

# Computational inference of genetic ancestry in cancer

Jenea Adams

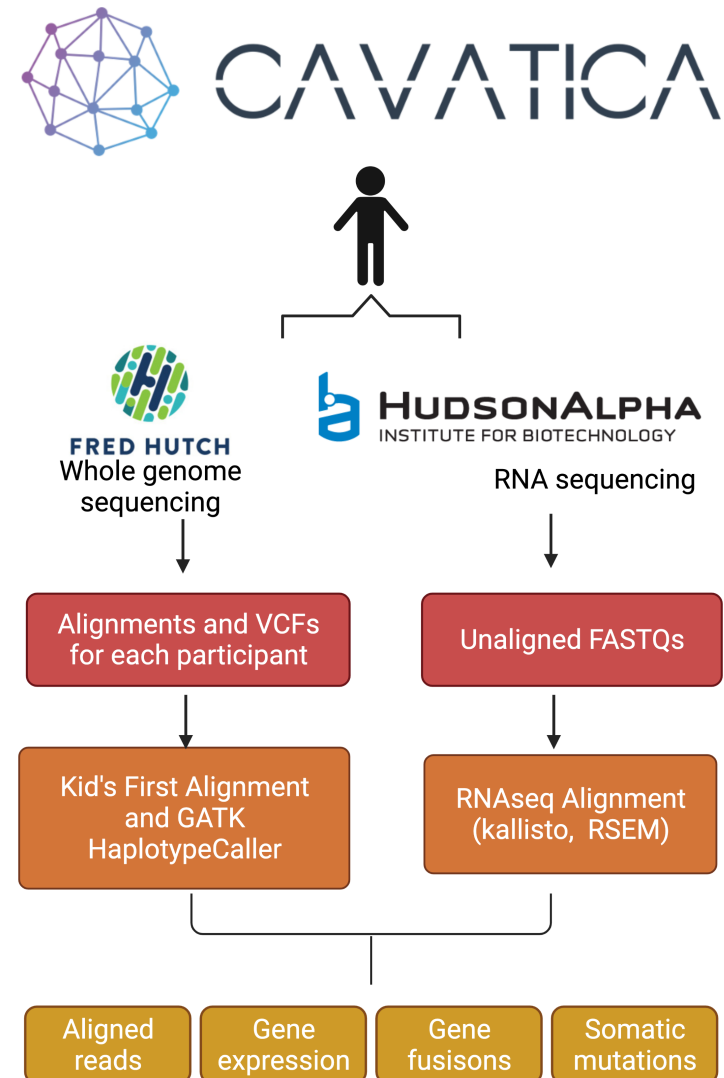
Xing Lab Roundtable

October 12, 2021

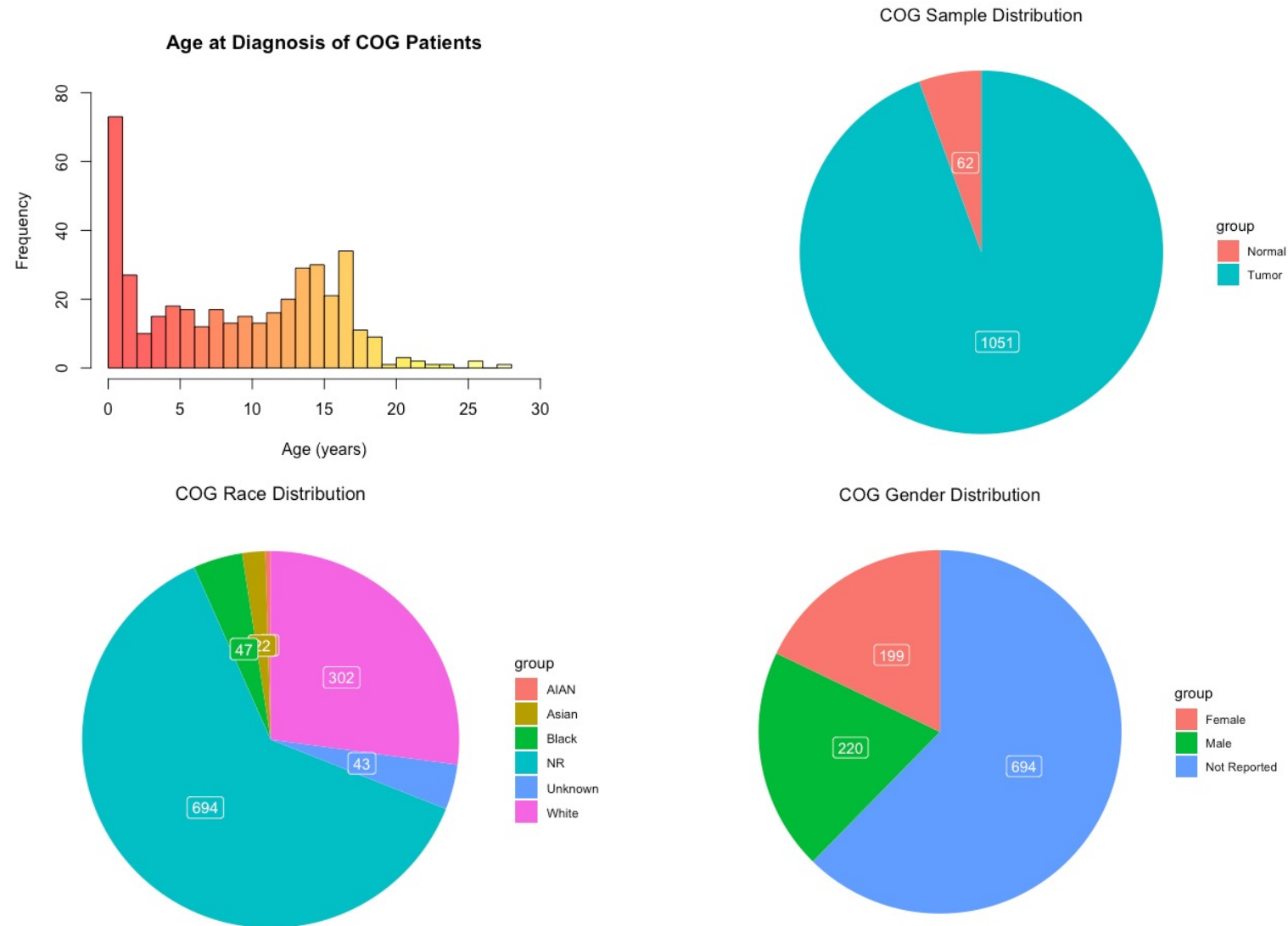
# CAVATICA: CHOP data-analysis platform of raw and harmonized multiomic data

## Overview of data used for this analysis

- De-novo AML, DS-AML, APL-AML
- Data from Children's Oncology Group Clinical Trial (No. AAML1031)
- 1,113 RNA-seq files (aligned and quantified)
  - Both kallisto and RSEM were used
  - Gene fusions also quantified
- 408 whole-genome sequencing (WGS) files

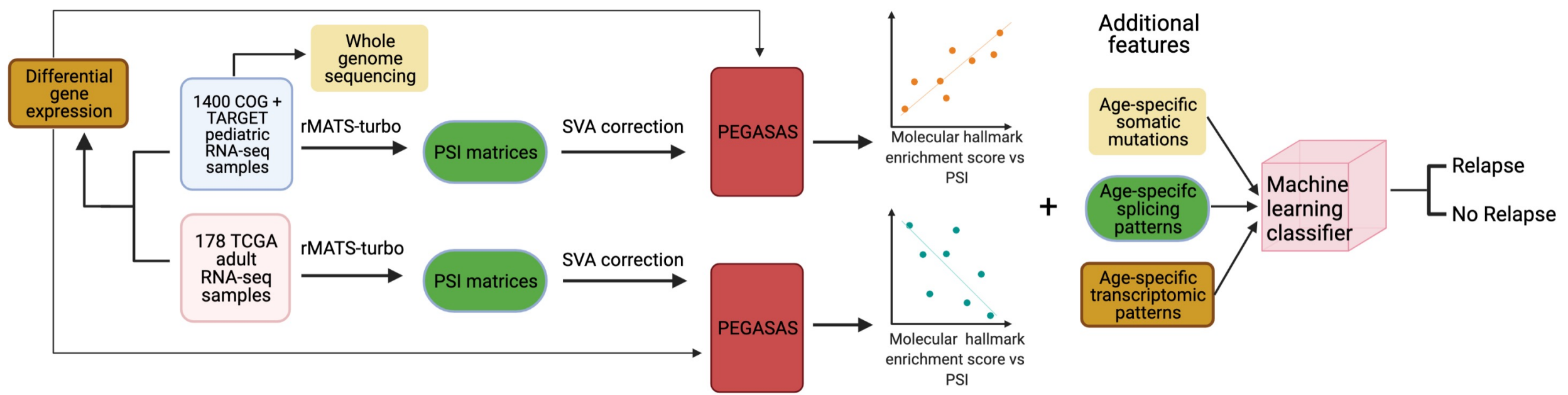


# CAVATICA allows easy metadata extraction



# In Progress

- Analysis: Sourcing metadata
  - **Treatment history of COG patients**
  - Use of whole genome seq data for ancestry association with molecular subtypes and disease progression
  - RNA-seq data for gene expression analysis of molecular subtypes



Genetic ancestry can be used to analyze factors affecting clinical outcomes

- Ancestries influence germline genetics
- Ancestries tend to carry different disease exposures
- Multi-omic features can contribute to the prevalence of relapse in pediatric AML, including ancestry

## Article

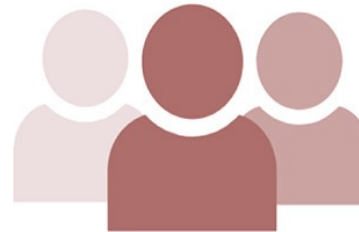
# Comprehensive Analysis of Genetic Ancestry and Its Molecular Correlates in Cancer

Jian Carrot-Zhang,<sup>1,2,3,24</sup> Nyasha Chambwe,<sup>4,24</sup> Jeffrey S. Damrauer,<sup>5,24</sup> Theo A. Knijnenburg,<sup>4,24</sup> A. Gordon Robertson,<sup>6,24</sup> Christina Yau,<sup>7,8,24</sup> Wanding Zhou,<sup>9,24</sup> Ashton C. Berger,<sup>1,2,24</sup> Kuan-lin Huang,<sup>10,24</sup> Justin Y. Newberg,<sup>11,24</sup> R. Jay Mashl,<sup>12,25</sup> Alessandro Romanel,<sup>13,25</sup> Rosalyn W. Sayaman,<sup>14,15,25</sup> Francesca Demichelis,<sup>13</sup> Ina Felau,<sup>16</sup> Garrett M. Frampton,<sup>11</sup> Seunghun Han,<sup>2,3</sup> Katherine A. Hoadley,<sup>5</sup> Anab Kemal,<sup>16</sup> Peter W. Laird,<sup>9</sup> Alexander J. Lazar,<sup>17</sup> Xiuning Le,<sup>18</sup> Ninad Oak,<sup>19,20</sup> Hui Shen,<sup>9</sup> Christopher K. Wong,<sup>21</sup> Jean C. Zenklusen,<sup>16</sup> Elad Ziv,<sup>14</sup> Cancer Genome Atlas Analysis Network, Andrew D. Cherniack,<sup>1,2,3,26,\*</sup> and Rameen Beroukhi<sup>1,2,3,22,23,\*</sup>

# Paper outline

- Determining a biological basis of health disparities in cancer
- Central question: *if you quantify ancestry, could you see differences in molecular characteristics from their tumors?*

10,678 cancer cases of diverse genetic ancestries



## Ancestry-differential molecular features



Somatic alterations



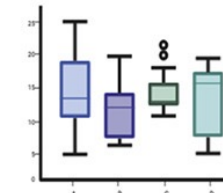
Methylation



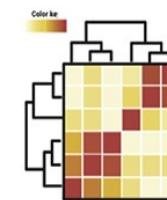
Gene expression



miRNA & expression



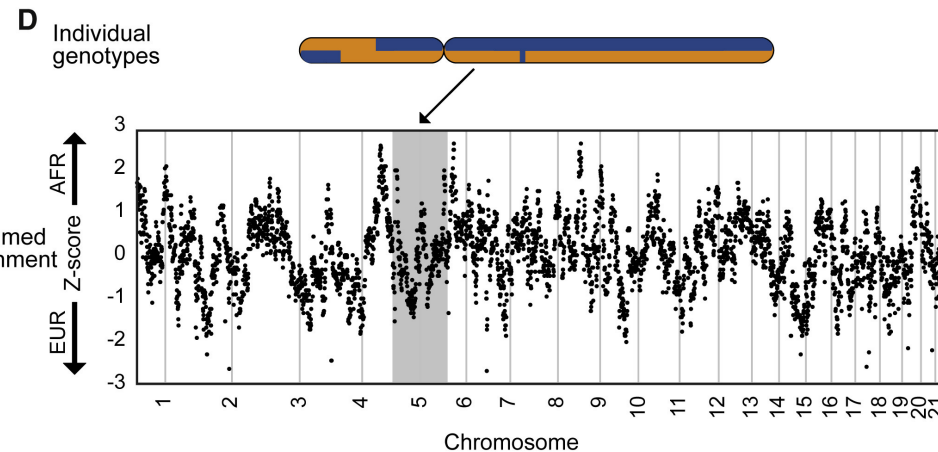
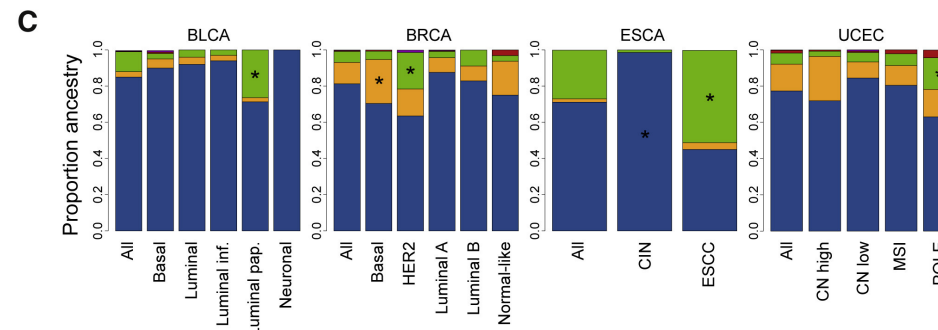
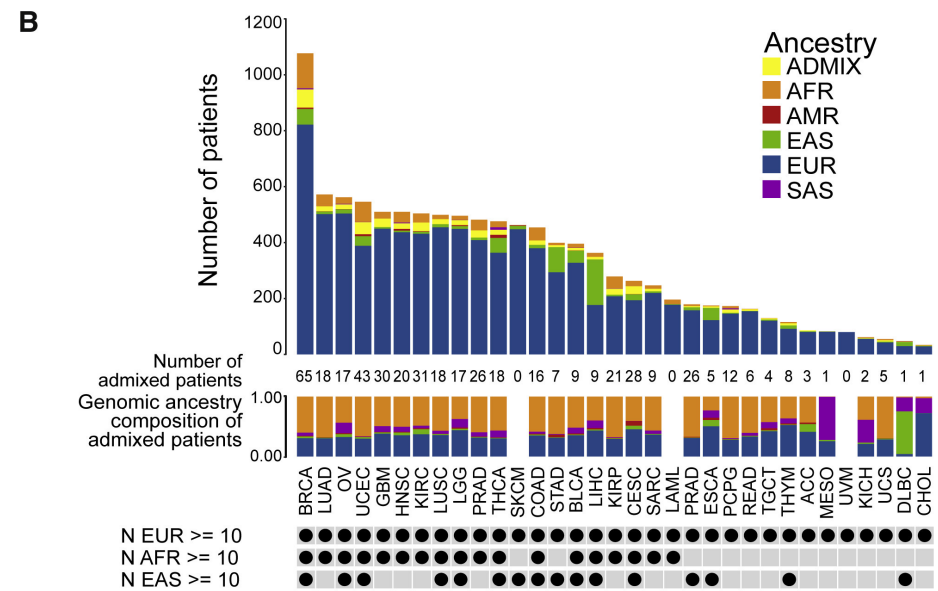
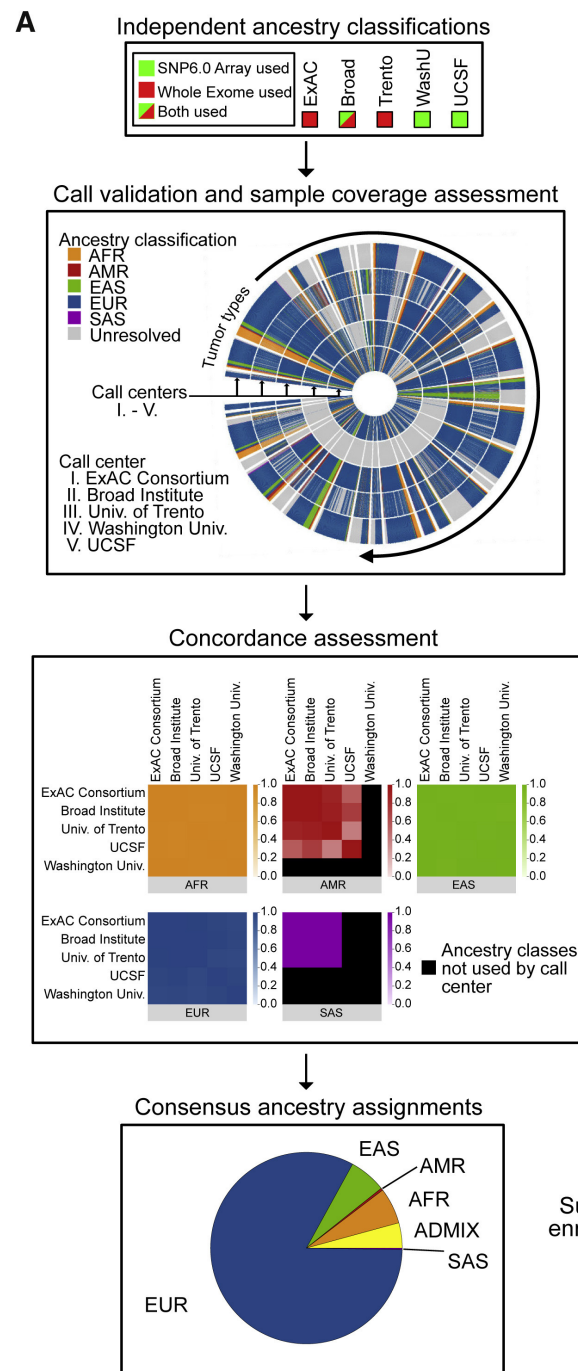
Expression QTLs



Cancer pathways and immunogenicity

Primary analyses

Integrative analyses





# (Some) methods used for ancestry inference

- Broad – SNP and Exome based calls
  - SNP → smartpca with AFR, EUR, EAS, AMR, SAS in first 3 PCs
  - Exome → from ExAC, PCA on 5400 SNPs to cluster continental ancestry
- WashU – SNP based calls
  - PCA in PLINK – use self-reported as reference within first 2 PCs
- UCSF – SNP based calls
  - PAM and k means clustering of PCs with PLINK (w/o LD pruning)

# Limitations and questions

- Reference populations
  - We need a diverse reference population (besides 1000 genomes)
    - What other sources exist?
- Most examples are from consortia-level science, what does this mean for this project?
- This paper concluded that biologically differences between ancestries were tissue-specific, not cancer specific