# The Biorepository for Integrative Genomics initiative



BIG

Lebonheur Children's Hospital
Regional One Health Hospital
East Tennessee Health Science Center
DaVita CKD Cohort
Family Resilience Initiative

Sequenced — 13,375
Biospecimen — 16,537
Electronic Health Records — 40,754

UटHSC

# Ancestry inference in BIG



- Four continental level ancestries

- Considerable proportion of admixed individuals

- High genetic diversity geographically concentrated

*Buonaiuto* , Marsico*, Nat Comm, 2025*

# Identity by Descent (IBD)



IBD segment

# Self-reported race and genetic ancestry through the lens of relatedness



**b** **Self-reported race**

**c** Inferred ancestry

same          Different

UTHSC.

# From pairwise IBD to population structure



IBD based network

Nodes: 17,234
Edges: 2,078,007
Global density: 0.014

Identity-by-descent captures Shared Environmental

Factors at Biobank Scale

Franco Marsico[1, 2, †], Silvia Buonaiuto[1], Ernestine K Amos−Abanyie[1], Lokesh K
Chinthala[3], Akram Mohammed[3], Regeneron Genetics Center[4], Robert J
Rooney[1,5], Robert W Williams[1,6], Robert L Davis[3], Terri H Finkel[7], Chester W
Brown[1, 8], Pjotr Prins[1], and Vincenza Colonna[1, 2, 5, †]

[1]Dept of Genetics, Genomics and Informatics, UTHSC, USA
[2]Institute of Genetics and Biophysics, National Research Council, Naples, 80111, Italy
[3]Center for Biomedical Informatics, UTHSC, USA
[4]Regeneron Genetics Center, Tarrytown, NY, USA
[5]Dept of Pediatrics, UTHSC, USA
[6]Center for Integrative and Translational Genomics, UTHSC, USA
[7]Dept of Pediatrics, Division of Rheumatology, UTHSC, USA
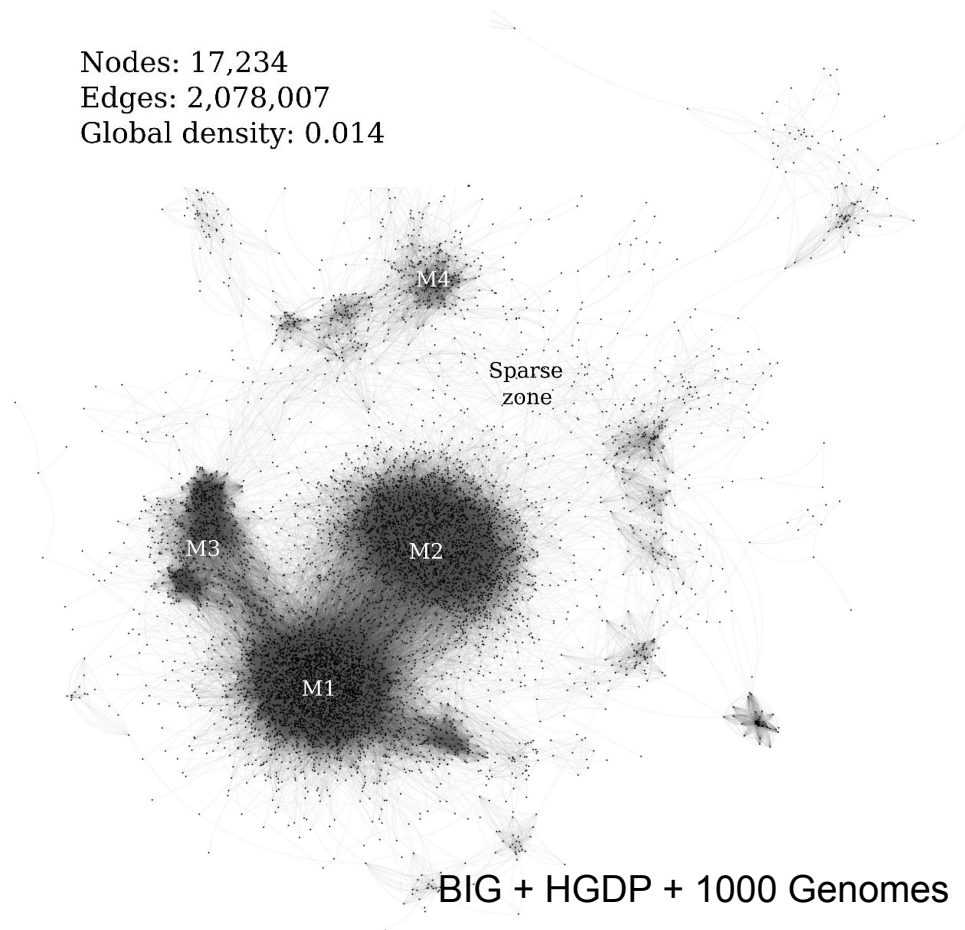[8]Dept of Pediatrics, Division of Genetics, UTHSC, USA
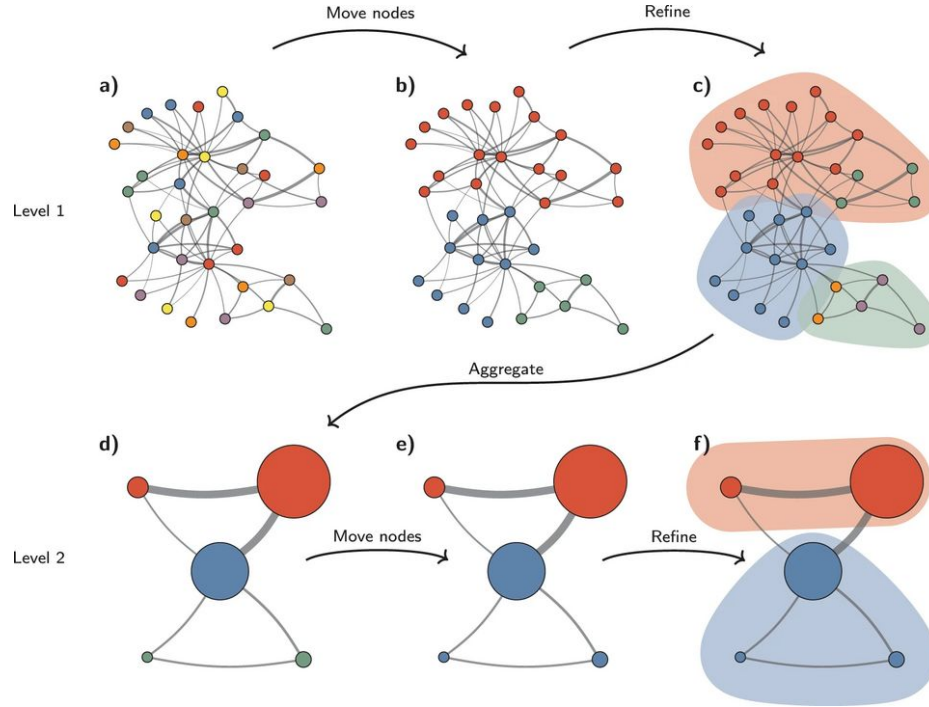[†]Corresponding authors: vcolonna@uthsc.edu, fmarsic1@uthsc.edu

https://www.biorxiv.org/content/10.1101/2025.05.03.652048v1.full

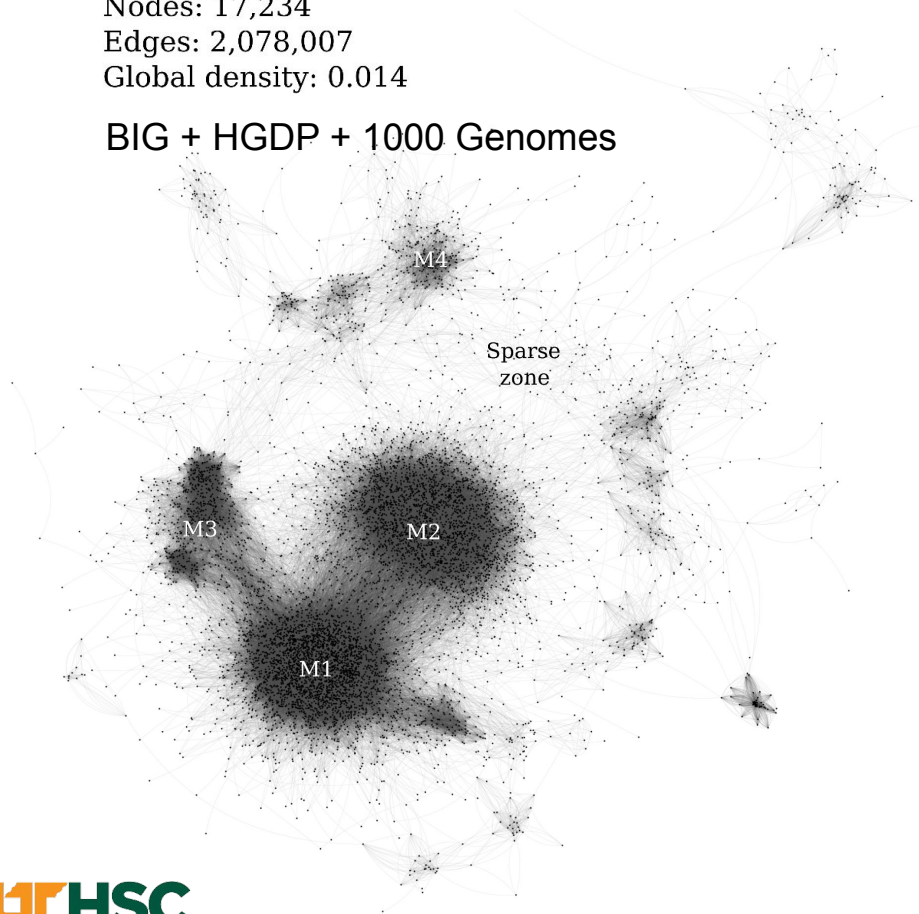BIG + HGDP + 1000 Genomes

# Leiden algorithm for community detection

IBD based network

Nodes: 17,234
Edges: 2,078,007
Global density: 0.014

BIG + HGDP + 1000 Genomes
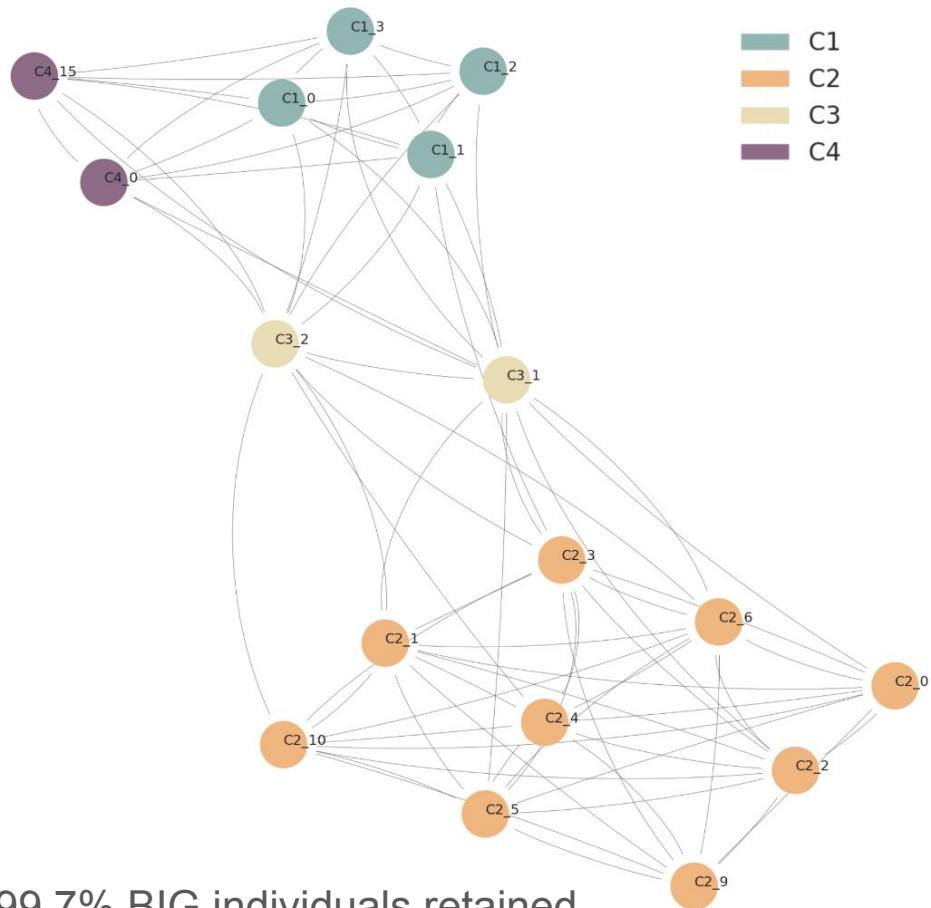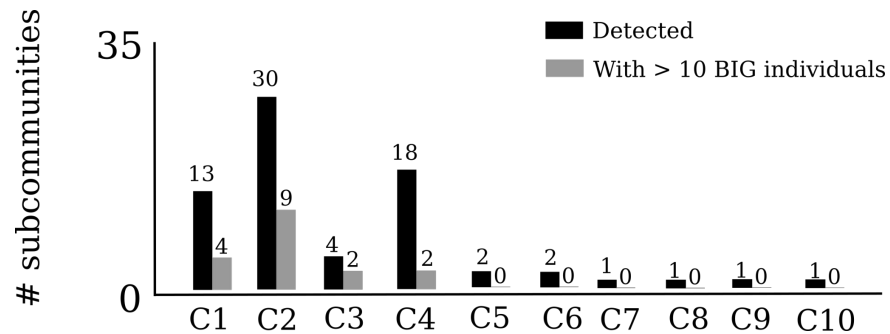
M4

Sparse zone

M3    M2

M1

Beyond ancestry: grouping people based on how much genome they actually share

EUR
EUR-AFR
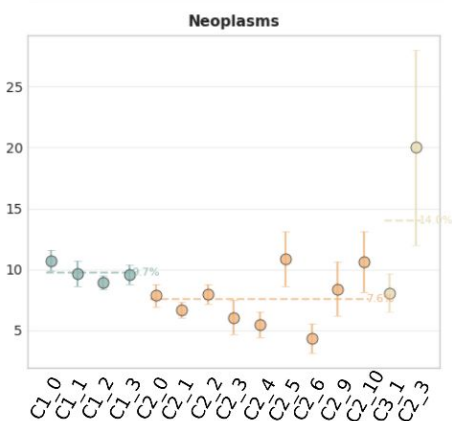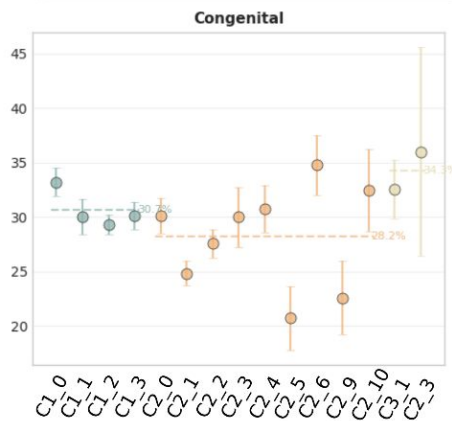AFR
EUR-AMR
AMR
EUR-EAS
EAS
Multiway

# IBD *versus* ancestry

99.7% BIG individuals retained

# Identity-by-descent captures shared environment



A - Genealogical and geographic distance



PM 2.5 EJ Index

Smoking prevalence

Can a similar approach be

useful on HPRCv2?

Preliminary/pilot results

Identity-by-descent (IBD) detection
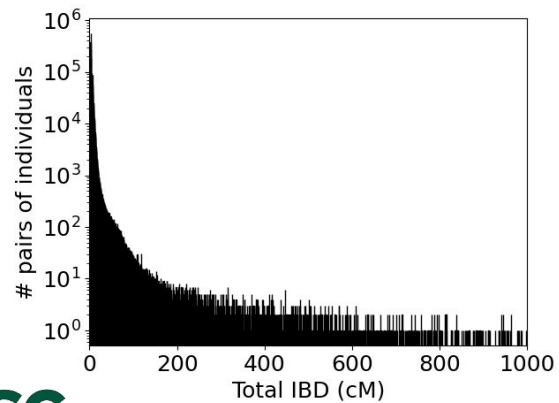
HGDP 1kGP (4,091)   HPRCv2 Samples (232)
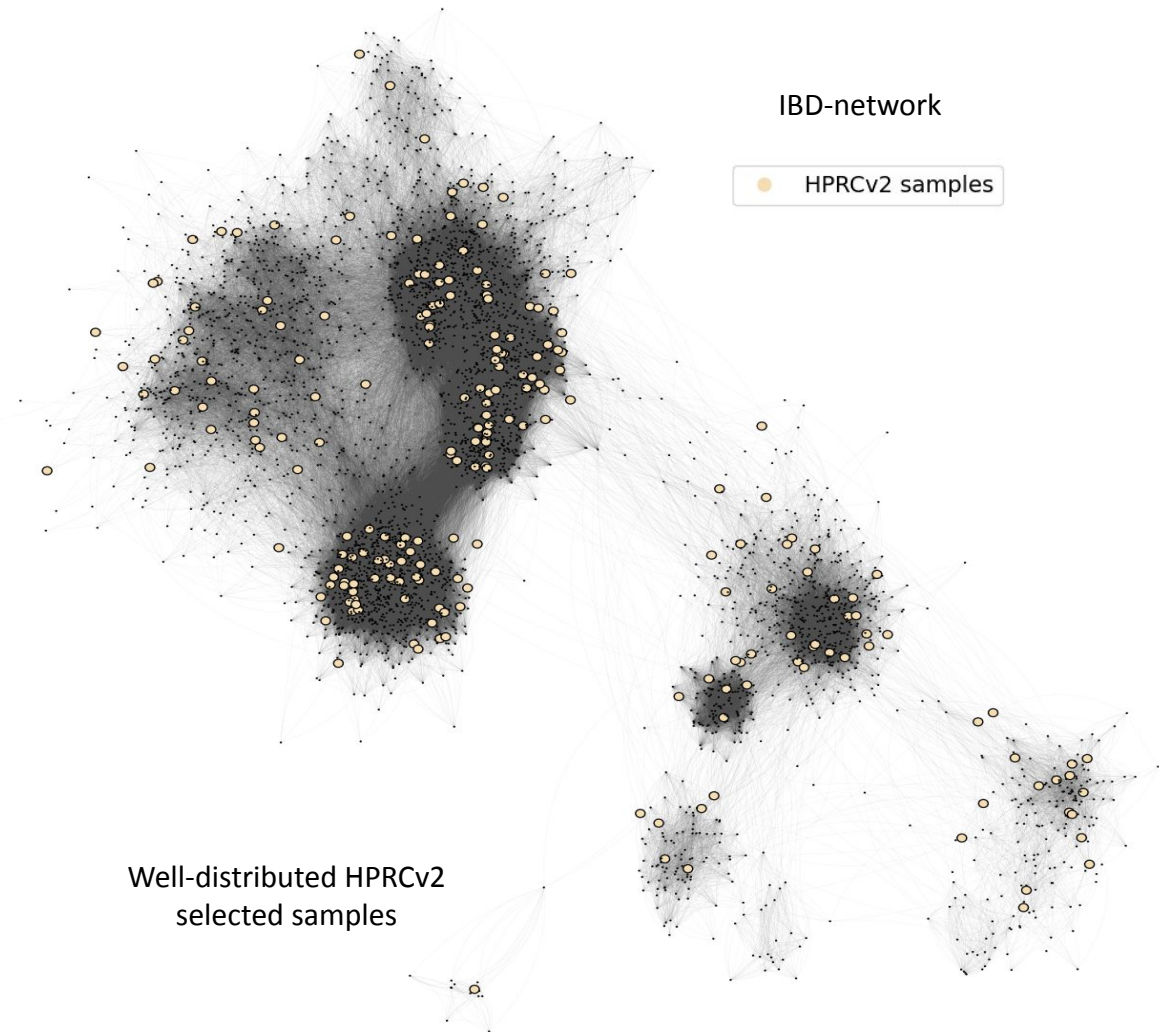
Labelled data (4,091)
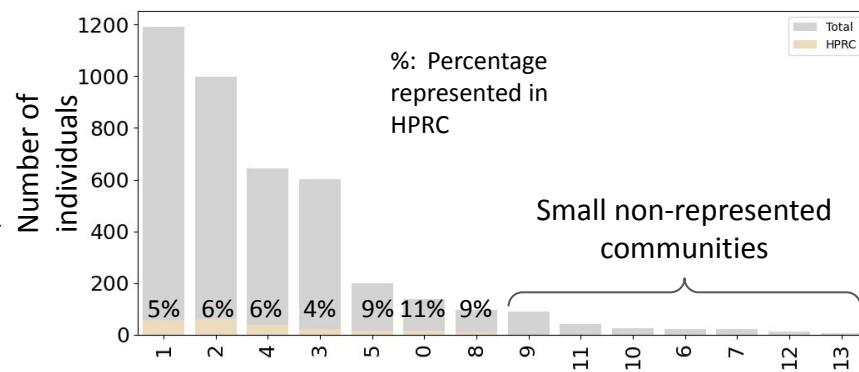
VCF QC
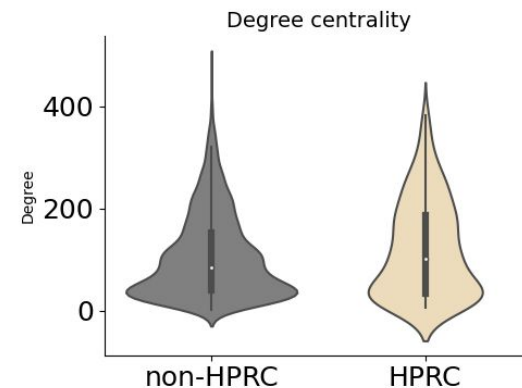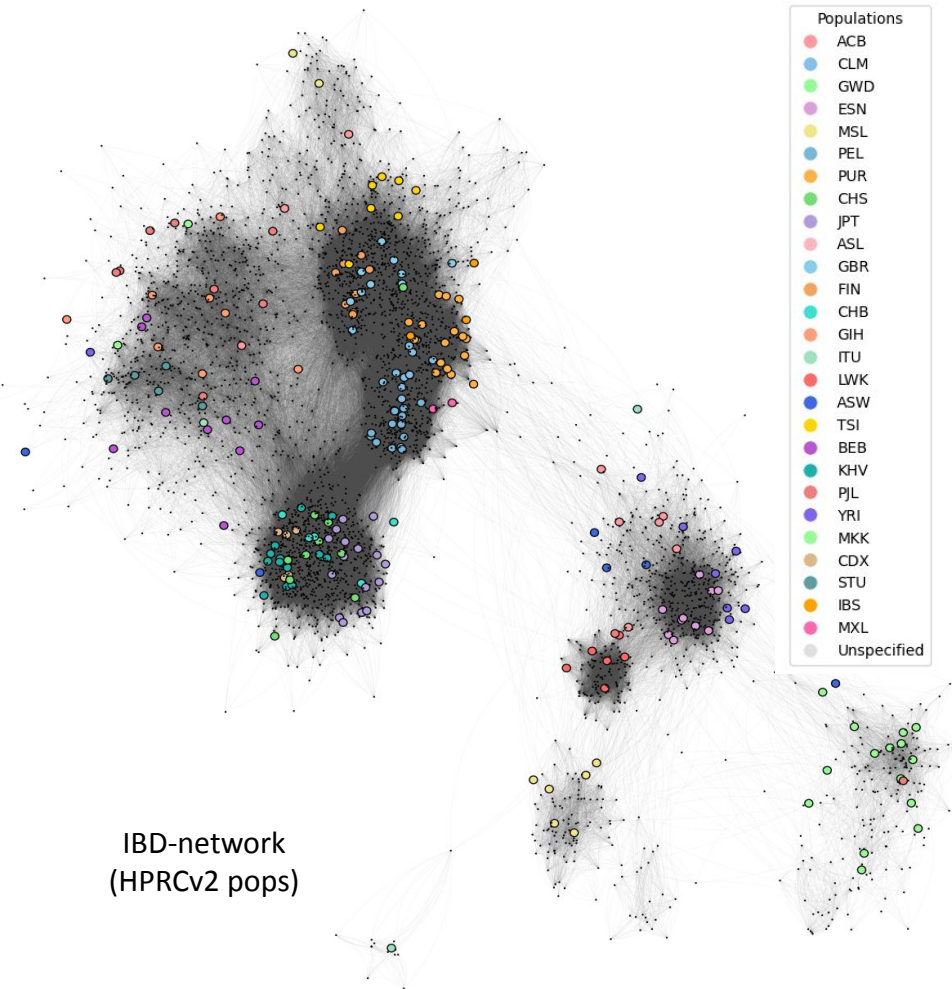
IBD segments

IBD QC

Total IBD per pair

IBD-network

HPRCv2 samples

Well-distributed HPRCv2 selected samples

Populations
- ACB
- CLM
- GWD
- ESN
- MSL
- PEL
- PUR
- CHS
- JPT
- ASL
- GBR
- FIN
- CHB
- GIH
- ITU
- LWK
- ASW
- TSI
- BEB
- KHV
- PJL
- YRI
- MKK
- CDX
- STU
- IBS
- MXL
- Unspecified

IBD-network
(HPRCv2 pops)

Degree centrality

non-HPRC    HPRC

%: Percentage represented in HPRC

Small non-represented communities

Some small HGDP-1kGP communities (detected by Leiden) are yet unexplored in HPRCv2
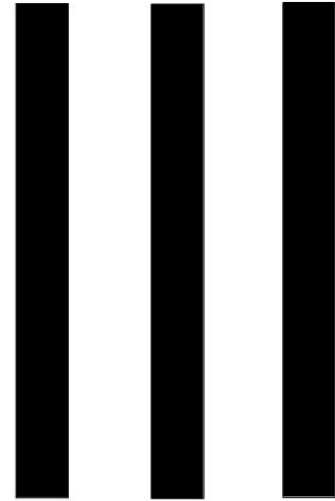
5%  6%  6%  4%  9%  11%  9%

UTHSC

# Other IBD applications in sample selection



**Haplotype informativeness**: Compute the amount of cM explained by IBD sharing between the selected samples and public resources (HGDP-1kGP).

**Proportion covered**: Compute the percentage (amount of cM) already explored by IBD sharing of a potential HGDP-1kGP sample to be incorporated to HPRCv2.

Thank you