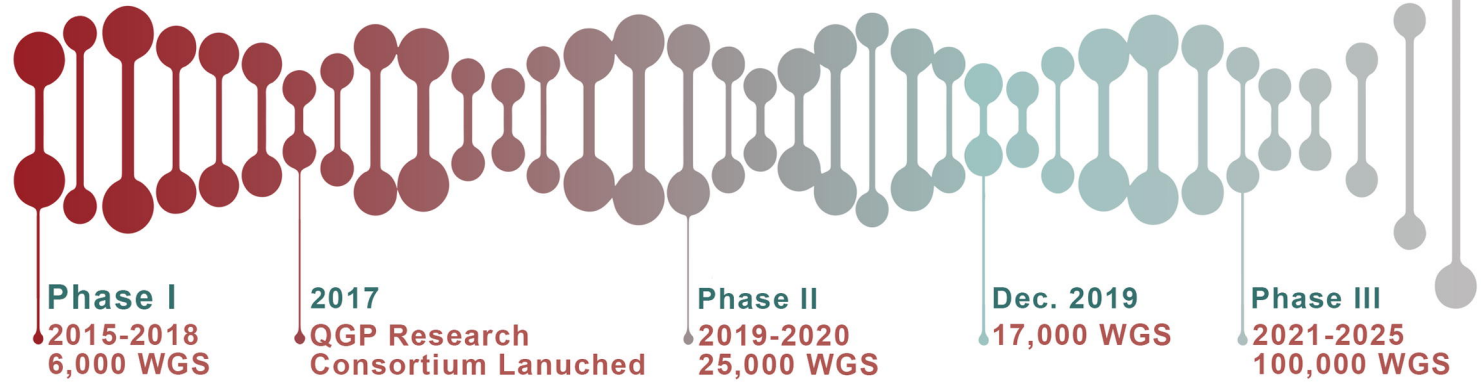
Chrom	Pos	ID	QGP_MAF	GME_MAF	GENE	HGVS_C	HGVS_P	Annotation	HGMD	CLINVAR	Disease Phenotype
2	47604159	rs606231204	0.00223251	0.000503525	EPCAM	c.583dupC	p.Gln195fs	frameshift variant	DM	Pathogenic	Congenital Tufting Enteropathy (CTE)
2	73676380	rs746640196	0.000165399	0.000503525	ALMS1	c.2723C>G	p.Ser908*	stop gained	DM	Likely pathogenic	Alstrom syndrome
2	172305304	rs797045038	0.00727634	0.001007049	DCAF17	c.436delC	p.Ala147fs	frameshift variant	DM	Pathogenic	Woodhouse-Sakati syndrome
3	113119479	rs866096259	0.00330961	0.000503525	WDR52	c.1387G>T	p.Glu463*	stop gained	DM	Pathogenic	Spermatogenic failure
4	108866582	rs397514513	0.00223251	0.001510574	CYP2U1	c.947A>T	p.Asp316Val	missense variant	DM	Pathogenic	Spastic paraplegia
4	119736287	rs730882211	0.00148834	0.001007049	SEC24D	c.700G>C	p.Gly234Arg	missense variant	DM	Likely pathogenic	Intellectual disability/Seizures
6	135776888	rs121434350	0.000165399	0.001009082	AHI1	c.1328T>A	p.Val443Asp	missense variant	DM	Pathogenic/Likely pathogenic	Joubert syndrome
8	145741257	.	8.27E-05	0.000503525	RECQL4	c.1149G>A	p.Trp383*	stop gained	DM	Pathogenic	Rothmund-Thomson syndrome
9	111899809	rs878853280	0.000413497	0.001007049	FRRS1L	c.961C>T	p.Gln321*	stop gained	DM	Pathogenic	Epileptic encephalopathy
15	65295453	rs863224897	8.27E-05	0.000503525	MTFMT	c.1116delT	p.Pro373fs	frameshift variant	DM	Likely pathogenic	Moyamoya disease
16	77369781	rs148319220	8.27E-05	0.001007049	ADAMTS18	c.1731C>G	p.Cys577Trp	missense variant	DM	Pathogenic	Microcornea, myopic chorioretinal atrophy, and telecanthus(MMCAT)
19	11304215	.	0.00256325	0.000504032	KANK2	c.541A>G	p.Ser181Gly	missense variant	DM	Pathogenic	Nephrotic syndrome,16
21	47805894	rs387906928	8.27E-05	0.000503525	PCNT	c.3460G>T	p.Glu1154*	stop gained	DM	Pathogenic	Microcephalic osteodysplastic primordial dwarfism type 2 (MOPD2)
22	27012112	rs1064793935	0.000248057	0.000504541	CRYBB1	c.171delG	p.Asn58fs	frameshift variant	DM	Pathogenic	Cataract

Table 2. Median number of variant sites per genome. Novel SNV and INDELS: Variants, which are not reported in dbSNP or gnomAD or 1000G project. GERP (Genomic Evolutionary Rate Profiling) score: Scores >3 represent highly conserved positions. ADM- Admixed, AFR- Africans, GAR- general Arabs, PAR- Peninsular Arabs, SAS- South Asians, WEP- Arabs of Western Eurasia and Persia.

Annotation	QGP (n=6,045, Depth = 32.4x)	ADM (n=1,180, Depth =32.2x)	AFR (n=92, Depth=31.9x)	GAR (n=2,311, Depth =32.2x)	PAR (n=1,052, Depth=32.2x)	SAS (n=38, Depth =31.9x)	WEP (n=1,372, Depth=32.2x)
SNV	3,467,270	3,596,354	3,967,082	3,466,051	3,391,850	3,492,506	3,458,604
INDELS	1,107,288	1,128,043	1,207,016	1,105,836	1,094,173	1,113,900	1,101,075
Singlet ons	1,336	9,242	17,606	2,484	408	20,056	3,193
Novel SNV	18,453	21,311	25,993	16,419	12,022	23,814	20,788
Novel INDELS	45,756	46,263	48,406	45,752	46,061	46,195	45,107
Synonymous	10,657	11,094	12,285	10,643	10,372	10,768	10,635
Missense	10,681	11,241	12,464	10,921	10,684	10,997	10,895
Intron	1,617,713	1,659,957	1,833,502	1,618,061	1,586,985	1,632,466	1,613,603
Intergenic	1,760,161	1,806,271	1,989,241	1,759,321	1,726,938	1,774,147	1,756,828
Conserved: GERP>3	3,751	3,859	4,233	3,755	3,667	3,788	3,738



A. QATAR GENOME



B.



C.

