# GWAS and Population Genetics
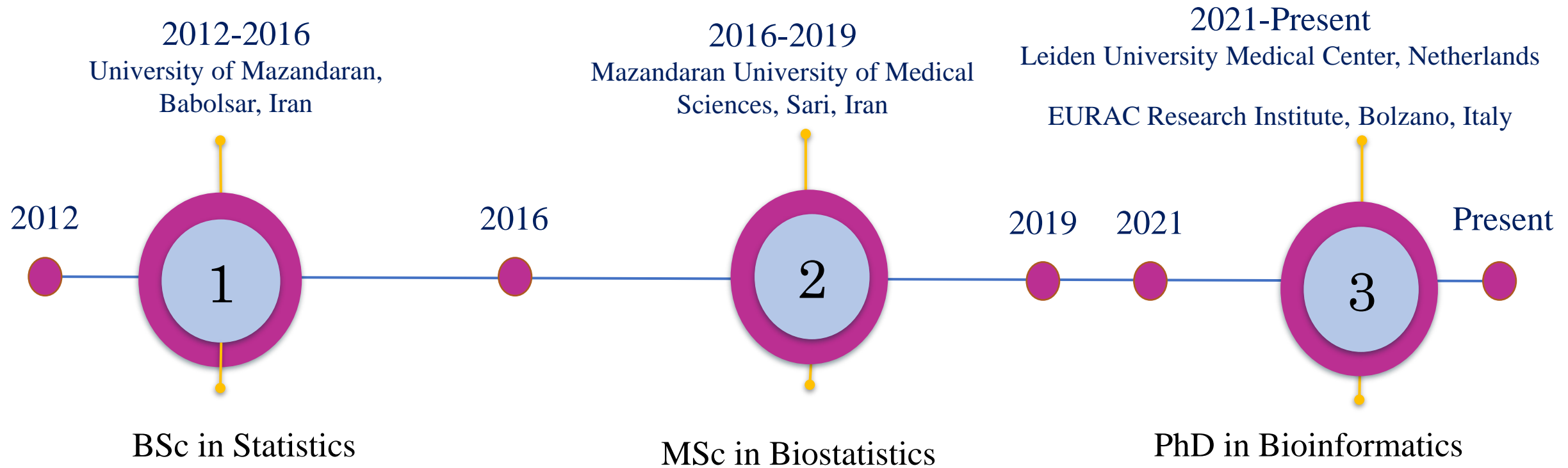
Dariush Ghasemi
PhD Fellow at
Leiden University Medical Center
EURAC Research Institute

July 2022

# Education

2012-2016
University of Mazandaran,
Babolsar, Iran

2016-2019
Mazandaran University of Medical
Sciences, Sari, Iran

2021-Present
Leiden University Medical Center, Netherlands

EURAC Research Institute, Bolzano, Italy

2012

2016

2019    2021

Present

1

2

3

BSc in Statistics

MSc in Biostatistics

PhD in Bioinformatics

# Genome

## DNA

ACTGACCTAG**ATCAGT**GTAGC**GATC**GTATACGAG**ACCGAT**TCATCGGCAT

↓ transcription

## RNA

**AUCAGU****CGAUC****ACCGAU**

↓ translation

## protein

# Structure of a DNA



## Bases:

**A**denine

**T**hymine

**C**ytosine

**G**uanine

**1953** Watson & Crick

# Genome



Human karyotype
Male

DNA molecule
Gene 1
Gene 2
Gene 3

DNA strand (template)
3′  A C C A A A C C G A G T  5′

TRANSCRIPTION

mRNA
5′  U G G U U U G G C U C A  3′

Codon

TRANSLATION

Protein

Trp  Phe  Gly  Ser

Amino acid

©1999 Addison Wesley Longman, Inc.

# Twenty years of technological progress



Mapping the human genome
2001



Characterizing common variation
2003 - 2007



Millions of features possible on each GeneChip® array

DNA probes in one corner of an Affymetrix array

Genome-wide association (GWAS)
2005 - present



Next generation sequencing
2009 - present

# Technologies



| | Arrays | Exomes | Genomes |
|---|---|---|---|
| Upside | Hits + epidemiology | Gene identification | Comprehensive capture |
| Downside | Hit interpretation | Limited Scope | Cost (small N) |

# Single Nucleotide Polymorphism (SNP)



- Mutation that arose at some point in demographic history
- Typically, each SNP has two alleles (bases)
- Each SNP is eventually given an "rs" number rs214621

# Structural Mutations

# Variants' Distribution in Population

# Simplest Regression Model of Association

$$Y_i = \alpha + \beta X_i + e_i$$

where

$$Y_i = \quad \text{trait value for individual i}$$

$$X_i = \quad \text{1 if allele individual i has allele 'A'}$$

$$\text{0 otherwise}$$

i.e., test of mean differences between 'A' and 'not-A' individuals

# Haplotypes

# Reference Panels for Genomic Imputation



## HapMap (haplotype map) Project

270 individuals:

30 parent-offspring trios of the Yoruba from Ibadan, Nigeria (YRI)
30 trios of Utah residents with European ancestry (CEU)
45 individuals from Beijing, China (CHB)
45 individuals from Tokyo, Japan (JPT)

The International HapMap Consortium (2005). A haplotype map of the human genome. *Nature*.

# Reference Panels for Genomic Imputation



## 1000 Genomes Project

Phase 1: 1,092 individuals from 14 populations..

Phase 3: 2,504 individuals from 26 populations (~500 samples form each 5 continental ancestry groups, with ~5 populations for each group)



The 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature*
The 1000 Genomes Project Consortium (2015). A global reference for human genetic variation. *Nature*

# Reference Panels for Genomic Imputation

The Haplotype Reference Consortium (HRC)

nature
genetics

A reference panel of 64,976 haplotypes for genotype imputation

The Haplotype Reference Consortium (2016). A reference panel of 64,976 haplotypes for genotype imputation. *Nature Genetics*.

# Reference Panels for Genomic Imputation

Trans-Omics for Precision Medicine (TOPMed)

# nature

# Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program

Taliun, D., Harris, D.N., Kessler, M.D. *et al.* (2021). Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*

# Step1: Variant Identification (sequencing)

Software:
- [freeBayes](#)
- [GATK](#)

# Step2: Population Stratification

Software:
- EIGENSOFT
- snpStats

# Step3: Statistical Tests

# Step3: Statistical Tests | Performing GWAS via EPACTS

```bash
#!/bin/bash

BASE=/home/dghasemi
PHENO=$BASE/phenodf4_NEW_scheme_W.ped
OUT=/home/dghasemi
KIN=/
DATA=
BIN=

for t in SerumCreatinine.Stdw.Res          eGFRw.Res          eGFRw.log.Res
do
for i in `seq 1 22`
do
 echo ${BIN}epacts single -vcf $DATA/chr$i.vcf.gz \
 -ped ${PHENO} --pheno $t \
 -out ${OUT}/$t.chr$i \
 -test q.emmax \
 -kinf ${KIN} --chr $i \
 -field DS \
 --run 24 --mosix-nodes \"\"
done
done
```

# Step4: Examine local region

Software:
- [PLINK](#)
- [Annotating Genomic Variants Workflow](#)

# Multiple Testing Adjustments

- Measure 10,000 genes

- Calculate 10,000 p-values

- Call genes "significant" if p-value < 0.05

- Expected Number of False Positives:

    10,000 × 0.05 = 500 False Positives

- Bonferroni Correction:

    P-values less than α/m are significant

# PhD Project: Kidney function Biomarkers in CHRIS Study Participants

• Database: CHRIS study

Cooperative Health Research In South Tyrol study

N= 10,500

10,500 GWAS

3,600 WES

Integrated via genotype imputation

**CHRIS**
eurac research

Südtiroler Gesundheitsstudie
Studio sulla salute in Alto Adige

. 13 municipalities
. 28,000 inhabitants (adults)
. 13,393 participants (2011-2018)
. Mean age 46 (18-93)
. 55% females

Schlanders: one hospital
A single site study

South Tyrol, Italy

Phenotypes: Serum Creatinine,      Serum Albumin,
Urinary Creatinine,      Urinary Albumin,
UACR,       eGFR

# What is imputation? (Marchini & Howie 2010)

# Genomic Imputation:
## 1- Starting Data

**Genotyped sample**

.   .   **C**   .   .   **G**   .   **C**   .

**Reference haplotypes**

A  G  A  T  C  T  C  C  T

A  G  C  T  C  T  C  A  T

A  G  A  T  C  G  C  C  T

A  G  A  T  C  T  A  C  T

# Genomic Imputation:
## 2- Identify shared regions of chromosome

**Genotyped sample**

| . | . | **C** | . | . | **G** | . | **C** | . |
|---|---|---|---|---|---|---|---|---|

**Reference haplotypes**

A G A T C T C C T

**A G C T** C T C A T

A G A T C **G** C **C** T

A G A T C T A C T

# Genomic Imputation: 3. Fill in missing genotypes

**Genotyped sample**

| A | G | **C** | T | C | **G** | C | **C** | T |
|---|---|---|---|---|---|---|---|---|

**Reference haplotypes**

| A | G | A | T | C | T | C | C | T |
|---|---|---|---|---|---|---|---|---|
| A | G | **C** | T | C | T | C | A | T |
| A | G | A | T | C | **G** | C | **C** | T |
| A | G | A | T | C | T | A | C | T |

# Genomic Imputation: Minimac4

- https://github.com/statgen/Minimac4

- Building on the work from Gonçalo Abecasis, Christian Fuchsberger and colleagues

- Analysis options
  - SAIGE
  - BoltLMM
  - plink2

## Next-generation genotype imputation service and methods

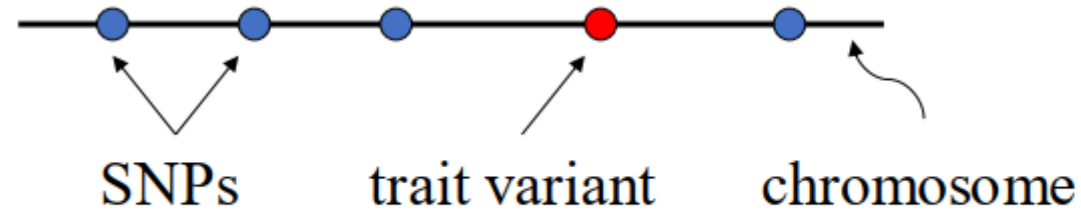Sayantan Das, Lukas Forer, Sebastian Schönherr, Carlo Sidore, Adam E Locke, Alan Kwong, Scott I Vrieze, Emily Y Chew, Shawn Levy, Matt McGue, David Schlessinger, Dwight Stambolian, Po-Ru Loh, William G Iacono, Anand Swaroop, Laura J Scott, Francesco Cucca, Florian Kronenberg, Michael Boehnke, Gonçalo R Abecasis ✉ & Christian Fuchsberger ✉
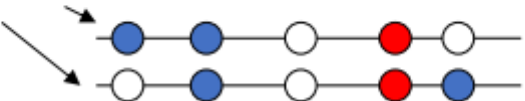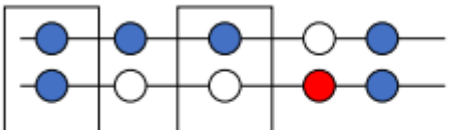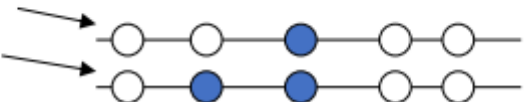
# Definitions

# Allelic Associations

# Biostatistics & Epidemiology

Cristian Pattaro
Group Leader

-Biostatistics
-Epidemiology
-Genetics

Luisa Foco
Senior Researcher

-Statistical genetics
-Genetic epidemiology
-Biomedical statistics

Roberto Melotti
Senior Researcher

-Epidemiology
-Biomedical statistics

Fabiola Del Greco M.
Senior Researcher

-Causal inference
-Statistical genetics
-Epidemiology

Ryosuke Fujii
Post Doc Researcher

-Epidemiology
-Genetics

Martin Gögele
Researcher

-Biodemography
-Epidemiology
-Study design

Giulia Barbieri
PhD Student

-Epidemiology
-Biomedical statistics

Maeregu W. Arisido
Post Doc Researcher

-Biostatistics
-Statistical methods

Daniele Bottigliengo
Post Doc Researcher

-Causal inference
-Biostatistics

Dariush Ghasemi S.
PhD Student

-Biostatistics
-Machine Learning
-Statistical genetics

Laura Barin
Post Doc Researcher

-Epidemiology
-Public Health
-Biostatistics

Rebecca Lundin
Senior Researcher

-Public Health
-Epidemiology

Daniele Giardiello
PhD Student

-Biostatistics
-Cancer epidemiology

Thank You for the Attention