# Large Genomic Data Sets in Autism Research: SFARI DOMA Practices

Natalia Volfovsky

November 16, 2017

# SFARI

SIMONS FOUNDATION
AUTISM RESEARCH INITIATIVE

*Mission: improve the understanding, diagnosis and treatment of autism spectrum disorders (ASD) by funding innovative research of the highest quality and relevance*

https://www.sfari.org

Identify risk factors whether genetic, environmental or epidemiological.

Use non-human organisms to understand how these risk factors alter brain function and animal behavior.
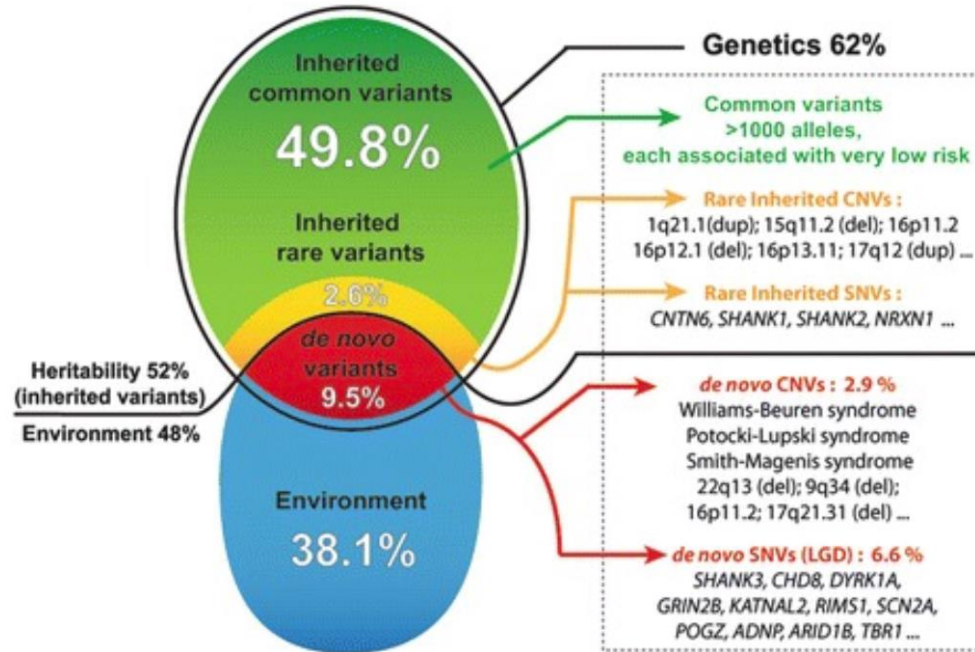
Promote preclinical and clinical investigations to improve autism diagnosis & therapy.

SIMONS FOUNDATION

# Autism spectrum disorder(ASD) is extremely heterogeneous
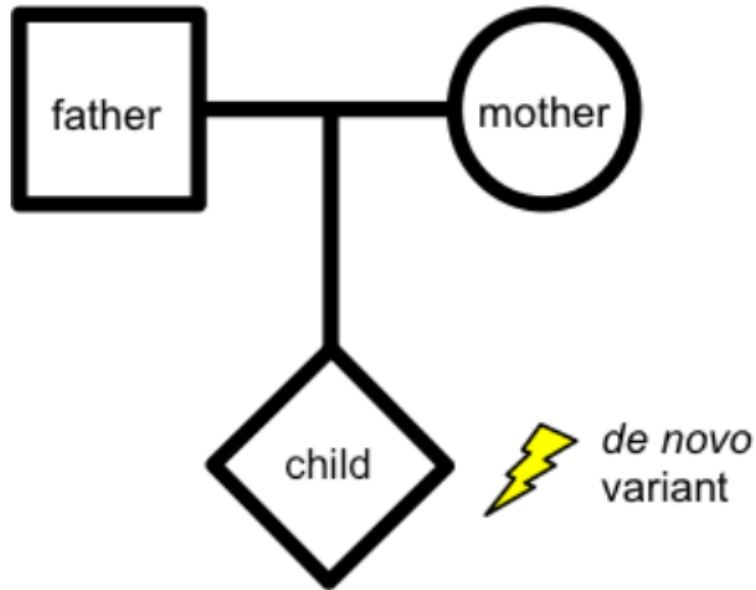


Pam Feliciano, 2017

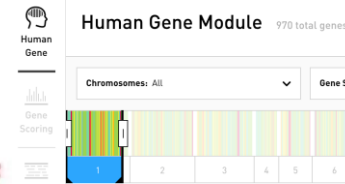# Autism risk is complex



Huguet, Benabou, and Bourgeron, 2016
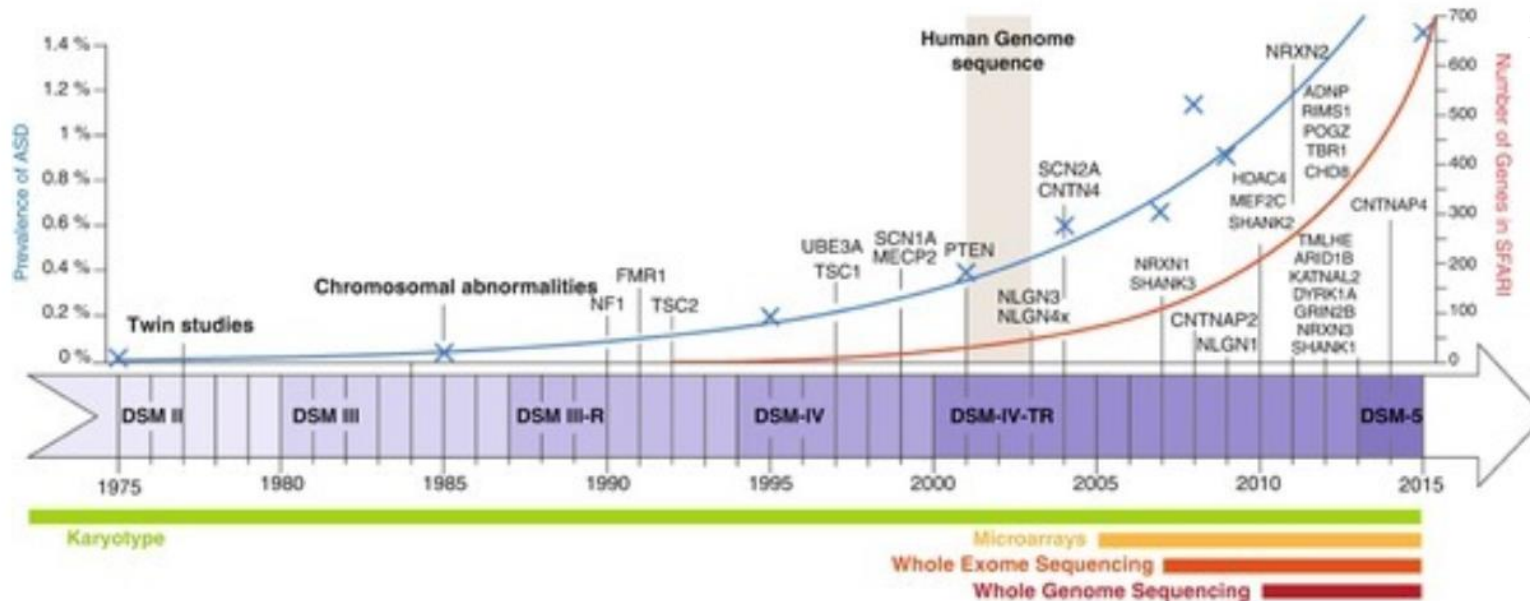
# De Novo Mutations



- Genotype is present in neither parent, but usually heterozygous in the child
- Spontaneous genetic mutations

# The history of the genetics of autism



gene.sfari.org

Huguet, Benabou, and Bourgeron, 2016

SIMONS FOUNDATION

# SFARI Research Cohorts

https://www.sfari.org



Simons Simplex Collection
(SSC)
~10,000 individuals

SIMONS VIP
Simons Variation in Individuals Project

~ 1,500 individuals

SPARK
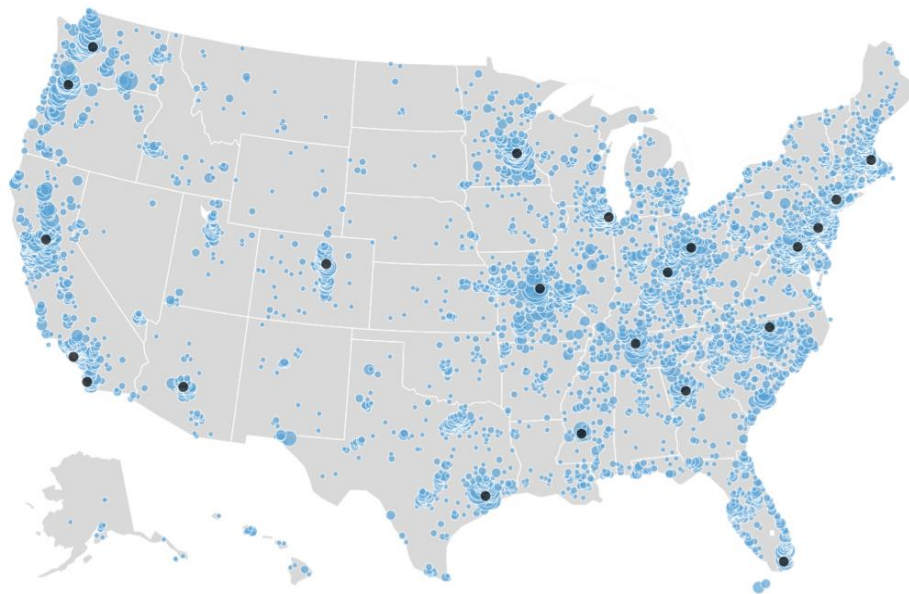Igniting autism research
Improving lives

SIMONS FOUNDATION

# SPARK
Igniting autism research
Improving lives

# Simons Foundation Powering Autism Research through Knowledge

https://sparkforautism.org

Recruit, engage, and retain 50
individuals with ASD
and their biological family
members to:

- identify causes of ASD
- enable genotype-driven res
- find better treatments to imp
  lives

Individuals with autism
    n= 28,501
Family members
    n= 52,070
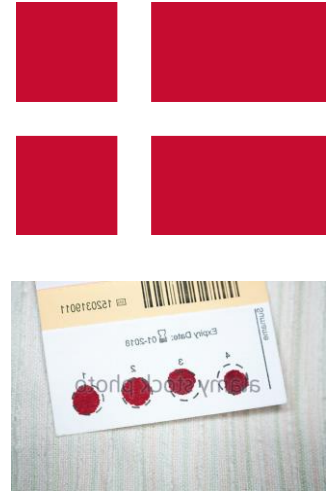Usable saliva samples
    n= 29,374
Complete trios + DNA
 n = 4,294 trios

Unaffected sibs + DNA   n= 2,086

Pam Feliciano, October 2017

SIMONS FOUNDATION

# Danish Neonatal Screening Biobank

The DNSB houses blood spots from all individuals born in Denmark Since 1981 (more than 2 million individuals). These samples can be linked to first-degree and other relatives, and there are now at least 15,000 diagnosed cases of autism in the biobank, as well as 15,000 cases of ADHD and 25,000 cohort-matched controls.

Exome sequencing of blood spot-derived DNA:
collaboration between Mark Daly (Broad Institute) and Preben Mortensen (Aarhus University, Denmark).

## Data Sharing Policy

SFARI aims to support reproducible scientific research of the highest quality.  This can be greatly facilitated by making all raw data, analysis methods, and, when applicable, computer code, available to the research community as quickly and transparently as possible.

Data generated from collections is considered "community resources" and, as such, all these data, including, but not limited to, alignment files and variant calls, will be made available to the entire research community, pre-publication, in real-time, as they are produced.

Besides being subject to any limitations by consent and IRB protocol, a publication embargo will also be enforced for 6 months after the defined project is complete while foundation-sponsored analyses are underway.

# SFARI Data Overview



Sequencing and preprocessing

Sequencing Center | Sequencing Center | Sequencing Center

Storage and secondary analysis

Flatiron Institute | AWS | Tape Archive at Fermi Lab

Data sharing and collaboration

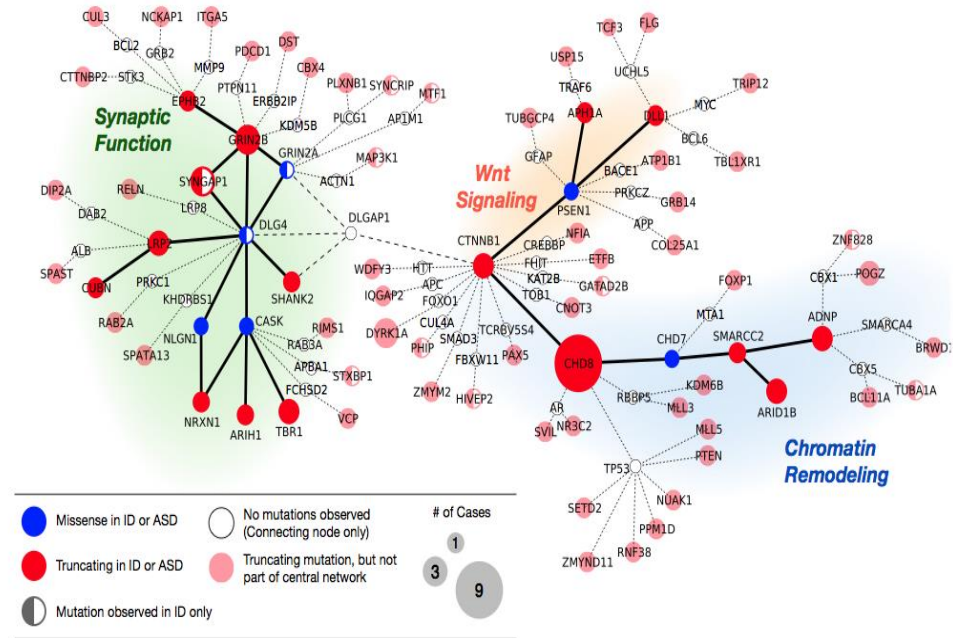Globus | Data portal | Data portal | Data portal

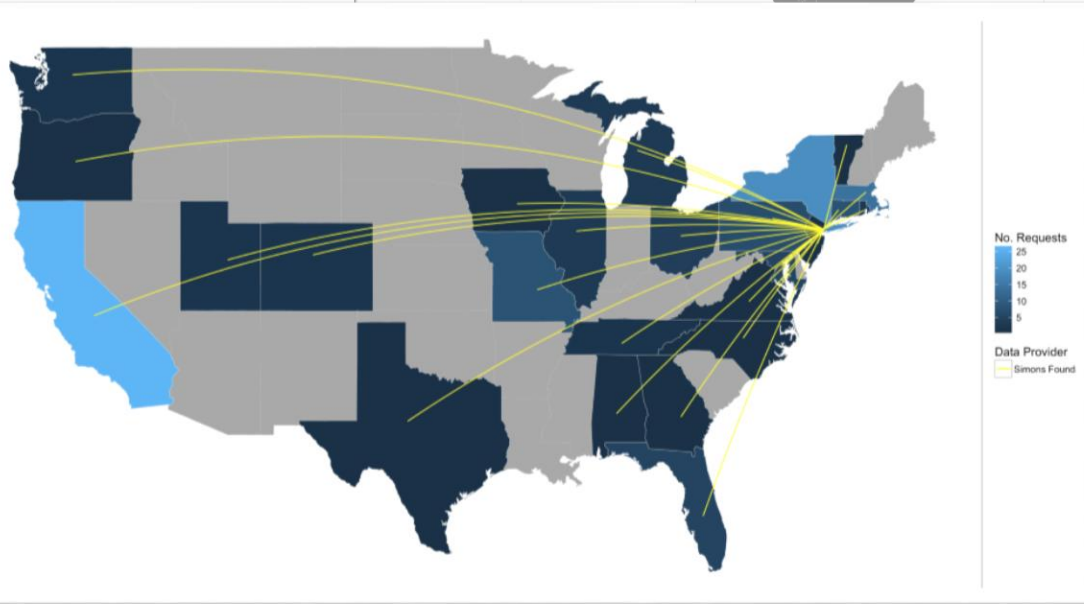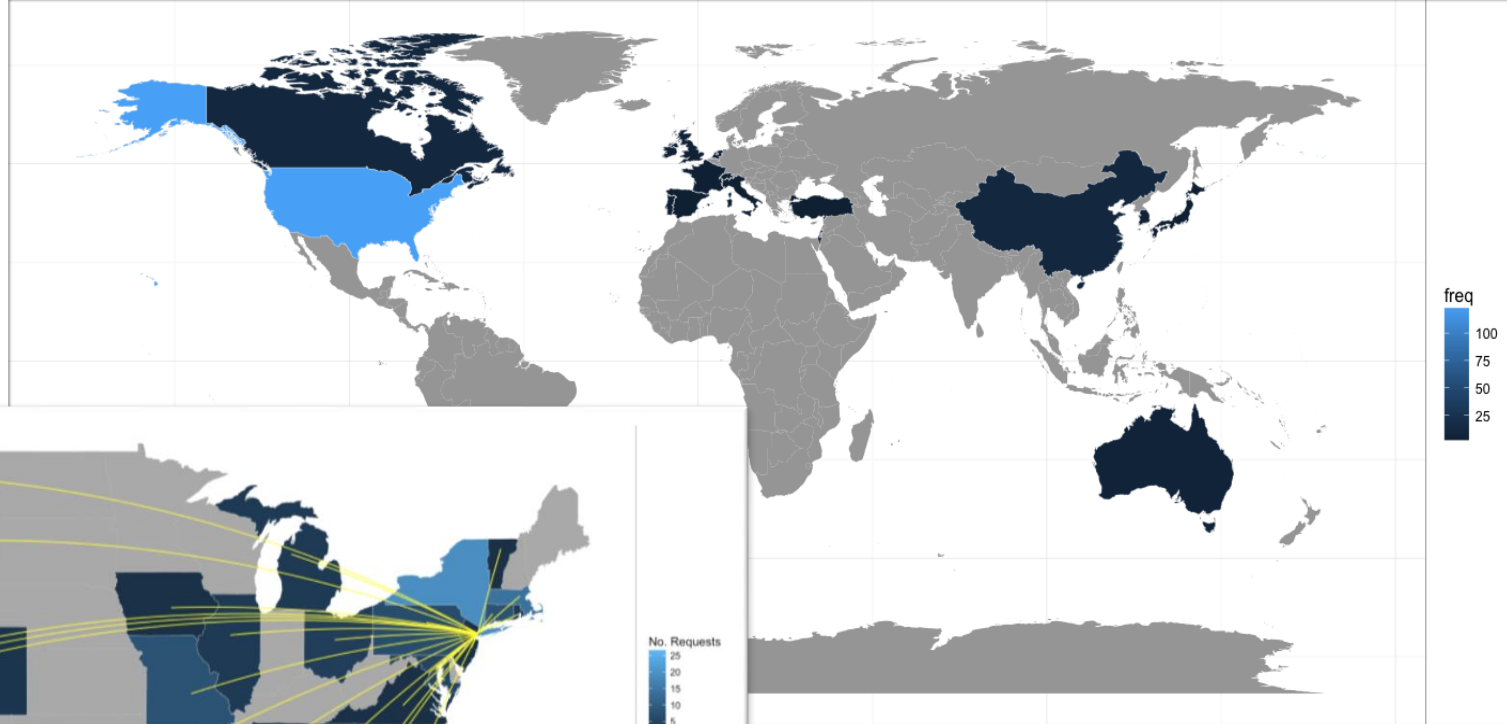SIMONS FOUNDATION

# Data Analysis

## Genomic Analysis Workflows



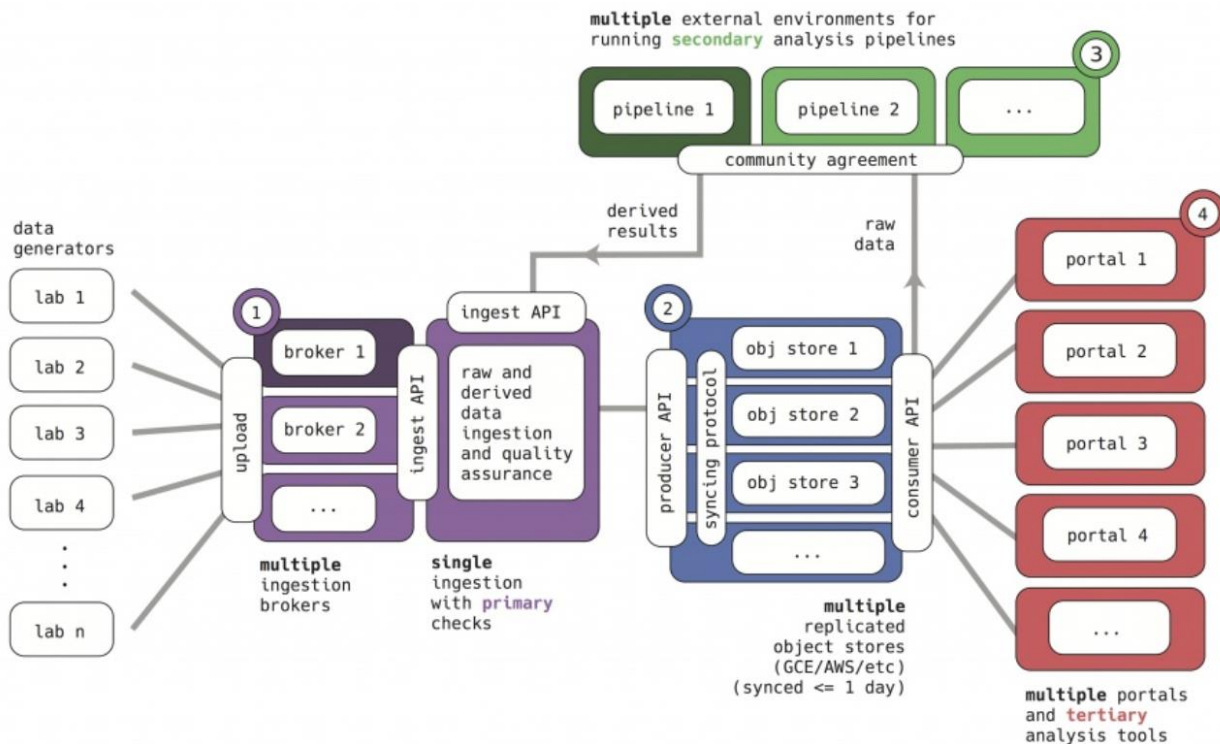## Interpretation of Genomic Results

# SFARI Data Requests



**SFARI** genomic data requests through **SFARIBase @ sfari.org**

SIMONS FOUNDATION

# Challenges of Data Sharing

- Data management
- Collaboration utilities
- Users' ability to deploy their own analysis tools
- Ability to harmonize data from different projects
- Ability to harmonize data from different archives of the same project
  - AWS, Google, local servers
- Price, speed
  - storage, data transfers, computin
- Support of international collaborations and ability to deal with data governance across countries, including local restrictions

# Human Cell Atlas (HCA) Data Coordination Portal

# Global Alliance for Genomics and Health (GA4GH): Data Biosphere



SIMONS FOUNDATION

# Genomics data portals

- Support
  - **collaboration**
  - **data sharing and**
  - **reproducibility** in scientific research
- Technology
  - **API access**
  - **docker containers and**
  - **standardized workflow description languages**
- Leverage experiences of scientific data management communities
- Invest in **the implementations of science specific standards**

# Genome and Phenotype Tool (GPF)

https://gpf.sfari.org/

## Query





Features:
- SPARK Content:
  - De novo variants
  - Rare and common transmitted variants
  - Phenotypic data
- Integration:
  - SSC exome and whole-genome
  - VIP
- Query: Interface
  - By gene
  - By set of genes (i.e. by pathway)
  - By predicted variant effect (i.e. LGDs)
  - By variant frequency
  - By transmission pattern (i.e. de novo or transmitted from mother)
  - By phenotypic properties (i.e. affected children with SCQ score larger than 20)
- Analysis tools:
  - Enrichment of de novo variant in a given gene set.

SIMONS FOUNDATION

# SFARI-IOBIO Data Federation Project

**bam.iobio**: http://bam.iobio.io
**vcf.iobio**: http://vcf.iobio.io
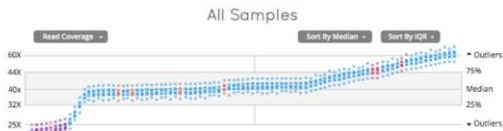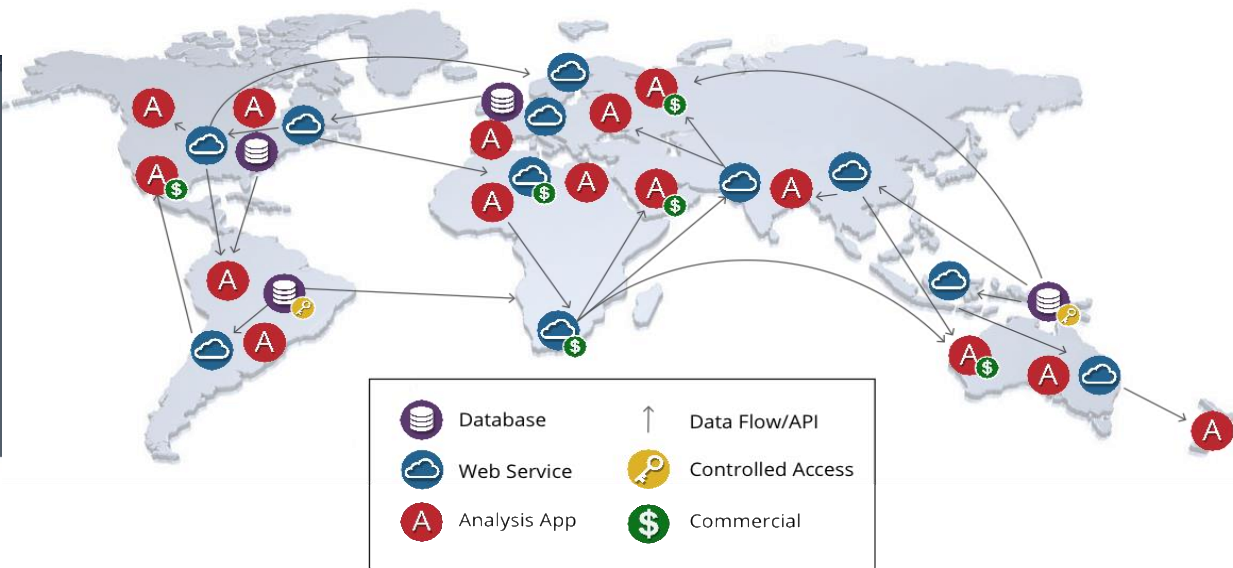**gene.iobio**: http://gene.iobio.io



Data Hub
register + find databases
hub.iobio/data

Service Hub
register + publish + find services
hub.iobio/services
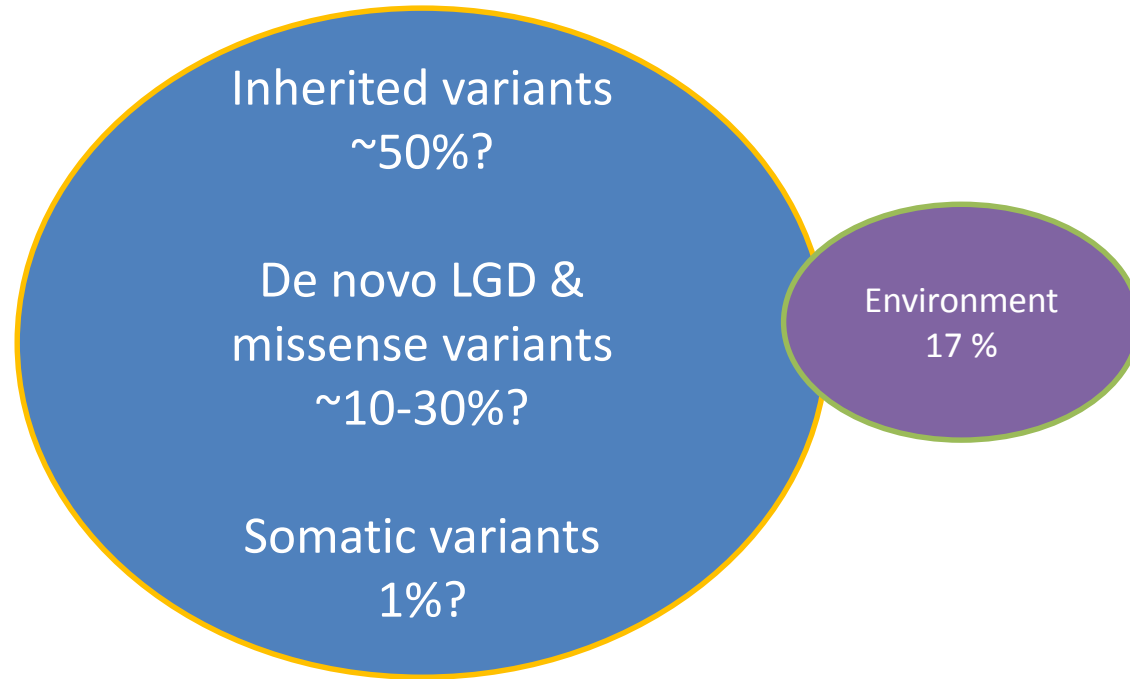
App Hub
register + find apps
hub.iobio/apps

Legend:
- Database
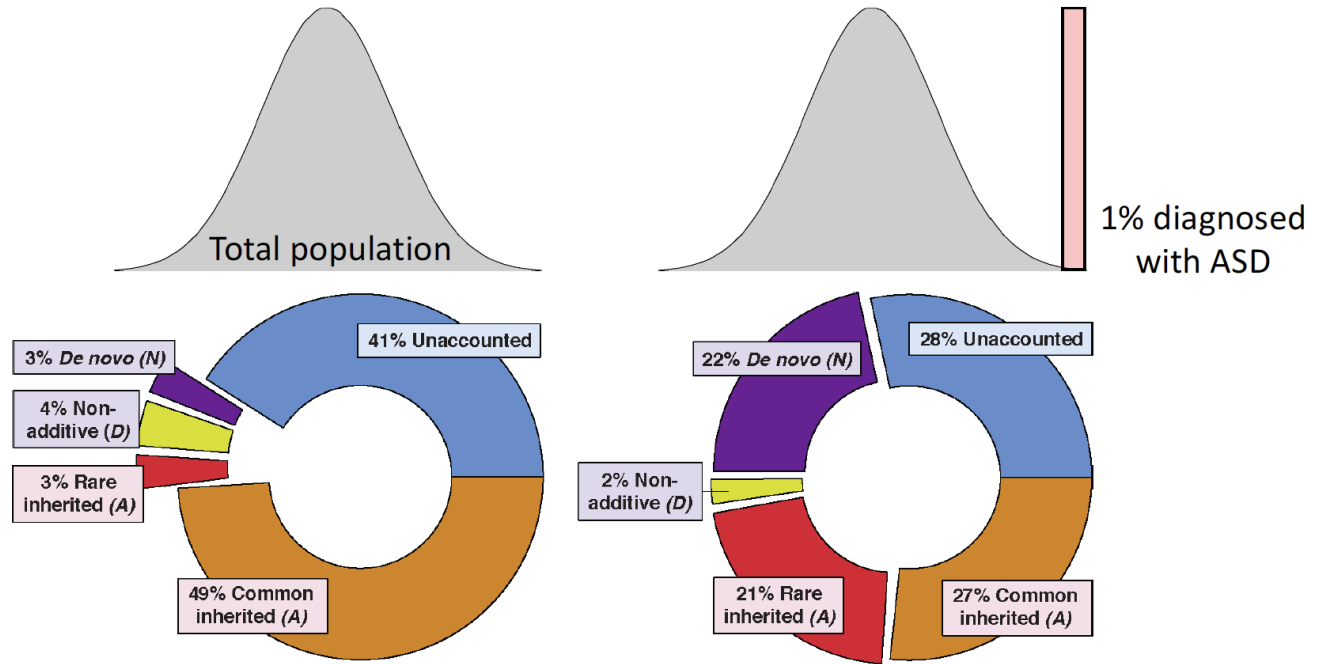- Web Service
- Analysis App
- Data Flow/API
- Controlled Access
- Commercial

# Acknowledgements

SIMONS FOUNDATION

Thank you!

# Autism risk is complex



Inherited variants ~50%?

De novo LGD & missense variants ~10-30%?

Somatic variants 1%?

Environment 17 %

Pam Feliciano, 2017

SIMONS FOUNDATION

# Distribution of genetic risk in autism



Gauglet at al, 2014

(from Stephan Sanders)

SIMONS FOUNDATION