# Webinar Instructions

*Welcome to the Gabriella Miller Kids First Pediatric Research Program's Public Webinar!*

- Every participant is muted upon entry.

- To ask public questions, use the **Q&A** bar (right side of your screen). We encourage you to save these for the question period.

- You can ask also use the "chat" service to send private messages to the host or presenters throughout the webinar.

- After the webinar, additional program-related questions can be emailed to: KidsFirst@od.nih.gov.

**This webinar will be recorded.
We will start at 3pm (EDT)**

# May 18ᵗʰ Webinar Agenda

- **3:00pm** - Introduction
- **3:05pm** - Kids First Orofacial Cleft Project Findings
- **3:40pm** - Kids First Data Resource Center
  - New Portal Features
  - Cavatica: Cloud User Workspace Introduction
  - User Workspace Demonstration
  - Kids First DRC Roadmap
- **4:30pm** - Kids First Program & Collaboration Update
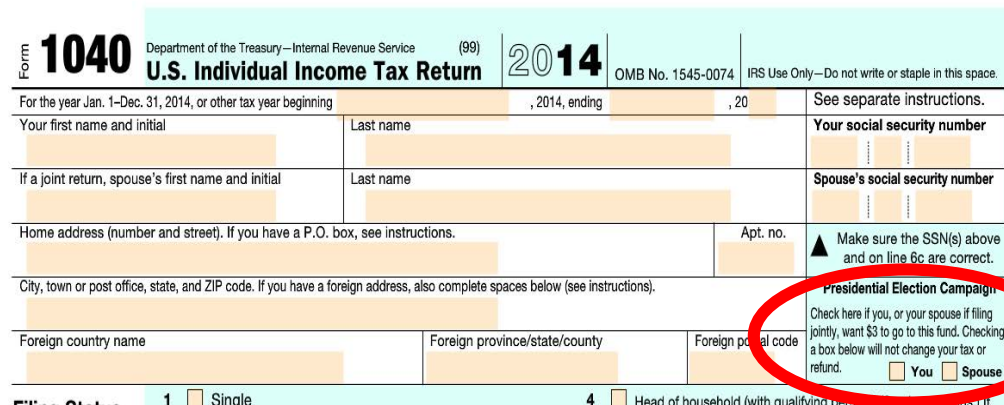- **4:50pm** - Questions and Answers

**Valerie Cotton**
Kids First Program Manager
*Eunice Kennedy Shriver* National Institute of
Child Health and Human Development (NICHD)

# *How did Kids First get started?*

- Initiated in response to the [2014 Gabriella Miller Kids First Research Act](#):
  - Signed into law on April 3, 2014
  - Ended taxpayer contribution to presidential nominating conventions
  - Transferred $126 million into the Pediatric Research Initiative Fund
  - Authorized appropriation of $12.6 million per year for 10 years to the NIH Common Fund for pediatric research; first appropriation was for FY2015
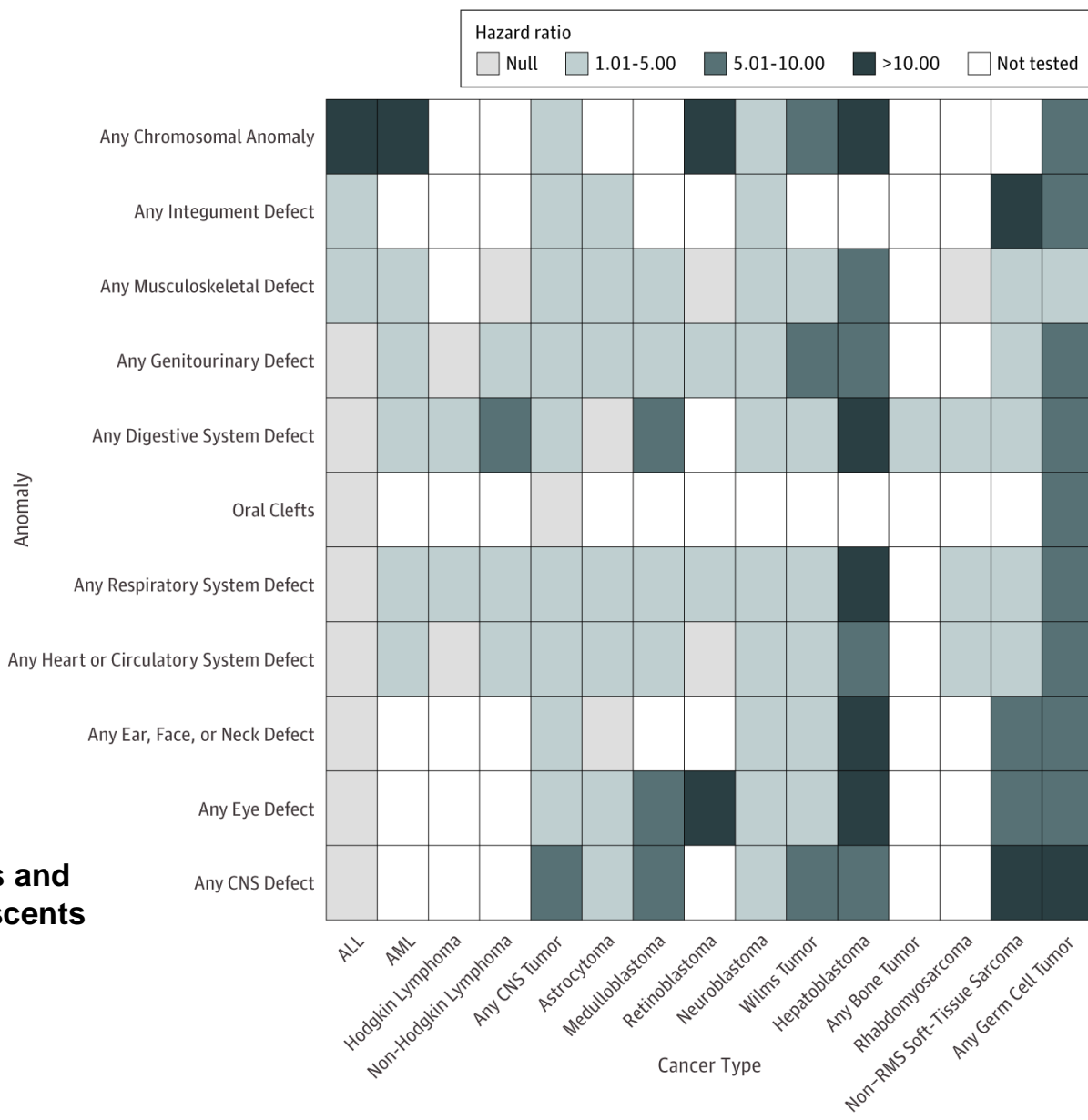
# Vision



Alleviate suffering from childhood cancer and structural birth defects by fostering **collaborative research** to uncover the etiology of these diseases and supporting **data sharing** within the pediatric research community.

# *Why study childhood cancer & structural birth defects together?*

*Birth defects associated with increased risk of cancer among children*

From: **Association Between Birth Defects and Cancer Risk Among Children and Adolescents in a Population-Based Assessment of 10 Million Live Births**

Hazard ratio
- Null
- 1.01-5.00
- 5.01-10.00
- >10.00
- Not tested

Anomaly (rows):
- Any Chromosomal Anomaly
- Any Integument Defect
- Any Musculoskeletal Defect
- Any Genitourinary Defect
- Any Digestive System Defect
- Oral Clefts
- Any Respiratory System Defect
- Any Heart or Circulatory System Defect
- Any Ear, Face, or Neck Defect
- Any Eye Defect
- Any CNS Defect

Cancer Type (columns):
ALL, AML, Hodgkin Lymphoma, Non-Hodgkin Lymphoma, Any CNS Tumor, Astrocytoma, Medulloblastoma, Retinoblastoma, Neuroblastoma, Wilms Tumor, Hepatoblastoma, Any Bone Tumor, Rhabdomyosarcoma, Non-RMS Soft-Tissue Sarcoma, Any Germ Cell Tumor

# NIH Kids First Working Group

Kids First is an NIH Common Fund program coordinated by a **trans-NIH Working Group**, which is chaired by four institutes:

> *Eunice Kennedy Shriver* National Institute of Child Health and Human Development (**NICHD**)

> National Human Genome Research Institute (**NHGRI**)

> National Heart, Lung, and Blood Institute (**NHLBI**)

> National Cancer Institute (**NCI**)

Other Working Group Representation:

**NIDCR**    **NIAAA**    **NIDDK**    **NEI**    **NIAID**    **ORIP**

**NIDA**    **NINDS**    **NIEHS**    **NIAMS**    **NCATS**    **CDC**

# Kids First Major Initiatives

**Through 2021:**

1. Identify & sequence cohorts of children with **childhood cancer and/or structural birth defects**.
2. Build the **Gabriella Miller Kids First Data Resource** to empower discovery

# The Kids First Dataset is Growing!

**39 projects | 37,000 genomes | 15,000 cases | 10 released datasets**



- Disorders of Sex Development
- Congenital Diaphragmatic Hernia
- Ewing Sarcoma
- Structural Heart & Other Defects
- Syndromic Cranial Dysinnervation Disorders
- Cancer Susceptibility
- Adolescent Idiopathic Scoliosis
- Neuroblastomas
- Enchondromatoses
- Orofacial Clefts in Caucasian, Latin American, Asian & African, Filipino populations
- Osteosarcoma
- Familial Leukemia
- Craniofacial Microsomia
- Hemangiomas, Vascular Anomalies & Overgrowth
- Nonsyndromic Craniosynostosis
- Patients with both childhood cancer and birth defects
- Kidney and Urinary Tract Defects

- Microtia
- Hearing Loss
- Bladder Exstrophy
- Cornelia de Lange Syndrome
- Intracranial & Extracranial Germ Cell Tumors
- Esophageal Atresia and Tracheoesophageal Fistulas
- Fetal Alcohol Spectrum Disorders
- Myeloid Malignancies + overlap with Down syndrome
- Congenital Heart Defects and Acute Lymphoblastic Leukemia in Children with Down Syndrome
- Structural Brain Defects
- Structural Defects of the Neural Tube (Spina Bifida: Myelomeningocele)
- CHARGE Syndrome
- Laterality Birth Defects
- T-cell Acute Lymphoblastic Leukemia
- Pediatric Rhabdomyosarcoma

# *More researchers are accessing Kids First data!*



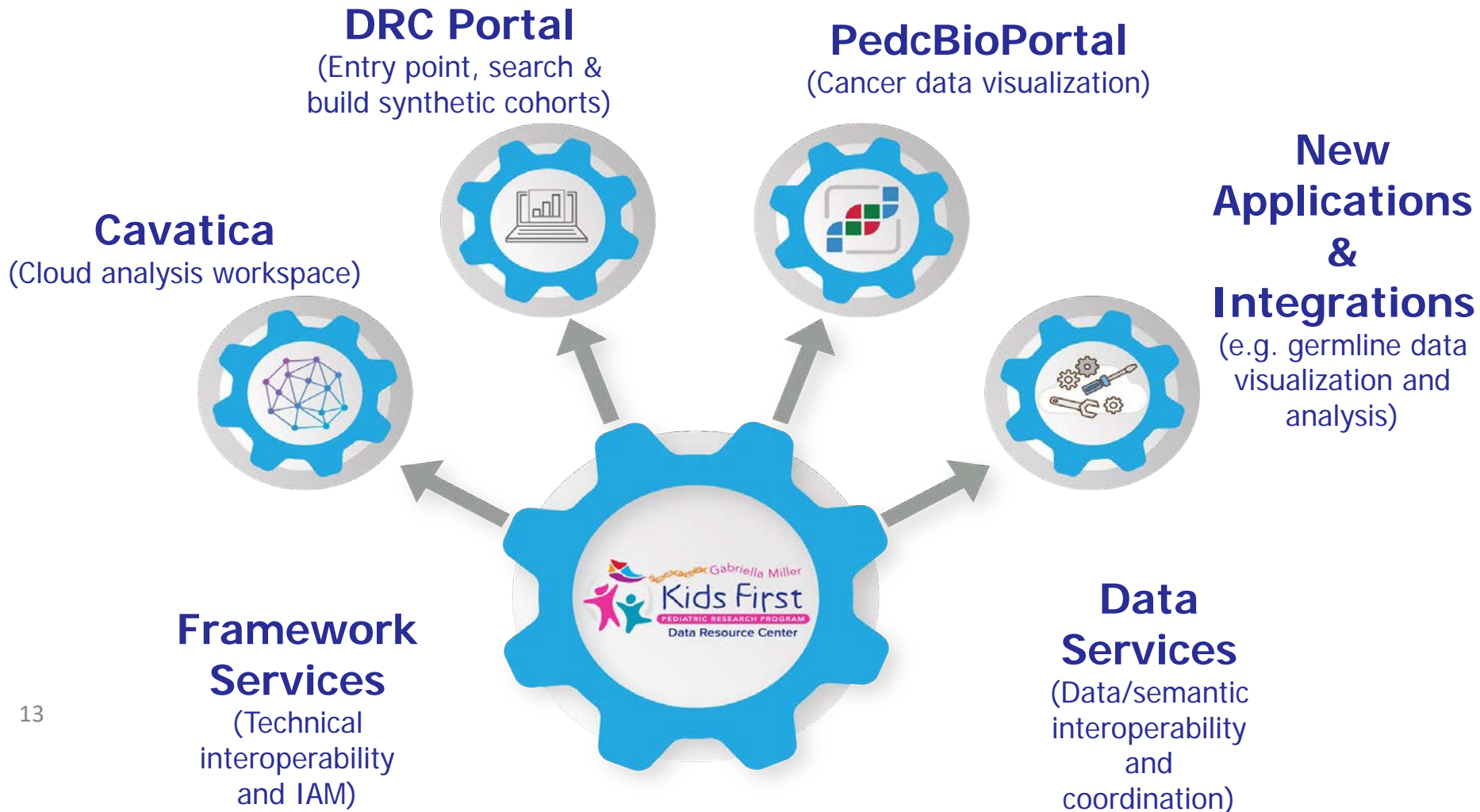**~1500 registered users since 2018 launch**

**Individual-level sequence data**
**>100** Data Access Requests approved by the Kids First Data Access Committee across **10** Kids First genomic datasets available

**NIH Kids First Data Access Committee**

dbGaP
GENOTYPES and PHENOTYPES

# The Kids First Data Resource for Collaborative Discovery



**DRC Portal**
(Entry point, search & build synthetic cohorts)

**PedcBioPortal**
(Cancer data visualization)

**Cavatica**
(Cloud analysis workspace)

**New Applications & Integrations**
(e.g. germline data visualization and analysis)

**Framework Services**
(Technical interoperability and IAM)

**Data Services**
(Data/semantic interoperability and coordination)

13

# Kids First X01 Investigators: Orofacial Clefts

**Mary Marazita, PhD**
University of Pittsburgh

**Eleanor Feingold, PhD**
University of Pittsburgh

**Elizabeth Leslie, PhD**
Emory University

**Harrison Brand, PhD**
Broad Institute

# Gabriel Miller Kids First: Orofacial Cleft (OFC) Studies

**Mary L. Marazita, Ph.D.; Eleanor Feingold, Ph.D.**

**and GMKF OFC team**

**Center for Craniofacial and Dental Genetics**
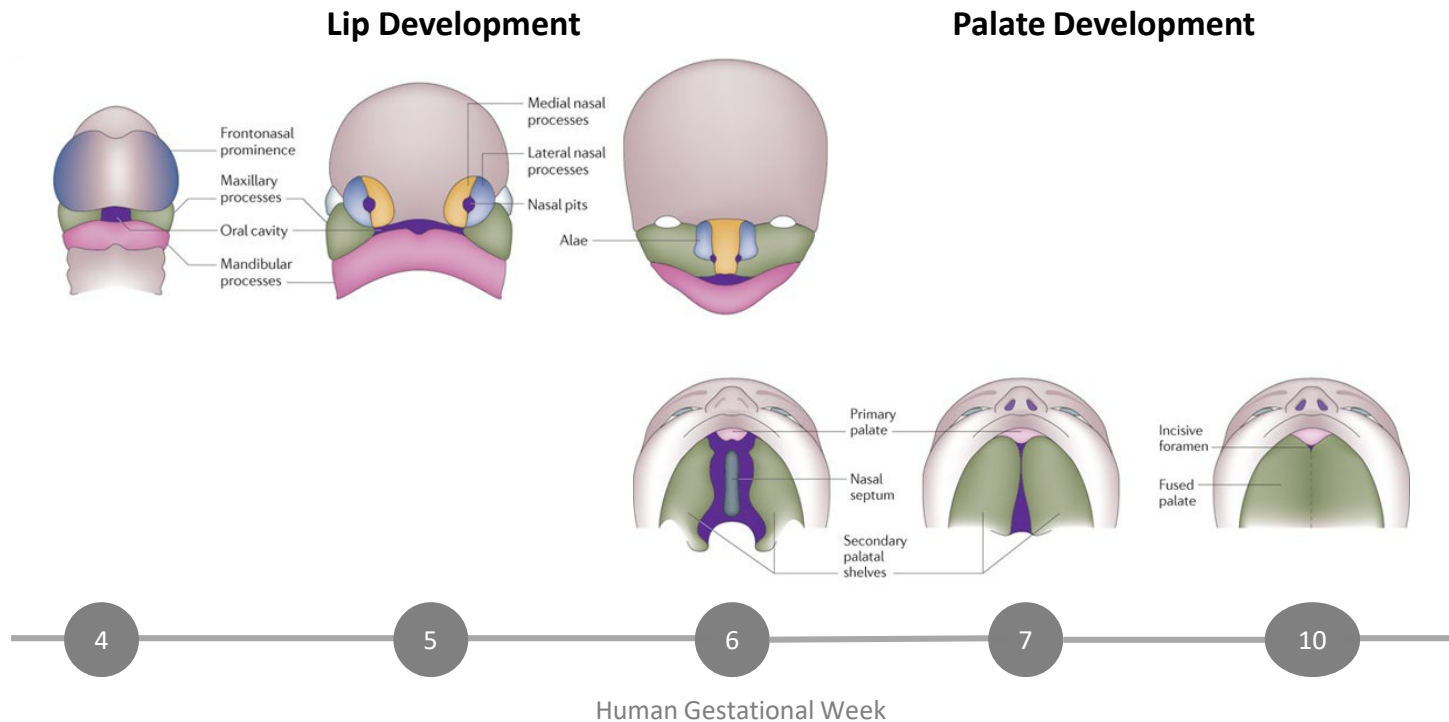
**University of Pittsburgh**

Kids First 2020 Spring Public Webinar
May 18, 2020

Goal is to elucidate the genetic basis of OFC, one of the most common structural birth defects in humans worldwide



*operationsmile.org*

# Lip and Palate Development



**Lip Development**                    **Palate Development**

Human Gestational Week

4   5   6   7   10

Dixon, Marazita, Beaty (2010), Nature Reviews Genetics

# Orofacial Clefts (OFCs), Sub-Phenotypes



CL/P                     CP

Cleft Lip (CL)    Cleft Lip and     Cleft Palate (CP)
                  Palate (CLP)

Jugessur et al. (2009),
Oral Diseases

# MULTI-ETHNIC, DEEP PHENOTYPING AND GENETIC RESOURCES



Houston, TX
Denver, CO
Iowa City, IA
St. Louis, MO
Puerto Rico
Guatemala
Colombia
Brazil
Argentina
Denmark
Spain
Hungary
India
China
Australia
Philippines
Nigeria

And more

# Deep Phenotyping Protocol

- **Questionnaires**
  - Demographics, personal and family medical history (eg birth defects, cancer, systems review), pregnancy history, developmental milestones

- **Physical examination**
  - Facial birth defects, Height/weight/head circumference, oral exam (plus imaging), limbs, hands/digits, speech sample for VPI perceptual screening

- **Imaging (to derive traits such as facial measurements)**
  - 3d facial image, intraoral photograph, palate video during speech (smCP), hand scans, ultrasound of upper lip region (OOM), dental casts plus 3d, upper and lower lip photos (lip pits)

# GMKF: OFC Whole Genome Sequencing

**APPROVED (total = 1,413 OFC trios to date):**

- **European descent:  447**

- **Latin American (Colombia) :  265**

- **African (Nigeria and Ghana) :  137 trios**

- **Asian (Taiwan) : 124 trios**

- **Asian (Philippines) : 373 trios (COVID-19 delay)**


- **In review: additional Latin American trios**

Nandita Mukhopadhyay

# AFRICAN AND ASIAN ANALYSES

(Butaliand Baneatalyy)s

**& 124 Taiwanese trios**

**gTDT analysis of 130 African trios**

**Chr8q.24 signal replicated in African TDT and at the same E-06 as was reported in the African only GWAS for CL/P**



Asian gTDT Analysis for CLP

African gTDT Analysis for CL/P

**MANY THANKS TO OUR PARTICIPANTS WORLDWIDE**

**U of Pittsburgh**:
Mary L. Marazita
Seth M. Weinberg
Eleanor Feingold
Nandita
  Mukhopadhyay
Ross Long
(Lancaster)

**Emory U**:
Elizabeth J. Leslie
Madison Bishop
Pankaj Chopra
Michael Mortillo
Dave Cutler
Michael Epstein

**U of Iowa**:
Jeffrey C. Murray
Azeez Butali
Lina M. Moreno
Luz Consuelo Valencia-Ramirez
George L. Wehby
Andrew Lidral

**Johns Hopkins University:**
Terri Beaty
Ingo Ruczinski
Margaret Taub
Allen Scott
Jackie Hetmanski
Debashree  Ray

**Taiwan**:
YH Wu-Chou
PK Chen

**Africa:**
NHGRI: Adebowale Adeyemo (NHGRI)
Kwame Nkurumah (Ghana)
Lord J.J Gowans (Ghana)
Lanre W Adeyemo (Nigeria)
Peter Mossey (U of Dundee)

**Other**:
Harrison Brand (Harvard)
Jacqueline T. Hecht (U of TX)
Frederic Deleyiannis (U of CO)
Carmencita Padilla (U of Manila)
Mauricio Arcos-Burgos
Andrew Czeizel
Eduardo Castilla
Ieda Orioli
Fernando Poletta

# Uncovering the Genome-Wide Contribution of *De Novo* Mutations in Orofacial Clefts

Elizabeth J. Leslie, PhD

Department of Human Genetics

Emory University

May 2020 Kids First Public Webinar

# *De novo* mutations

*de novo* mutations per genome:
~70-90 single nucleotide variants
~6 insertion/deletions
0.02 copy number variants

~1 *de novo* mutation per exome

1 out of 20,563 protein-coding genes are hit per generation

# *De novo* mutations are common causes of congenital and developmental anomalies



**Autism**
Neale et al., 2012; Sanders et al., 2012;
De Rubeis et al., 2014; Satterstrom et al., 2018

**Craniosynostosis**
Timberlake et al., 2017, 2018

**Congenital Heart Disease**
Homsy et al., 2015; Jin et al., 2017; Watkins et al., 2019

**Neural Tube Defects**
Lemay et al., 2015

**Congenital Diaphragmatic Hernia**
Yu et al., 2014; Qi et al., 2018

# What is the role of de novo mutations in OFCs?

Madison Bishop, PhD

**1** Are coding DNMs enriched in OFC cases?

**2** What is the biological relevance of DNMs in OFC cases?

**3** What is the clinical significance of DNMs in OFC cases?

# OFC de novos: by the numbers

- **756 trios (US/European, Colombian, and Taiwanese)**
  - 80 cleft lip only
  - 618 cleft lip and palate
  - 58 cleft palate only

- **73,027 DNMs genome-wide**

- **862 coding DNMs**

# Protein-Altering DNMs are enriched in OFCs

# DNMs are enriched in developmental genes

# Craniofacial-Specific Annotations?



Mesenchyme Clusters
(M1–M8)

Ectoderm
Clusters
(E1–E11)

Li et al., *Development*, 2019

# Excess of DNMs in genes expressed at point of fusion



Nasal process fusion zone

Olfactory epithelium

Li et al., *Development*, 2019

# De novo mutations in *IRF6*, *TFAP2A*, and *ZFHX4* are associated with OFCs

# Excess of DNMs in cranial neural crest genes

# Excess of DNMs in cranial neural crest genes

# DNMs in SOX2-interactome

**29 genes** with Loss of Function DNMs +
pLI > 0.95 + top 20% hNCC expression

**8 genes interact with SOX2 (FDR p = 9.5 x 10$^{-4}$)**
MACF1, RBM15, SETD2, CHD7, CTNND1, ZFHX4, IRF2BP1,TFAP2A

**126 genes** with Protein-Altering DNMs +
pLI > 0.95 + top 20% hNCC expression

**16 genes interact with SOX2 (FDR p = 5.1 x 10$^{-5}$)**
TCF20, RFX1, PPP2R5D, MDN1, NFIA, SPEN, NIPBL, ZNF292

# Towards a clinical gene panel for OFCs?

Curated an OFC gene list
containing 289 genes

- PubMed
- OMIM
- ClinVar
- Existing gene panels
  - NCBI Genetic Testing Registry: 6 genes
  - Fulgent: 24 genes
  - Prevention Genetics: 172 genes
  - Blueprint Genetics: 22 genes

# DNMs are enriched in clinically-relevant genes

# Summary

- We identified an excess of protein-altering DNMs in OFC trios (~1.2x more than expected)

- DNMs are in biologically relevant genes:
  - marker genes expressed in cells at the point of fusion of developing lip
  - genes in the top 20% in human cranial neural crest cells that are constrained to mutation

- DNMs are found in clinically relevant genes (~18x more than expected in AD OFC genes)

- 3 genes (*IRF6*, *TFAP2A*, *ZFHX4*) had individual excesses of DNMs

- ZFHX4 is a novel gene for OFCs

# Acknowledgments

Coding *de novo* mutations identified by WGS reveal novel orofacial cleft genes

bioRxiv

Madison R. Bishop, Kimberly Diaz Perez, Miranda Sun, Samantha Ho, Pankaj Chopra, Nandita Mukhopadhyay, Jacqueline B. Hetmanski, Margaret A. Taub, Lina M. Moreno-Uribe, Luz Consuelo Valencia-Ramirez, Claudia P. Restrepo Muñeton, George Wehby, Jacqueline T. Hecht, Frederic Deleyiannis, Seth M. Weinberg, Yah Huei Wu-Chou, Philip K. Chen, Harrison Brand, Michael P. Epstein, Ingo Ruczinski, Jeffrey C. Murray, Terri H. Beaty, Eleanor Feingold, Robert J. Lipinski, David J. Cutler, Mary L. Marazita, Elizabeth J. Leslie

bioRxiv 2020.04.01.019927; doi: https://doi.org/10.1101/2020.04.01.019927

The Leslie Lab

- Grace Carlock
- **Kim Diaz-Perez**
- Courtney Willett
- Dan Chang
- **Madison Bishop, PhD**

- Sarah Curtis, PhD
- Kelly Manning
- **Samantha Ho**
- Sydney Chung
- Shade Awoniyi

others

Gabriella Miller Kids First
PEDIATRIC RESEARCH PROGRAM

# SV Calling in Orofacial Clefts

Harrison Brand
Assistant Professor
MGH, Harvard Medical School, & Broad Institute
May 18th, 2020

# Introduction

- Impact of structural variation (SV) in non-syndromic forms of orofacial clefts (OFC) is largely uncharacterized

- We applied GATK-SV, our computational SV discovery pipeline, to 2,746 WGS samples that passed quality control
  - Includes 837 complete trios for *de novo* SV analysis

- OFC samples from 3 GMKF studies (140543, 136465, 132377) and 4 distinct populations (African, Asian, Latino, Caucasian)

- GATK-SV recently applied to 14,891 individuals in the gnomAD reference database (Collins*, Brand*, *et al. Nature,* in press)

CENTER FOR GENOMIC MEDICINE

# Structural Variation Background

# STRUCTURAL VARIATION

## Four basic classes of SV in the human genome



DELETION

DUPLICATION

INVERSION

INSERTION

# COMPLEX SVs

Complex SVs are comprised of combinations of the four basic SV classes

Paired-duplication inversion (dupINVdup)

Paired-deletion inversion (delINVdel)



Brand *et al.*, *Am. J. Hum. Genet.* (2014 & 2015)

Collins, Brand *et al.*, *Genome Biology* (2017)

# ABUNDANCE OF COMPLEX SVS IN THE GENOME

## Complex SVs are surprisingly abundant in the genome



Complex SV counts in the gnomAD SV cohort

# METHODS

# SV Discovery in Whole Genome Sequencing (WGS)

## Different classes of SVs leave distinct signatures in Illumina WGS data



Modified from:
Tattini *et al.*, *Front. Bioeng. Biotechnol.* (2015)
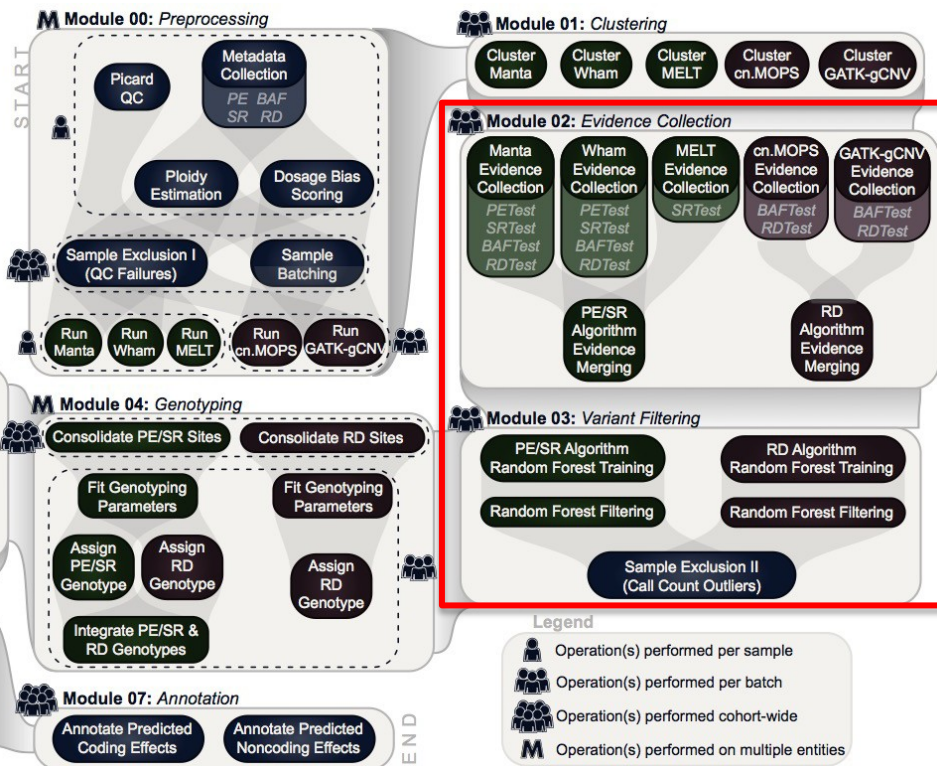
# GATK-SV: Cloud Enabled SV Pipeline

- Run several unfiltered algorithms to **maximize sensitivity**

- Re-evaluate evidence directly from BAMs to improve specificity

- Captures both unbalanced (CNV) and balanced (inversion, translocation) SV

- Integrates SV signatures to resolve complex events

- Has been adapted to work on Google Cloud via Broad Institute's Terra Platform
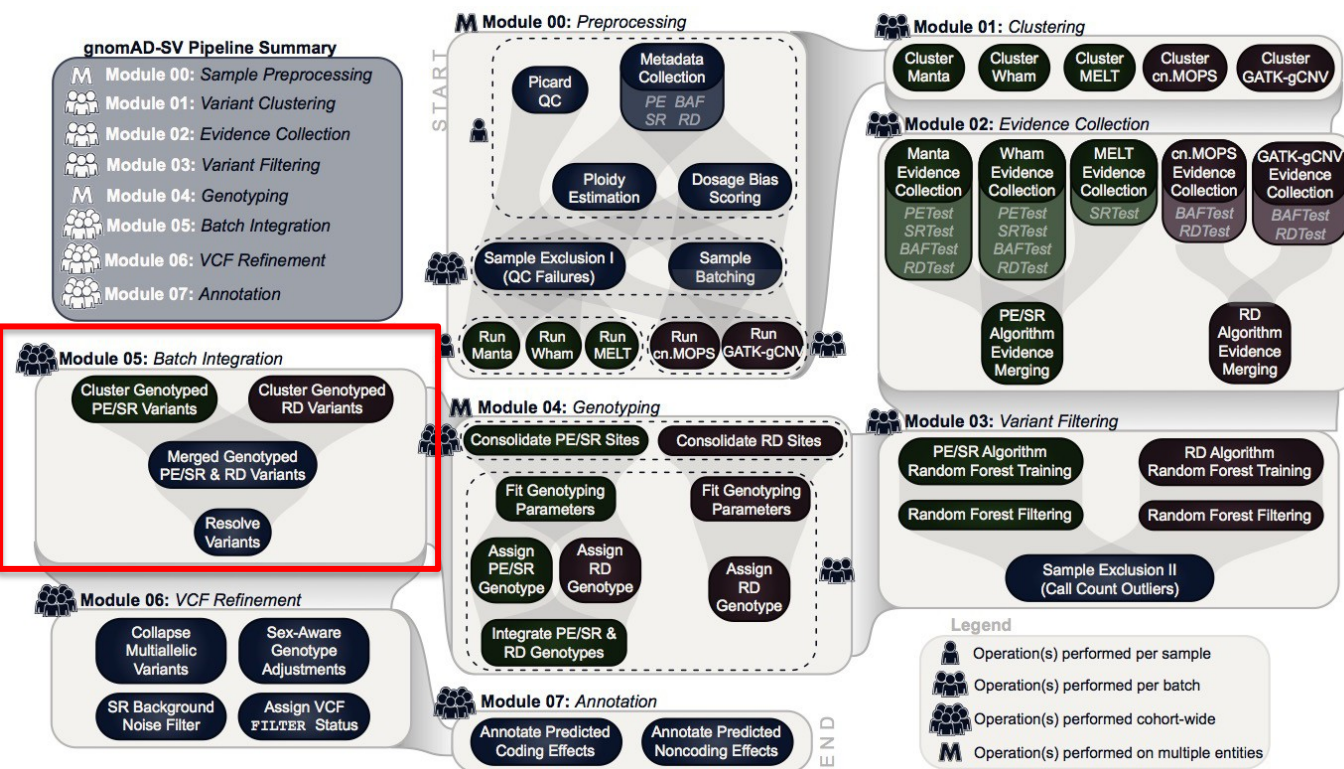
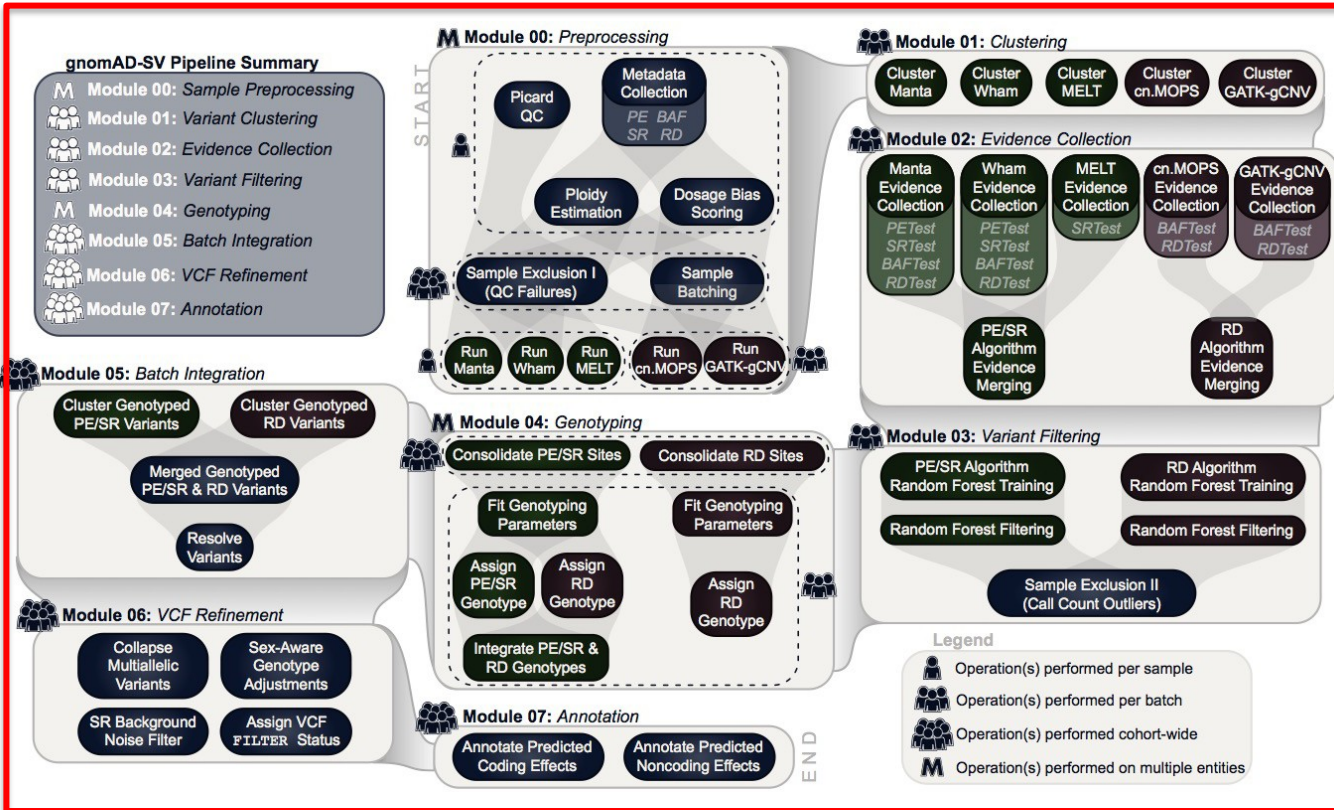# GATK-SV: CLOUD ENABLED SV PIPELINE

- Run several unfiltered algorithms to maximize sensitivity

- **Re-evaluate evidence directly from BAMs to improve specificity**

- Captures both unbalanced (CNV) and balanced (inversion, translocation) SV

- Integrates SV signatures to resolve complex events

- Has been adapted to work on Google Cloud via Broad Institute's Terra Platform
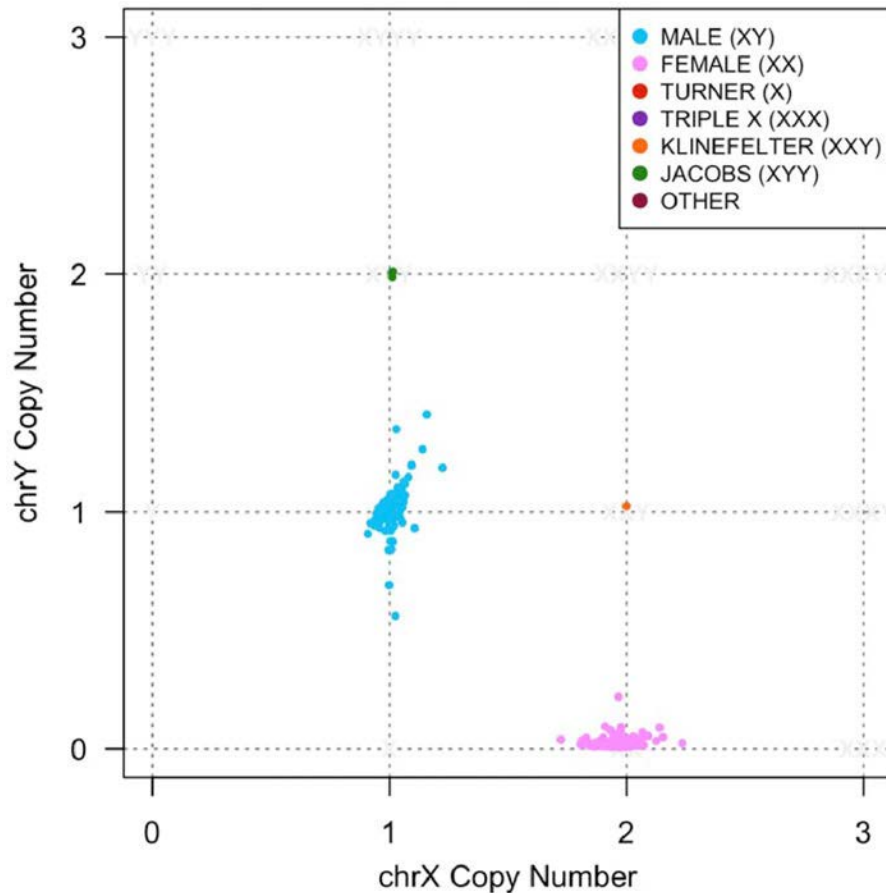
- Run several unfiltered algorithms to maximize sensitivity

- Re-evaluate evidence directly from BAMs to improve specificity

- Captures both unbalanced (CNV) and balanced (inversion, translocation) SV

- Integrates SV signatures to resolve complex events

- Has been adapted to work on Google Cloud via Broad Institute's Terra Platform

# GATK-SV; Cloud Enabled SV Pipeline



- Run several unfiltered algorithms to maximize sensitivity

- Re-evaluate evidence directly from BAMs to improve specificity

- Captures both unbalanced (CNV) and balanced (inversion, translocation) SV

- Integrates SV signatures to resolve complex events

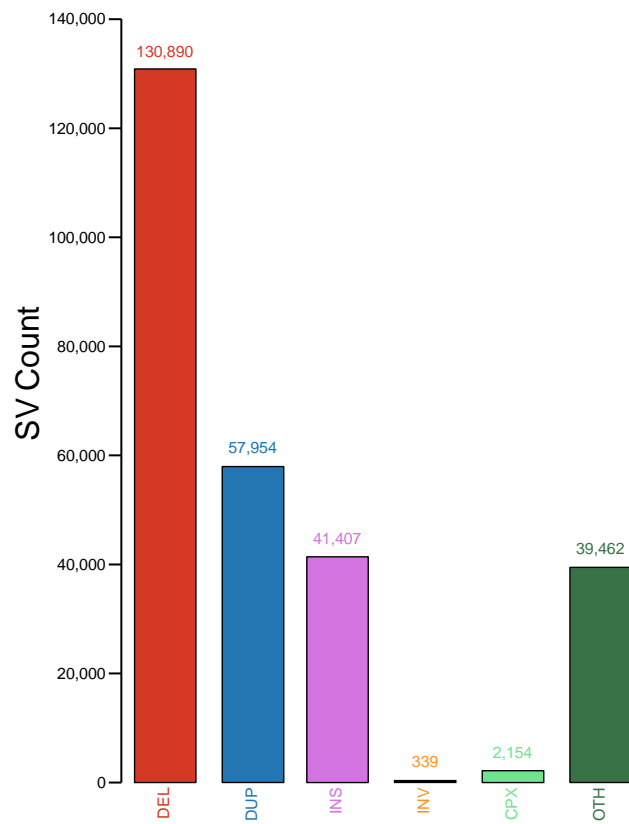- **Has been adapted to work on Google Cloud via Broad Institute's Terra Platform**
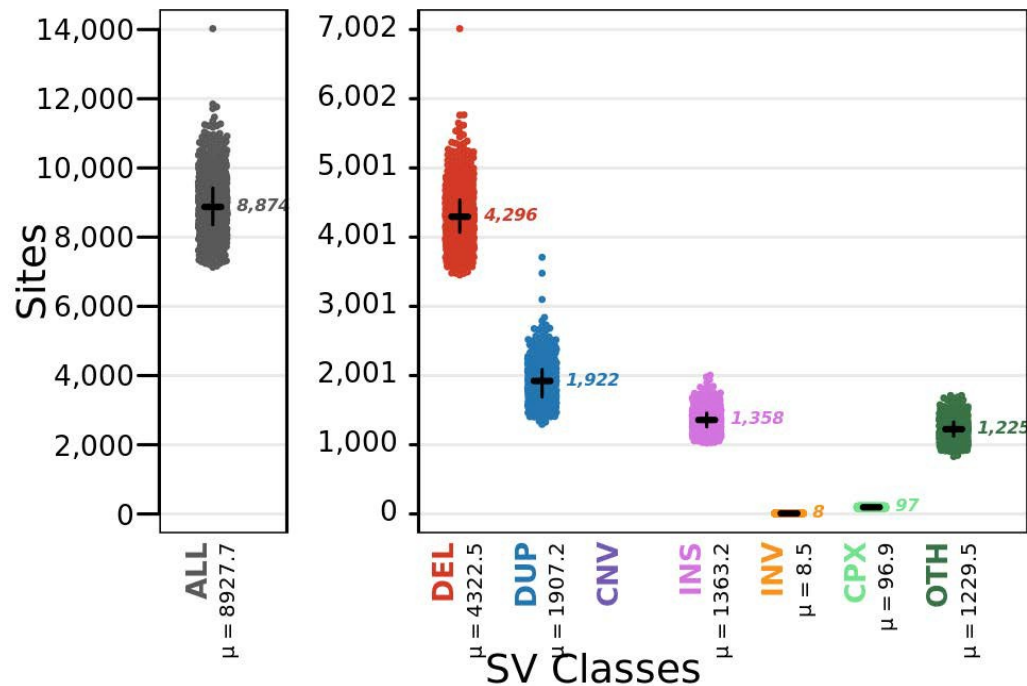
# Sex Chromosome Aneuploidies



- 1 sample of Klinefelter Syndrome (XXY)
  - ➢ 1 Proband

- 4 samples with Jacob Syndrome (XYY)
  - ➢ 3 Fathers
  - ➢ 1 Proband

# SV Counts from OFC Cohort

Total Variant Count

SV Per Sample

# DE NOVO SV DETECTION

- Similar to SNVs and indels, **careful filtering is required to identify** *de novo* SVs at high precision

- SVs can be **misclassified** due to **parental mosaicism**, **lack of phasing**, and/or **inconsistent evidence**

- We developed a *post hoc* GATK-SV *de novo* workflow to eliminate erroneous *de novo* events

- Sample phenotypes and ascertainment can have a huge effect on # of *de novo* SVs per cohort

# GATK-SV *De Novo* Workflow

## *Variant Level*

**Unfiltered VCF**
~4,000 *de novo* SV per person

**Remove mCNVs and BNDs**

**Exclude SV with >0.5% VF in parents**

**Variant Filtered**
~*100 de novo* SV per person

## Stringent Filtering

- Depth support in overlapping parental SV
- Require raw algorithm support
- PE/SR support on both sides
- Variable copy number or repetitive site
- Sample outliers excluded
- ROC GQ filters

**High Quality *de novo* SV**
8.1% of samples with ≥1 *de novo* SV

## Less Stringent Filtering

- Include filtered *de novo* SV from outliers

- Outlier *de novo SV* for future review

1. Reduce GQ filters for Private SV
2. CNV >500 KB
3. Recurrent Genomic Disorder SV
4. Repetitive Regions (i.e segdups)

Add mosaic variants and aneuploidies

**Manual Check (n = 1,105)**

- **TP: HQ rare inherited events with raw support**

- **FP: De novo variants no raw support or PE/SR support on both sides**

- **Both parent and child cutoffs investigated**

- **Depth & PE/SR assessed separately**

- **Variant level filtering for depth events derived from similar training set**

## *Sample Level*

***Final de novo SV (n = 165)***
***17.8% of samples with ≥1 de novo SV***

# GATK-SV *De Novo* Workflow

## *Variant Level*

**Unfiltered VCF**
~4,000 *de novo* SV per person

Remove mCNVs and BNDs

Exclude SV with >0.5% VF in parents

**Variant Filtered**
~*100 de novo* SV per person

## Stringent Filtering

Depth support in overlapping parental SV

Require raw algorithm support

PE/SR support on both sides

Variable copy number or repetitive site

Sample outliers excluded

ROC GQ filters

**High Quality *de novo* SV**
8.1% of samples with ≥1 *de novo* SV

## Less Stringent Filtering

Include filtered *de novo* SV from outliers

1. Reduce GQ filters for Private SV
2. CNV >500 KB
3. Recurrent Genomic Disorder SV
4. Repetitive Regions (i.e segdups)

Add mosaic variants and aneuploidies

**Outlier *de novo SV*** for future review

Manual Check (n = 1,105)

- **TP: HQ rare inherited events with raw support**

- **FP: De novo variants no raw support or PE/SR support on both sides**

- **Both parent and child cutoffs investigated**

- **Depth & PE/SR assessed separately**

- •**Variant level filtering for depth events derived from similar training set**

## *Sample Level*

**Final de novo SV (n = 165)**
**17.8% of samples with ≥1 de novo SV**

# GATK-SV *De Novo* Workflow

**Variant Level**

Unfiltered VCF
~4,000 *de novo* SV per person

Remove mCNVs and BNDs

Exclude SV with >0.5% VF in parents

Variant Filtered
~*100 de novo* SV per person

**Stringent Filtering**

Depth support in overlapping parental SV

Require raw algorithm support

PE/SR support on both sides

Variable copy number or repetitive site

Sample outliers excluded

ROC GQ filters

High Quality *de novo* SV
8.1% of samples with ≥1 *de novo* SV

**Less Stringent Filtering**

Include filtered *de novo* SV from outliers

1. Reduce GQ filters for Private SV
2. CNV >500 KB
3. Recurrent Genomic Disorder SV
4. Repetitive Regions (i.e segdups)

Add mosaic variants and aneuploidies

Outlier *de novo SV* for future review
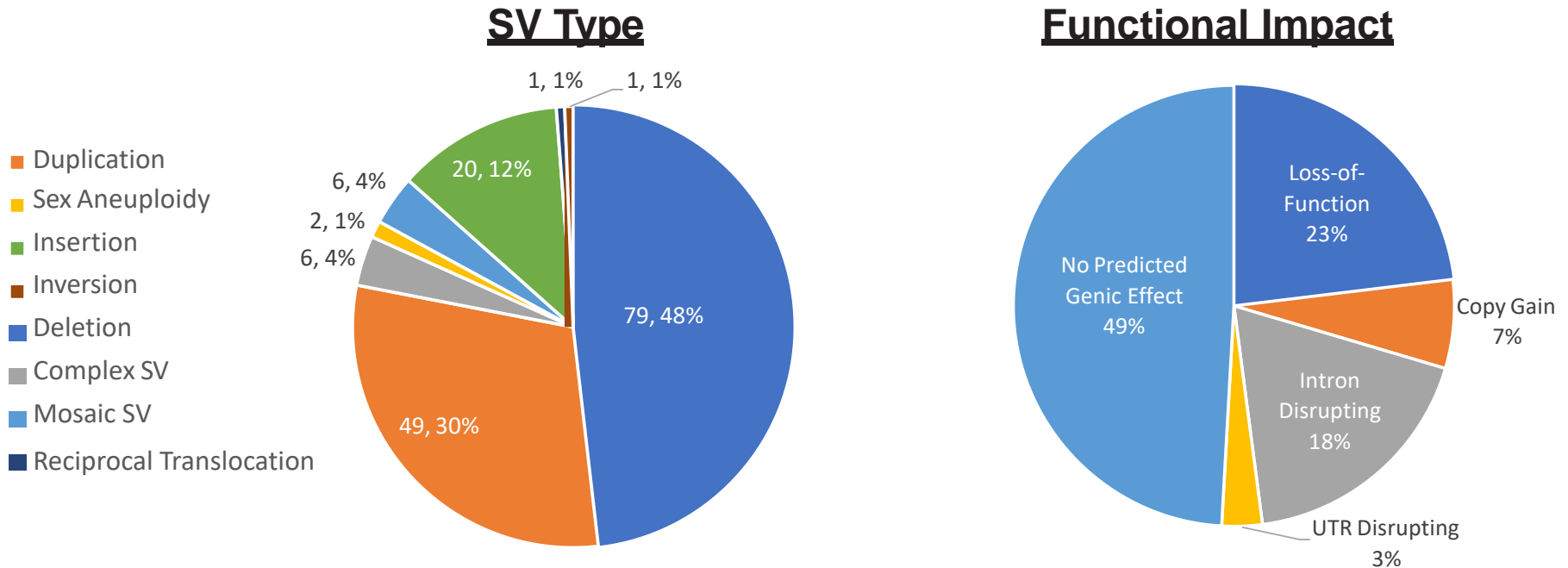
Manual Check (n = 1,105)

**Sample Level**

Final de novo SV (n = 165)
17.8% of samples with ≥1 de novo SV

- **TP: HQ rare inherited events with raw support**
- **FP: De novo variants no raw support or PE/SR support on both sides**
- **Both parent and child cutoffs investigated**
- **Depth & PE/SR assessed separately**
- **Variant level filtering for depth events derived from similar training set**

# GATK-SV *De Novo* Workflow

**Variant Level**

Unfiltered VCF
~4,000 *de novo* SV per person

Remove mCNVs and BNDs

Exclude SV with >0.5% VF in parents

Variant Filtered
~*100 de novo* SV per person

**Stringent Filtering**

- Depth support in overlapping parental SV
- Require raw algorithm support
- PE/SR support on both sides
- Variable copy number or repetitive site
- Sample outliers excluded
- ROC GQ filters

High Quality *de novo* SV
8.1% of samples with ≥1 *de novo* SV

**Less Stringent Filtering**

Include filtered *de novo* SV from outliers

1. Reduce GQ filters for Private SV
2. CNV >500 KB
3. Recurrent Genomic Disorder SV
4. Repetitive Regions (i.e segdups)

Add mosaic variants and aneuploidies

Outlier *de novo SV* for future review

Manual Check (n = 1,105)

**Sample Level**

**Final de novo SV (n = 165)
17.8% of samples with ≥1 de novo SV**

- **TP: HQ rare inherited events with raw support**
- **FP: De novo variants no raw support or PE/SR support on both sides**
- **Both parent and child cutoffs investigated**
- **Depth & PE/SR assessed separately**
- **Variant level filtering for depth events derived from similar training set**

# Summary of *De Novo* SVs

- After filtering, we observe 165 *de novo* SVs in 138 probands

- 17.8% of probands have at least one de novo event after accounting for 63 outlier samples that failed the de novo pipeline



**SV Type**

Duplication
Sex Aneuploidy
Insertion
Inversion
Deletion
Complex SV
Mosaic SV
Reciprocal Translocation

1, 1%   1, 1%
20, 12%
6, 4%
2, 1%
6, 4%
79, 48%
49, 30%

**Functional Impact**

Loss-of-Function
23%
No Predicted Genic Effect
49%
Copy Gain
7%
Intron Disrupting
18%
UTR Disrupting
3%

# HIGHLIGHTED DE NOVO SV

## De Novo SV in Recurrent Genomic Disorder Regions (n = 14 in cohort)

| Syndrome | Size | N | OFC Reported in Syndrome | Ethnicity | Case Phenotype |
|---|---|---|---|---|---|
| 1q21.1 proximal dup | 195 kb | 1 | Yes | Asian | CL/P |
| 7q11.23 dup | 1.4 Mb | 2 | Yes | Asian African | CL/P CL/P |
| 8p23.1 dup | 3.6 Mb | 1 | Yes | Latino | CL/P |
| 15q11.2 del (BP1-BP2) | 290 kb | 1 | Yes | Latino | CL/P |
| 15q11.2 dup (BP1-BP2) | 290 kb | 1 | No | Caucasian | CL |
| 16p11.2 distal del | 224 kb | 1 | Yes | Asian | CL/P |
| 16p11.2 distal dup | 224 kb | 1 | No | Caucasian | CP |
| 22q11.2 del | 2.6 Mb | 1 | Yes | Caucasian | CL/P |
| 22q11.2 dup | 2.6 Mb | 1 | Yes | Caucasian | CL/P |
| 22q11.2 distal deletion | 1.7 Mb | 2 | Yes | Caucasian African | CL/P CL/P |
| 22q11.2 distal del | 1.7 Mb | 1 | Yes | Asian | CL/P |
| Xp22.31 del (female) | 1.6 Mb | 1 | No | Caucasian | CL/P |

5.3 mb deletion
(chr4:49,952,415-55,275,096)

- Disrupts >30 protein coding genes
- Difficult to pinpoint single causative gene

78.6 kb deletion disrupts *BTBD18,*
***CTNND1**, SELENOH, TMX2*

3.2 kb deletion disrupts ***TFAP2A***

Mutations in the Epithelial Cadherin-p120-Catenin
Complex Cause Mendelian Non-Syndromic
Cleft Lip with or without Cleft Palate

Liza L. Cox,[1,2,3] Timothy C. Cox,[1,2,4,31,*] Lina M. Moreno Uribe,[5] Ying Zhu,[6,7] Chika T. Richter,[5] Nichole Nidey,[6] Jennifer M. Standley,[6] Mei Deng,[9] Elizabeth Blue,[10] Jessica X. Chong,[11] Yueqin Yang,[12] Russ P. Carstens,[12,13] Deepti Anand,[14] Salil A. Lachke,[14] Joshua D. Smith,[15] Michael O. Dorschner,[16,17] Bruce Bedell,[8] Edwin Kirk,[6,18] Anne V. Hing,[1,19] Hanka Venselaar,[20] Luz C. Valencia-Ramirez,[21] Michael J. Bamshad,[11,15] Ian A. Glass,[9,11] Jonathan A. Cooper,[3] Eric Haan,[22,23] Deborah A. Nickerson,[15] Hans van Bokhoven,[24,25] Huiqing Zhou,[24,26] Katy N. Krahn,[27] Michael F. Buckley,[6] Jeffrey C. Murray,[6] Andrew C. Lidral,[28] and Tony Roscioli[18,29,30,31,*]

**REPORT**

*TFAP2A* Mutations Result
in Branchio-Oculo-Facial Syndrome

Jeff M. Milunsky,[1,2,3,*] Tom A. Maher,[1] Geping Zhao,[1] Amy E. Roberts,[4] Heather J. Stalker,[5] Roberto T. Zori,[5] Michelle N. Burch,[5] Michele Clemens,[6] John B. Mulliken,[7] Rosemarie Smith,[8] and Angela E. Lin[9]

424.1 kb deletion disrupts ***PDGFC***

A specific requirement for PDGF-C in palate formation and PDGFR-α signaling

Hao Ding[1], Xiaoli Wu[1], Hans Boström[2], Injune Kim[3], Nicole Wong[4], Bonny Tsoi[4], Meredith O'Rourke[4], Gou Young Koh[3], Philippe Soriano[5], Christer Betsholtz[2], Thomas C Hart[6], Mary L Marazita[7], L L Field[8], Patrick P L Tam[4] & Andras Nagy[1,9]

9.0 kb deletion disrupts ***GRHL2***

***Grainyhead-like 2*** regulates neural tube closure and adhesion molecule expression during neural fold fusion

Christina Pyrgaki[1], Aimin Liu[2], and Lee Niswander[1,*]

[1] HHMI, Department of Pediatrics, Molecular Biology Graduate Program, University of Colorado School of Medicine, Aurora, CO 80045 USA

[2] Department of Biology, Eberly College of Science, The Pennsylvania State University, University Park, PA 16802 USA

34 kb duplication disrupts *RFC1*

Intragenic Exonic Duplication



Reduced folate carrier 1 (*RFC1*) is associated with cleft of the lip only

A.R. Vieira[1,2,3,4], M.E. Cooper[1,3], M.L. Marazita[1,3,4,5], E.E. Castilla[6,7] and I.M. Orioli[8]

RFC1 and non-syndromic cleft lip with or without cleft palate: An association based study in Italy

Ambra Girardi[a], Marcella Martinelli[a,*], Francesca Cura[a], Annalisa Palmieri[a], Francesco Carinci[b], Enrico Sesenna[c], Luca Scapoli[a]

## Constrained gene defined as gnomAD LOEUF < 0.4



8.1 kb deletion disrupts *NRF1*



https://gnomad.broadinstitute.org

- Mutations in *NRF1* not yet associated with OFC

- Transcription factor involved in several pathways that could contribute to OFC

De Novo Complex Events (n = 6 in cohort)

Complex Event on Chromosome 4

# De Novo Reciprocal Translocations (n = 1 in cohort )

46,XY,t(1;2)(q42.12;q12.3) -  Disrupts *CNIH3*

derivative chr 1

chr1:224448379    chr2:107350640

derivative chr 2

chr2:107350644    chr1:224448382

Breakpoint overlaps edge of **1q41-q42 deletion** syndrome

### PERINATAL/NEONATAL CASE PRESENTATION

A neonate with the Pelger–Huët anomaly, cleft lip and palate, and agenesis of the corpus callosum, with a chromosomal microdeletion involving 1q41 to 1q42.12

RD Christensen[1] and HM Yaish[2]

Ideograms modified from https://en.wikipedia.org/

103 kb mosaic deletion predicted to be **present in 40%** of white blood cells

582 kb mosaic duplication predicted to be **present in 70%** of white blood cells

# Conclusions

- Application of GATK-SV was able to discover a diverse set of SV in the OFC samples

- Adjudication with additional *de novo* filtering identified 165 *de novo* SV in 17.8% of probands

- We find both established OFC genes disrupted and novel candidate genes for further follow-up

- WGS has the resolution to detect complex SV and balanced SV not easily detectable by exome sequencing or microarrays

# Future Directions

- Exploration of rare inherited SVs

- Examination of noncoding SV

- Integration of results with the SNV/Indel callset presented by Elizabeth Leslie

- Investigation of recessive and compound heterozygous variation

- Applying GATK-SV in additional GMKF cohorts to build an aggregated SV map of congenital birth defects

# Acknowledgments

**Broad GMKF & Broad-SV
Sequencing and Analysis Teams**

**GMKF OFC Working Group**

**Harrison Brand**
**Michael Talkowski**
**Stacey Gabriel**
**Daniel MacArthur**
Xuefang Zhao*
Stacey Mano
Ben Weisburd
Ryan Collins
Harold Wang
Mark Walker
Chris Wheelan
Candace Patterson

University of Iowa
**Azeez Butali**
Jeff Murray
Lina Moreno
Luz Consuelo
  Valencia-Ramirez
George Wehby
Andrew Lidral

Emory University
**Elizabeth J. Leslie**
Madison Bishop
Pankaj Chopra
Michael Mortillo
Dave Cutler
Michael Epstein

University of Pittsburgh
**Mary L. Marazita**
Seth M. Weinberg
Eleanor Feingold
Nandita Mukhopadhyay

Johns Hopkins University
**Terri Beaty**
Ingo Ruczinski
Margaret Taub
Alan Scott
Jacqueline Hetmanski
Debashree Ray

Other
Jacqueline Hecht
Andrew Czeizel
Yah-Huei Wu Chou
Frederic Deleyiannis
Adebowale Adeyemo
Mauricio Arcos-Burgos

Lord Gowans
Peter Mossey
Lanre Adeyemo
Philip Chen

# Kids First Data:

# Kids First Data - By the Numbers . . .

**Studies**
10

# Kids First Data - By the Numbers . . .

**Studies**
10

**Participants**
10,560

# Kids First Data - By the Numbers . . .



**Studies**
10

**Participants**
10,560

**Families**
3,679

# Kids First Data - By the Numbers . . .

**Studies**
10

**Participants**
10,560

**Families**
3,679

**Diagnoses**
7,502

# Kids First Data - By the Numbers . . .

**Studies**
10

**Participants**
10,560

**Families**
3,679

**Diagnoses**
7,502

**Phenotypes**
70,916

**Genomes**
10,901

**Size**
~1 PB

# Outline

Two new major Kids First DRC portal feature developments

- Enhanced ontology data model and search tool

- The germline variant data warehouse

Next steps/future directions

# Ontologies within Kids First DRC

**Kids First DRC makes extensive use of ontologies**

Human Phenotype Ontology (HPO), Mondo, NCIT, SNOMED

**Ontologies provide both controlled vocabularies and "parent-child" relationships**

**E.g.** *Oral cleft* **(HP:0000202) IS AN** *Abnormal oral cavity morphology* (HP:0000163)

**New**  **The portal now integrates relationships in participant search queries**

Users can now find participants with a specific term *and* all its descendant

E.g. Searching for participants with *Abnormal oral cavity morphology* **will return** *Oral cleft* **participants**

*https://portal.kidsfirstdrc.org/explore*

# Variant Data within Kids First DRC

Currently available in gVCF files:

- Files can be searched using the portal's File Repository
- Selected files can be pushed to Cavatica for in-depth analyses

# The KFDRC Variant Data Warehouse Workspace Environment

New

A performant and scalable variant database that can be queried directly from the portal

### Comprehensive set of variant annotations

Genes, allele frequencies, gene panels, inheritance, functional impact predictions, pathways, external references, etc.

### Individual-level clinical data integration to enable multi-dimensional queries

E.g. find all **rare missense** variants with **high functional impact** in **low grade glioma** patients affected by any **cardiovascular abnormalities**

### Web-based variant data analytics and visualisation tools

### Security and privacy rules enforcement

Users can only access variant datasets they have been authorized to

# A *Big Data* Challenge

- **High number of germline variants to process from whole genomes**

- **Current version**
  - 8 studies, 8,100 participants, **251,801,242** unique variants, **42,513,213,093** occurrences

- **For comparison/context**
  - NCI Genomic Data Commons (GDC): 3.1 M somatic variants for ~10,000 cases
  - International Cancer Genome Consortium (ICGC): 82 M somatic variants for 19,700 cases

- **Challenge: Complex data to query through responsive web interfaces**
  - Link to extended individual-level clinical data
  - Integrate rich variant annotations

# KFDRC Variant Data Processing Workflow

gVCFs on AWS

Participant data

Variant annotation files E.g TopMed

## Spark Cluster

Parquet Variant DB

Zeppelin Notebook

## ElasticSearch Cluster

Data indices

**Spark, Parquet & ElasticSearch** Technologies that can scale with data growth.

# Phase I (First Release, Beta)
## Foundation & Zeppelin notebooks

**Objectives**

- Build and deploy the foundational infrastructure of the KFDRC variant warehouse database
- Implement the data extraction, annotation and loading workflow
- Annotate variants with a limited (initially) set of annotations
- Implement the data security framework
- Provide researchers with the Zeppelin data analytic environment for querying and analysing the variant database
- Link the variant data analytic environment to the Cohort Builder, enabling researchers to analyse variants from their virtual patient cohorts

Spark Cluster

Parquet Variant DB

Zeppelin Notebook

# The Zeppelin Data Analytic Environment

**Provides programmatic access to the variant database from web browsers**

- Accessible from the Portal

- User notebook workspace

- Private/Individual Spark clusters on AWS

- Support for various programmatic languages (SQL, Python, R, Scala)

*https://portal.kidsfirstdrc.org/variantDb*

# Demo

# Next Steps

- *Performance tests and data quality control (QC)*
- *Beta release to KFDRC user groups and X01 Investigators*
- *User testing and feedback integrations*

**Additional Short Term Development Road map**

- **Indexing variant data warehouse using Elasticsearch**
- **Build data querying interfaces within the portal (integrating notebook use cases)**
  - **GA4GH-like Beacon service (return yes/no answers on variant occurrences)**
  - **Gnomad-like Summary interface (mainly allele frequencies)**
  - **Direct integration within the Cohort Builder allowing complex queries that return *both* participant and variant lists**
- **More annotations supporting variant prioritization**

# Special Thanks To

**CHU Ste-Justine Research Center**

- **Jeremy Costanza**, Lead software architect and developer
- Developers
    - Adrian Paul
    - Evans Girard
    - Francis Lavoie
- UX
    - Lucas Lemonnier

**CHOP**

- DevOps Lead
    - Alex Lubneuski
- Bioinformatics Lead
    - Yuankun Zhu

# CAVATICA: Cloud User Workspace Introduction

# CAVATICA: Integrated Cloud-Based Workspace

# Data Inputs - Kids First

**Immediately bring over files you're authorized to use into CAVATICA**

# Data Inputs - Kids First

# Data Inputs - Kids First

**New projects created in CAVATICA start as private**

# Data Inputs - Kids First

CAVATICA continues to check permissions for integrated datasets

# Data Inputs - Own Data on AWS or GCP



**Via the Data->Volumes Features**

# Data Inputs - Own Data on AWS or GCP



**After adding own files to same private project with Kids First files, can now utilize existing or bring your own workflows**

# Workflows - Existing Examples

# Workflows - Port Your Own



Tool docker image

Tool details, basic command line, input/output arguments, parameters etc

Final command line preview

# Data Cruncher - Interactive Analysis

# Data Cruncher - Interactive Analysis

# Data Cruncher - RStudio and Shiny Apps

# Data Cruncher - RStudio and Shiny Apps



## Beta Feature!

# All Features are Collaborative

You maintain control of access to your projects and data.

# Germline mutations in cancer susceptibility genes occur in 8-10% of pediatric cancers
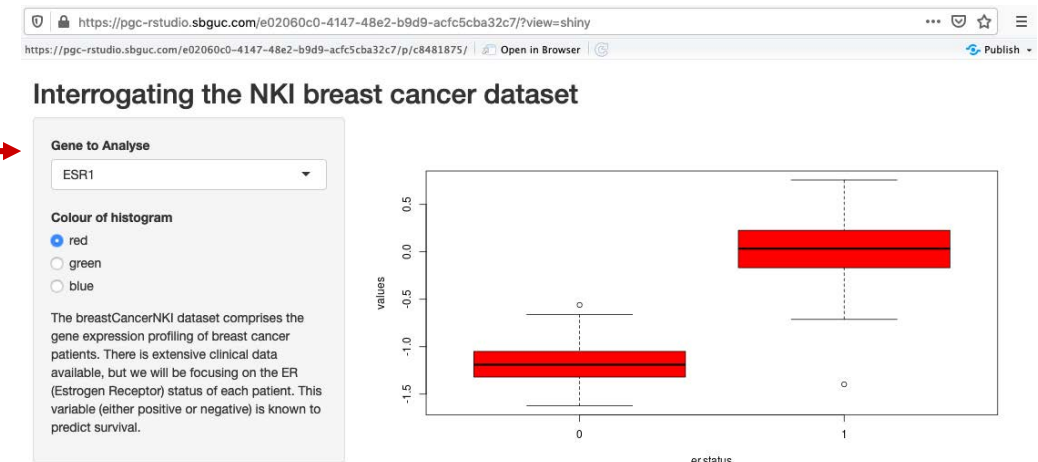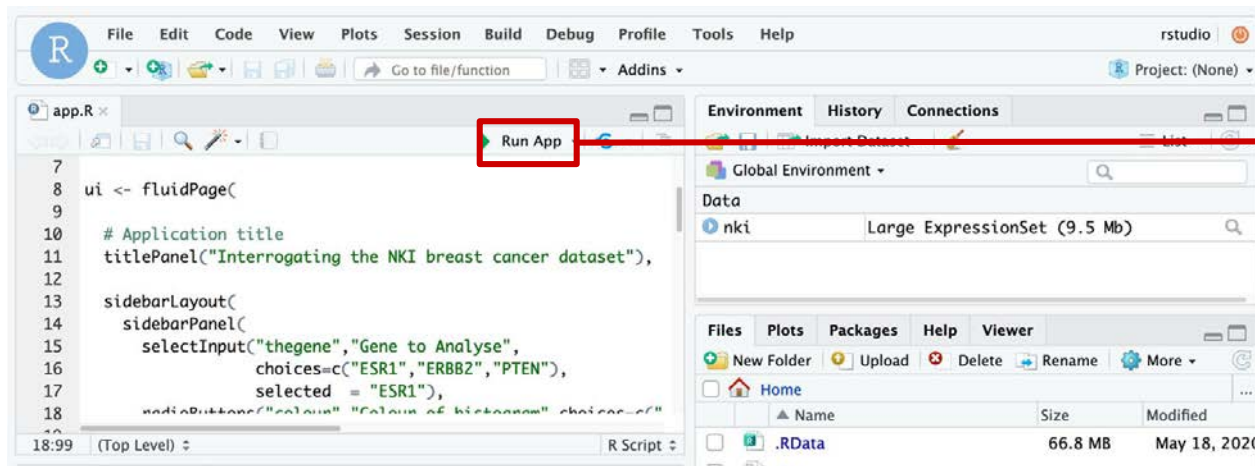
Germline susceptibility genes include: RB1, NF1, WT1 etc.

**Susceptibility genes and Wilms Tumor**

1991

2019

WT1 first
cloned

multiple other susceptibility genes
discovered

4 new genes discovered from
WES

**Not all cancer predisposition or susceptibility genes have been identified**

# Probands from BASIC3 have undergone germline and somatic WES

**Goal:** characterize the diagnostic yield of combined tumor and germline WES for 287 children with solid tumors

- Not enriched for specific cancer type between CNS and non-CNS tumors

- Found pathogenic variants in
    - Genes with associated with specific cancers
    - Genes not previously associated with specific cancers

# 120 probands-parents from BASIC3 selected for germline WGS

**Goal 1:** Identify *de novo* Single Nucleotide Variants (SNVs) and Structural Variants (SVs)

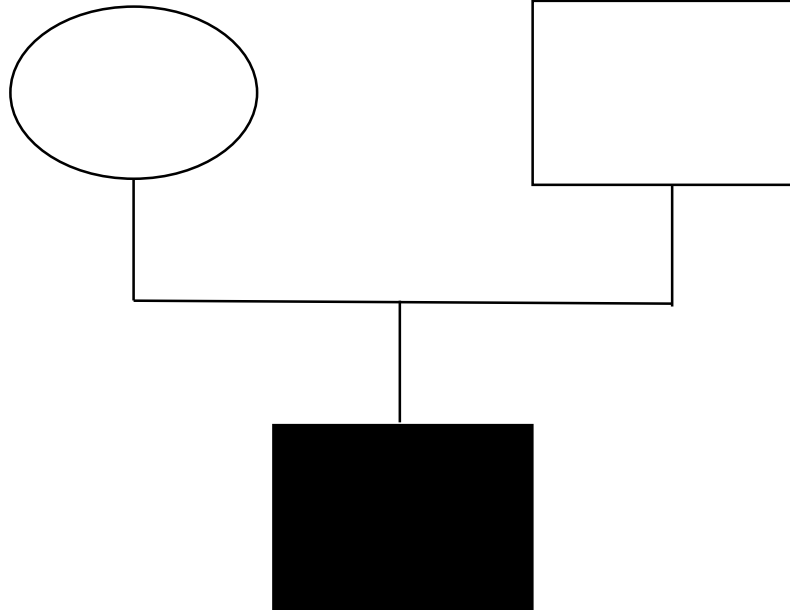**Goal 2:** Identify putative pathogenic variants in known cancer genes that may have been missed by WES



Breakdown of trios, duos, and singletons in cohort

- Trios: N=63
- Duos: N=52
- Singletons: N=5

# Use of Platypus for *de novo* SNV calling on Cavatica

## Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications

Andy Rimmer[1,na1], Hang Phan[1,na1], Iain Mathieson[1], Zamin Iqbal[1], Stephen R F Twigg[2], WGS500 Consortium, Andrew O M Wilkie[2], Gil McVean[1,3,na1] & Gerton Lunter ✉[1]

## 🅣 Platypus

Created by vojislav_varjacic on Mar. 12, 2018 06:46 • Last edited by vojislav_varjacic on Aug. 15, 2018 06:41
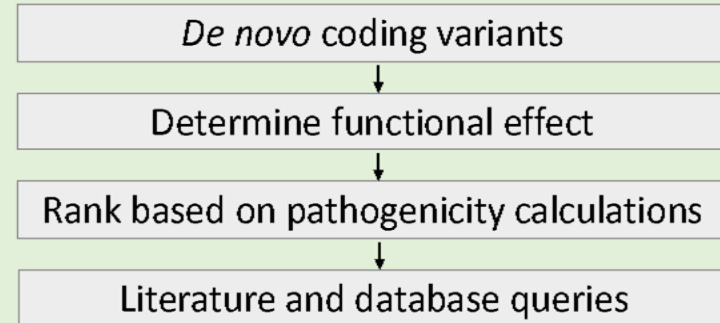Revision note: "typo in JS fixed"

### Description

**Platypus** is a tool designed for efficient and accurate variant-detection in high-throughput sequencing data.

**Platypus** reads data from **BAM files**, and outputs a **single VCF file** containing a list of identified variants, and genotype calls and likelihoods for all samples.

# Analysis on Cavatica expedited *de novo* variant discovery



**Prioritization of *de novo* variants:**

- *De novo* coding variants
- Determine functional effect
- Rank based on pathogenicity calculations
- Literature and database queries

**Outcome:**

- SNV analysis completed on 54 proband-parent trios
- The pipeline resulted in an expected number of variants per trio

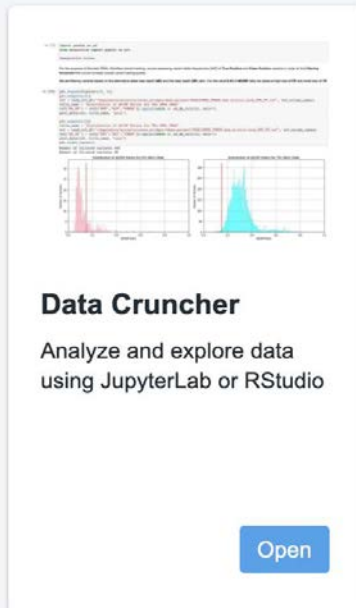| Variant Type | Frequency |
|---|---|
| Genome-wide *de novo* | 60 to 190 |
| Coding *de novo* | 0 to 4 |

# *De novo* SV analysis on Cavatica



**Caller A, B, C, D, & E:**
**Lumpy, Manta, Delly,**
**Breakdancer, & CNVnator**

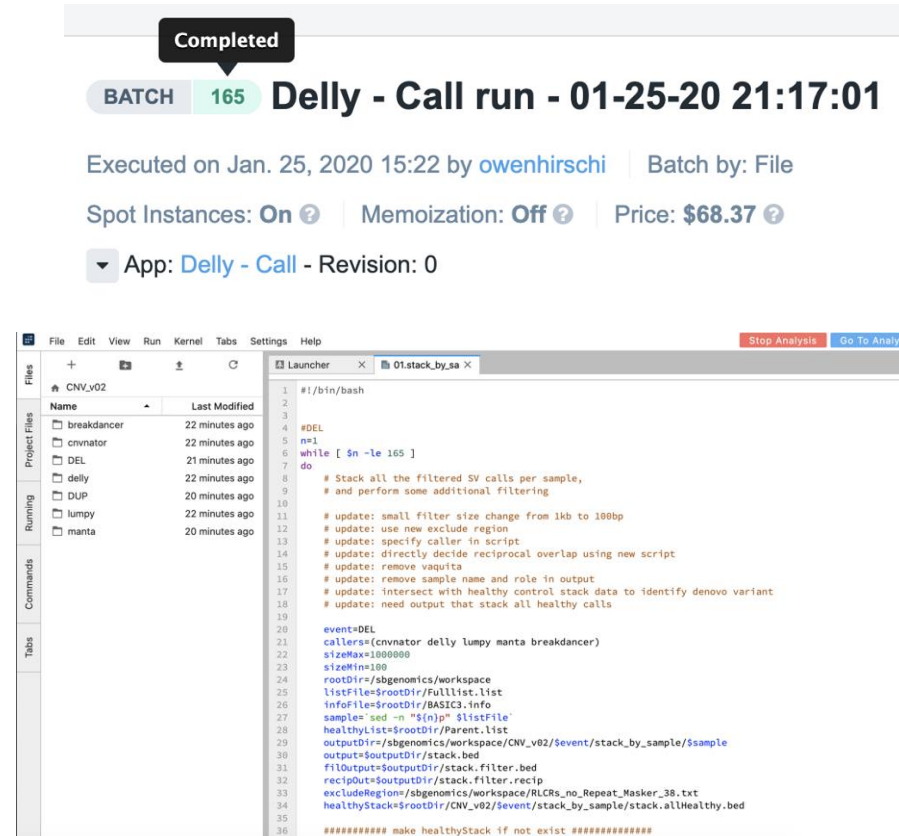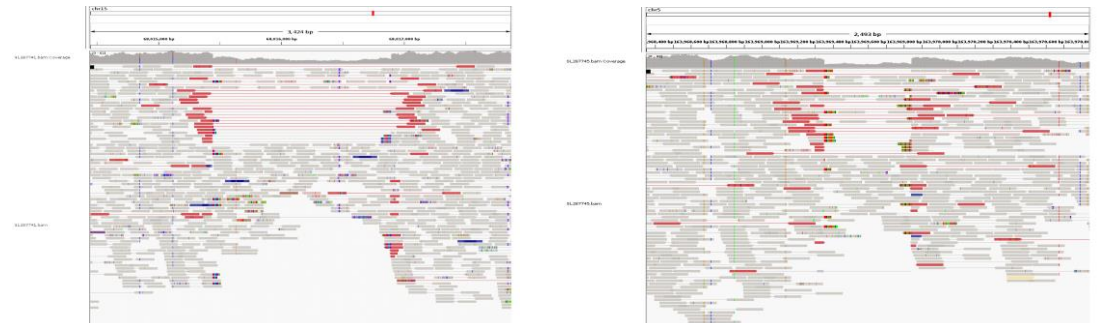# Analysis of SV on Cavatica requires multiple features of the platform



IGV images

# Analysis of miRNA variation on Cavatica

## Framework for microRNA variant annotation and prioritization using human population and disease datasets

Ninad Oak, Rajarshi Ghosh, Kuan-lin Huang, David A. Wheeler, Li Ding, Sharon E. Plon ✉

T **ADmiRE** PLONLABBCM

Created by owenhirschi on Feb. 7, 2020 11:02 • Last edited by owenhirschi on Feb. 7, 2020 13:57

**Description**

Annotative Database of miRNA Elements, ADmiRE, combines multiple existing and new biological annotations to aid the prioritization of causal miRNA variation.

ADmiRE Highlights: Annotation wrapper for adding comprehensive miRNA annotations to a user-supplied list of variants (tab-separated format) Adds information for miRNA domains, gnomAD mean allele frequency percentiles, evolutionary conservation, etc.

perl annotate_admire.pl [--input INPUT_FILE] [--output OUTPUT_FILE] [--admire_path PATH] [--chr NUMBER] [--pos NUMBER]

--input: INPUT_FILE [REQUIRED]

--output: OUTPUT_FILE (Default: INPUT_FILE.ADmiRE.tab) [OPTIONAL]

--admire_path: Path to ADmiRE.tab database. (Default: same directory with annotate_admire.pl) [OPTIONAL]

--chr: Column number in the INPUT_FILE with chromosome information. (Default: 1 -1st column) [OPTIONAL]

--pos: Column number in the INPUT_FILE with base position information. (Default: 2 -2nd column) [OPTIONAL]

# Acknowledgments

**Plon Lab members:**

Sharon Plon, MD, PhD

Saumya Sisoudiya

Adam Weinstein

Deborah Ritter, PhD

Xi Luo, PhD

Ryan Zabriskie

Ninad Oak, PhD- former

**Baylor HGSC:**

Hurley Li, PhD

**BASIC3 Co-PI:**

William Parson, MD, PhD

Texas Children's Cancer Center

HGSC
HUMAN GENOME SEQUENCING CENTER

Baylor College of Medicine

# Kids First DRC - The Model We Follow



**Japan win silver in the 4-x-100-meter relay at the Rio de Janeiro Games**
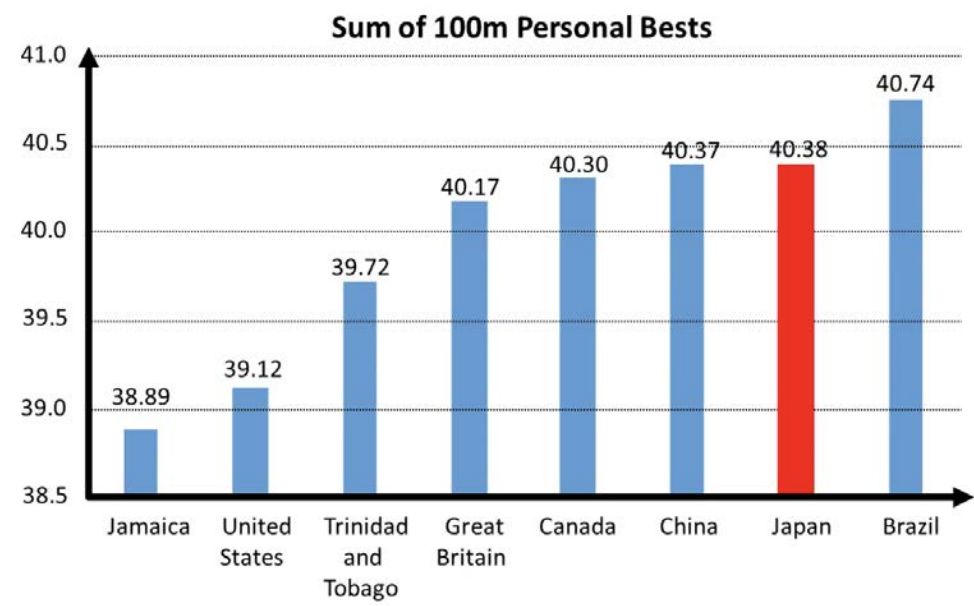


**None of their team having ever run 100m in under 10 seconds**

# Kids First DRC - The Model We Follow



**Japan win silver in the 4-x-100-meter relay at the Rio de Janeiro Games**



Sum of 100m Personal Bests

**None of their team having ever run 100m in under 10 seconds**

# Kids First DRC - The Model We Follow

# Addressing Scale - A Model

# Addressing Scale - A Model

**Danyelle Winchester, PhD**
Health Specialist
Office of Strategic Coordination
Division of Program Coordination, Planning, and Strategic Initiatives
Office of the Director, National Institutes of Health (NIH)

# 10 Released Datasets

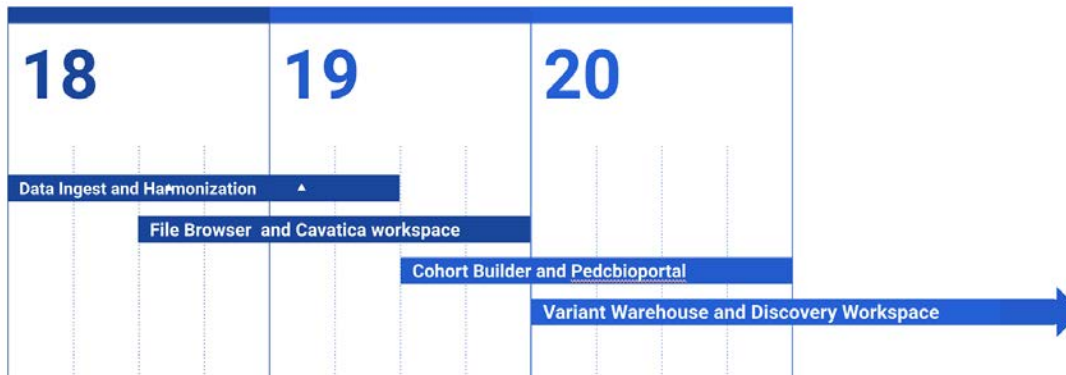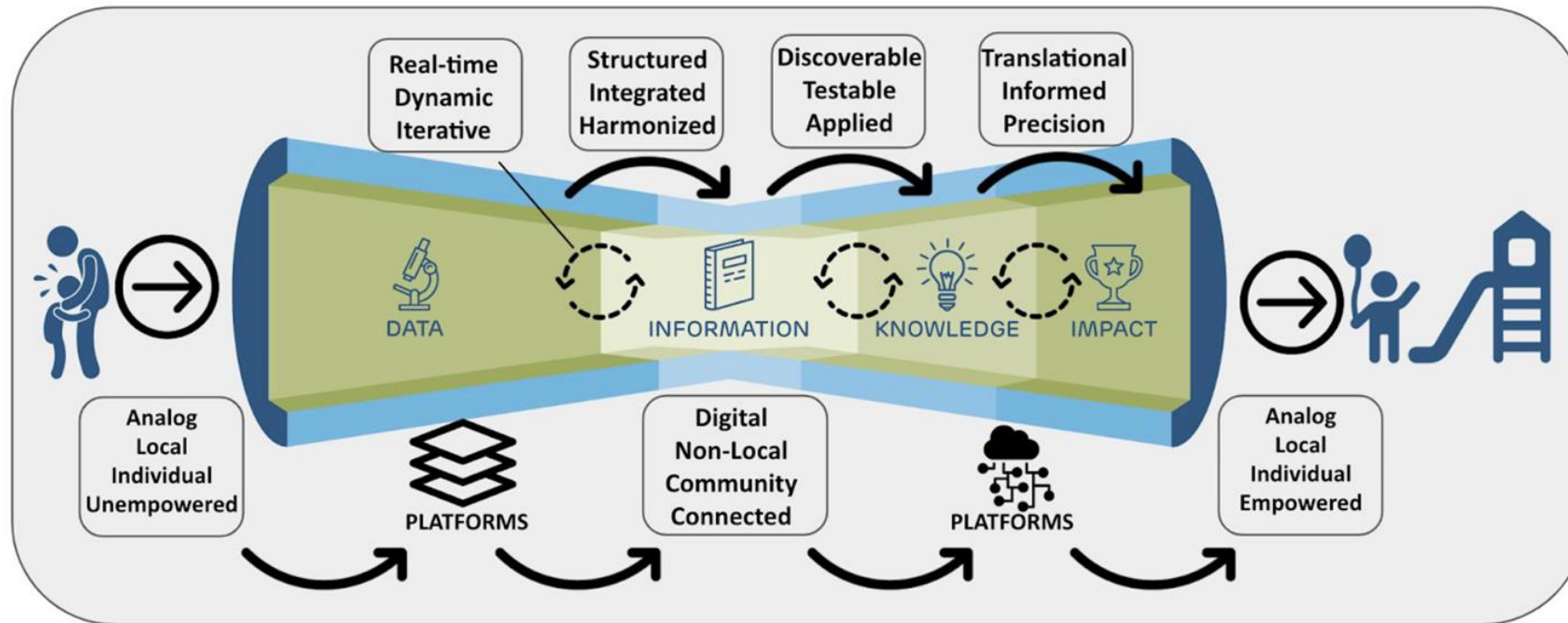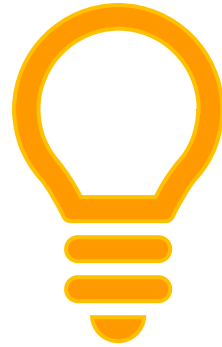- Disorders of Sex Development.           PI: Eric Vilain
- Congenital Diaphragmatic Hernia         PI: Wendy Chung
- Ewing Sarcoma                           PI: Joshua Schiffman
- Orofacial Clefts: Caucasian families    PI: Mary Marazita
- Orofacial Clefts: Latin American families   PI: Mary Marazita
- Structural Heart & Other Defects        PI: Christine Seidman (PCGC)
- Cranial Dysinnervation Disorders        PI: Elizabeth Engle
- Adolescent Idiopathic Scoliosis         PI: Jonathan Rios
- Neuroblastoma                           PI: John Maris
- Enchondromatoses                        PI: Nara Sobreira

- **Kids First DRC website:** https://kidsfirstdrc.org/support/studies-and-access/

- **NIH Kids First Umbrella BioProject:** https://www.ncbi.nlm.nih.gov/bioproject/338775 > dbGaP links

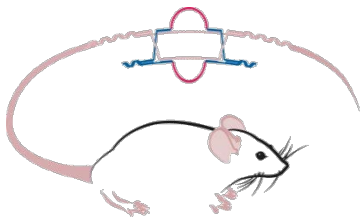- **X01 Abstracts**: https://commonfund.nih.gov/kidsfirst/x01projects

*The value of Kids First datasets will be amplified when researchers use and analyze these data to make discoveries that will ultimately improve prevention, diagnostics, and therapeutic interventions for these conditions*

# *Researchers are using Kids First data to answer new scientific questions*

➢ ***13 awards*** *for R03 for analyses of Kids First data* (PAR-16-348 ; PAR-18-733; PAR-19-069, PAR-19-375)

➢ ***1 award for NIDCR R03*** (PAR-16-070)

➢ ***2 awards for R01s*** (PA-13-302, PAR-17-236)

➢ Spurred **new collaborations** with KOMP2 & INCLUDE



*Knockout Mouse Phenotyping Project (KOMP2)*



*INvestigation of Co-occurring conditions across the Lifespan to Understand Down syndrome (INCLUDE)*

# Kids First Publications

**ELSEVIER**
Guide for Authors | About | Explore this Journal

**AJHG**

## Germline 16p11.2 Microdeletion Predisposes to Neuroblastoma

Laura E. Egolf,[1,2,3] Zalman Vaksman,[2,3,4] Gonzalo Lopez,[2,3,4] Jo Lynne Rokita,[2,3,4] Apexa Modi,[2,3,5] Patricia V. Basta,[6,7] Hakon Hakonarson,[8,9] Andrew F. Olshan,[6,7] and Sharon J. Diskin[1,2,3,4,5,10,*]

▸ Author information ▸ Article notes ▸ Copyright and License information Disclaimer

---

**Human Mutation**
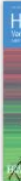Variation, Informatics, and Disease

**HGVS**
OFFICIAL JOURNAL

RESEARCH ARTICLE 🔓 Full Access

## Deleterious de novo variants of X-linked *ZC4H2* in females cause a variable phenotype with neurogenic arthrogryposis multiplex congenita

Suzanna G.M. Frints ✉, Friederike Hennig, Roberto Colombo, Sebastien Jacquemont, Paulien Terhal, Holly H. Zimmerman, David Hunt, Bryce A. Mendelsohn, Ulrike Kordaß ... See all authors ∨

---

**Total number of publications
from Kids First ~13
Average RCR: 1.32**

---

**Springer** Link

## Whole genome sequencing of orofacial cleft trios from the Gabriella Miller Kids First Pediatric Research Consortium identifies a new locus on chromosome 21

Nandita Mukhopadhyay, Madison Bishop, Michael Mortillo, Pankaj Chopra, Jacqueline B. Hetmanski, Margaret A. Taub, Lina M. Moreno, Luz Consuelo Valencia-Ramirez, Claudia Restrepo, George L. Wehby, Jacqueline T. Hecht, Frederic Deleyiannis, Azeez Butali, Seth M. Weinberg, Terri H. Beaty, Jeffrey C. Murray, Elizabeth J. Leslie, Eleanor Feingold & Mary L. Marazita ✉

---

**PLOS GENETICS**

🔓 OPEN ACCESS 📄 PEER-REVIEWED

RESEARCH ARTICLE

## *De novo* variants in congenital diaphragmatic hernia identify *MYRF* as a new syndrome and reveal genetic overlaps with other developmental disorders

Hongjian Qi co, Lan Yu co, Xueya Zhou co, Julia Wynn, Haoquan Zhao, Yicheng Guo, Na Zhu, Alexander Kitaygorodsky, Rebecca Hernan, Gudrun Aspelund, Foong-Yen Lim, Timothy Crombleholme, Robert Cusick, [ ⋯ ], Yufeng Shen ✉ [ view all ]

# Kids First Publications: How to Acknowledge Kids First Data

- Secondary users (end users) must acknowledge the dataset(s) they use by listing dbGaP accession numbers and the databases from which the data were accessed (e.g. link to the Kids First Data Resource Center or Portal). The acknowledgement statement can be found at the bottom of the dbGaP study page and in the Data Use Certification.

  - **See Frequently Asked Questions for X01 Cohorts Selected for Sequencing #5**: https://commonfund.nih.gov/kidsfirst/FAQ#X01%20selected

- **Principal Investigator**
  - Wendy Chung, MD, PhD. Columbia University Medical Center, New York, NY, USA
- **Co-Principal Investigator**
  - Yufeng Shen, PhD. Columbia University Medical Center, New York, NY, USA
- **Funding Sources**
  - X01 HL132366. National Institutes of Health, Bethesda, MD, USA
  - X01 HL136998. National Institutes of Health, Bethesda, MD, USA
  - X01 HL140543. National Institutes of Health, Bethesda, MD, USA
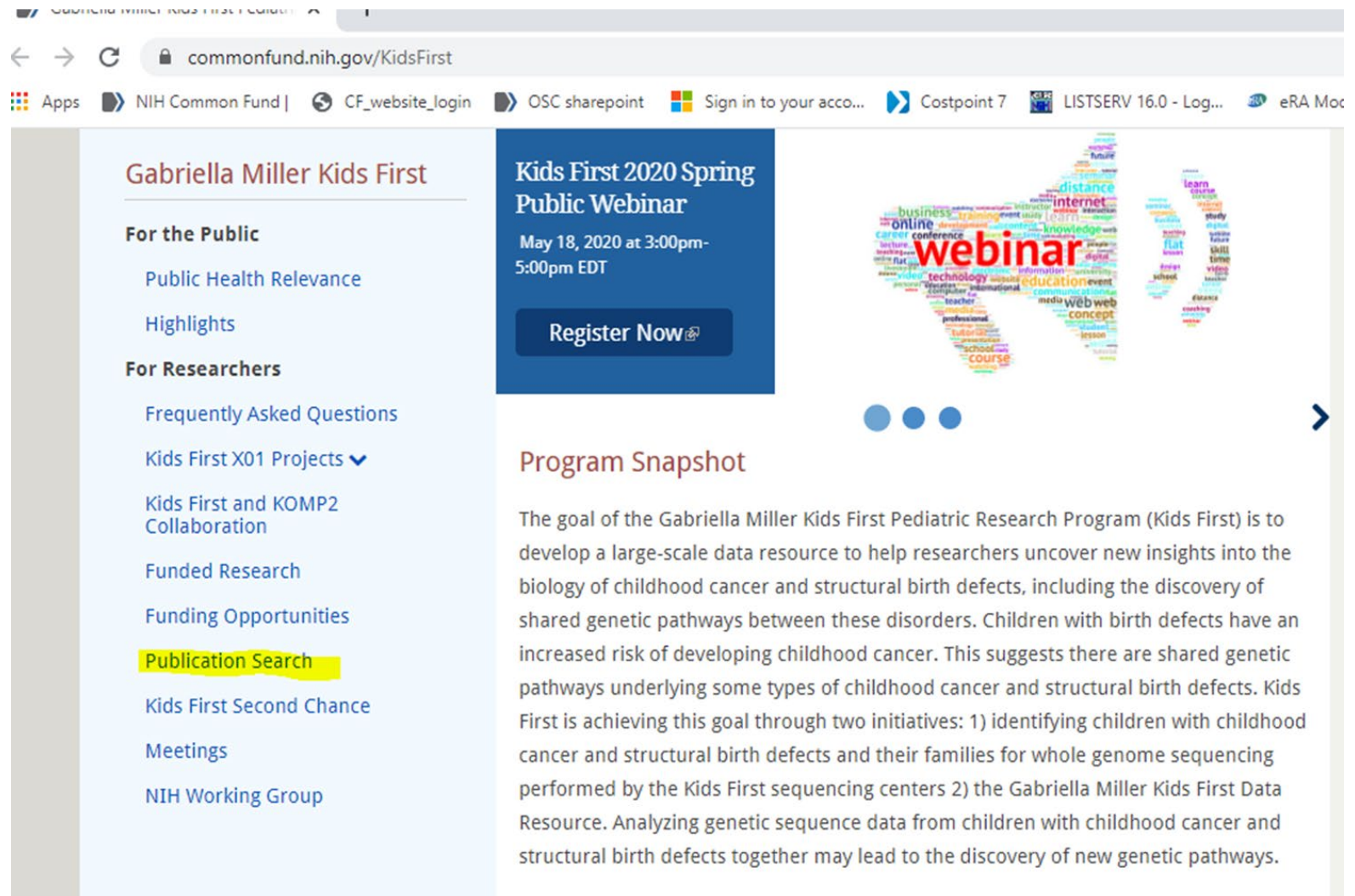  - R01 HD057036. National Institutes of Health, Bethesda, MD, USA

**Acknowledgement Statement:** Please cite/reference the use of dbGaP data by including the dbGaP accession phs001110.v2.p1. Additionally, use the following statement to acknowledge the submitter(s) of this study:

The results analyzed and <published or shown> here are based in whole or in part upon data generated by Gabriella Miller Kids First Pediatric Research Program projects <insert phs accession number(s)>, and were accessed from the Kids First Data Resource Portal ( https://kidsfirstdrc.org and/or dbGaP (www.ncbi.nlm.nih.gov/gap).

# Kids First Publications Search Page

# Kids First Publications Search Page

https://commonfund.nih.gov/publications?pid=40

## Publications Search by Program

## Search Result

The search results on this publication page are automated on a monthly schedule based on acknowledgement of NIH Common Fund award numbers and intramural awards. Therefore, this list is not an exhaustive or error-free account of the program's publications.
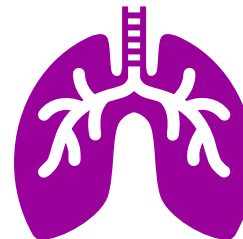
▾ **Gabriella Miller Kids First (13)**

Show [50 ▾] entries

Search: [          ]

| Publication Title | Authors | Journal | Publication Date | Page No | PubMedID |
|---|---|---|---|---|---|
| Whole genome sequencing of orofacial cleft trios from the Gabriella Miller Kids First Pediatric Research Consortium identifies a new locus on chromosome 21 | Mukhopadhyay N, Bishop M, Mortillo M, Chopra P, Hetmanski JB, Taub MA, Moreno L, Valencia-Ramirez LC, Restrepo C, Wehby GL, Hecht JT, Deleyiannis F, Butali A, Weinberg SM, Beaty TH, Murray JC, Leslie EJ, Feingold E, Marazita ML. | Human genetics. | 2019 Dec 17 | | 31848685 |
| Germline microsatellite genotypes differentiate children with medulloblastoma. | Rivero-Hinojosa, Samuel; Kinney, Nicholas; Garner, Harold R; Rood, Brian R | Neuro-oncology. | 2020 Jan 11; | | 31562520 |
| Germline 16p11.2 Microdeletion Predisposes to | Egolf, Laura E; Vaksman, Zalman; Lopez, Gonzalo; Rokita, Jo | American | 2019 Sep | | 31474320 |

# Kids First Investigators: Past Presentations

- **Congenital Diaphragmatic Hernia**, Wendy Chung (April 2019): https://www.youtube.com/watch?v=3CS6AphmCp0&t=978s

- **Neuroblastoma**, Sharon Diskin (September 2019): https://www.youtube.com/watch?v=Gq8kK2UGI4s

# Strategic Planning

## Progress on Addressing Key Challenges

# 7 Consensus Recommendation Themes

1. **Innovation: Resource, infrastructure, or tool development.**
   *Activities: Data Visualization tools; other tools for clinical/phenotypic data*

2. **Clinical/phenotypic data extraction, harmonization, & curation.**
   *Activities: Collect, extract, organize, curate, harmonize, and submit deep clinical and phenotypic data; annotate variants with pathogenicity, ClinGen scores.*
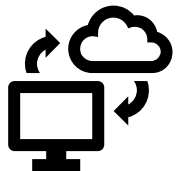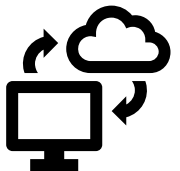
3. **Collaborative validation and discovery.**
   *Activities: Building synthetic cohorts; identify structural variants; test pipelines.*
   *Engage trainees in data analysis projects**Bring users to the platform*

4. **Integration and interoperability of external pediatric datasets.**
   *Activities: Using DRC workflow and best practices to harmonize external pediatric datasets; Building tools that can operate across multiple spaces*

5. **Consent and data sharing.**
   *Activities:  Re-consenting cohorts in line with our data sharing expectations*

6. **Validation with model organisms.**
   *Activities: validating KF findings/variants, deep phenotyping of animal models*

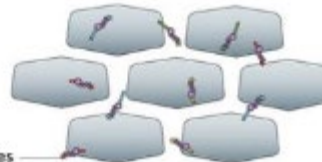7. **Continue WGS & data generation**, invest in long-read, consider other – omics. Reissues of: https://grants.nih.gov/grants/guide/pa-files/PAR-19-104.html

# 7 Consensus Recommendation Themes

**1. Innovation: Resource, infrastructure, or tool development.**
  *Activities: Data Visualization tools; other tools for clinical/phenotypic data*

**2. Clinical/phenotypic data extraction, harmonization, & curation.**
  *Activities: Collect, extract, organize, curate, harmonize, and submit deep clinical and phenotypic data; annotate variants with pathogenicity, ClinGen scores.*

**3. Collaborative validation and discovery.**
  *Activities: Building synthetic cohorts; identify structural variants; test pipelines.*
  *Engage trainees in data analysis projects**Bring users to the platform*

**4. Integration and interoperability of external pediatric datasets.**
  *Activities: Using DRC workflow and best practices to harmonize external pediatric datasets; Building tools that can operate across multiple spaces*

**5. Consent and data sharing.**
  *Activities:  Re-consenting cohorts in line with our data sharing expectations*

**6. Validation with model organisms.**
  *Activities: validating KF findings/variants, deep phenotyping of animal models*

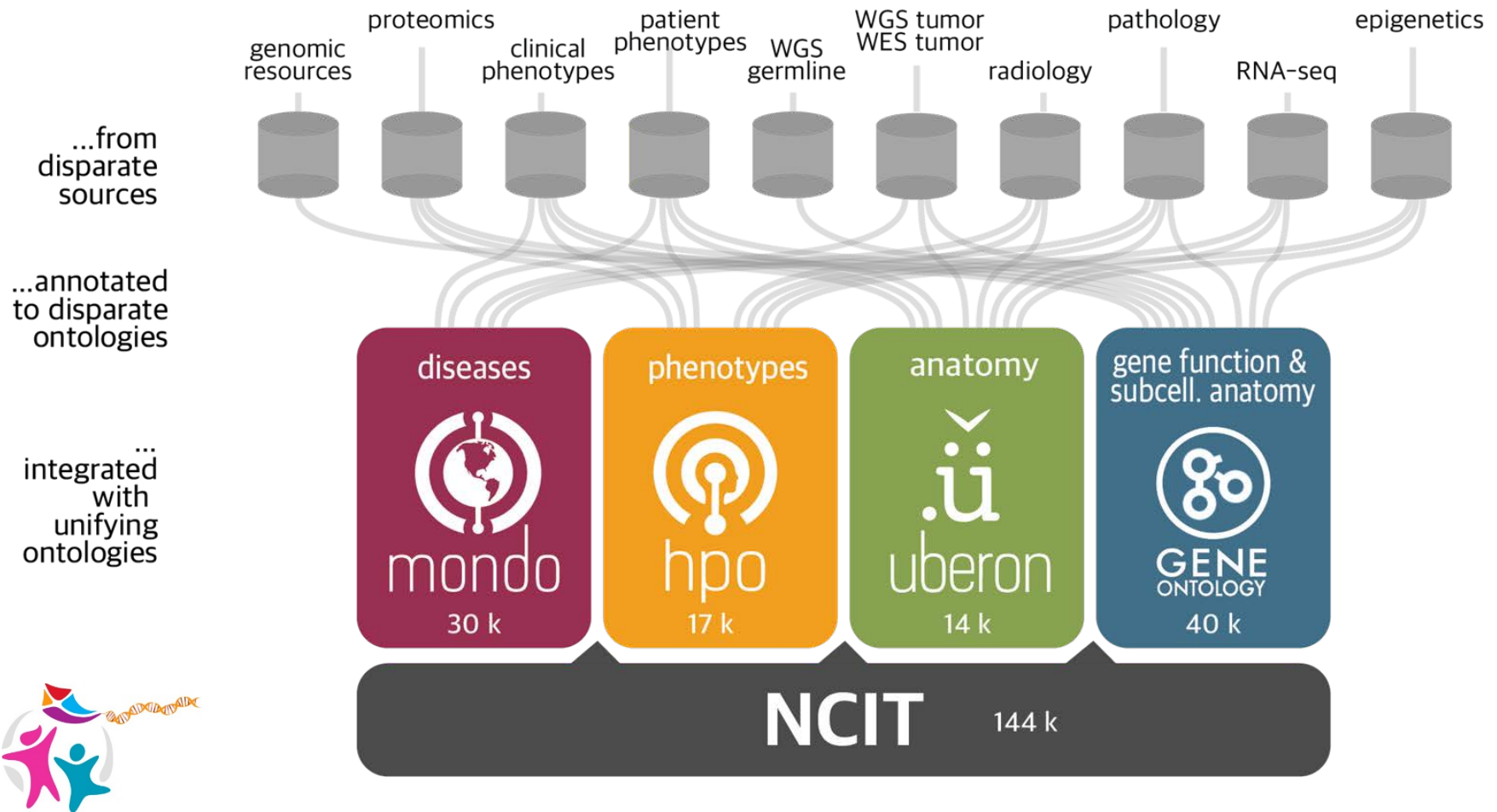**7. Continue WGS & data generation**, invest in long-read, consider other –
  omics. Reissues of: https://grants.nih.gov/grants/guide/pa-files/PAR-19-104.html

# 2020 Kids First X01
# Long Read Sequencing Pilot



PACIFIC BIOSCIENCES

OXFORD NANOPORE

# 7 Consensus Recommendation Themes Emerged

1. **Innovation: Resource, infrastructure, or tool development.**
   *Activities: Data Visualization tools; other tools for clinical/phenotypic data*

2. **Clinical/phenotypic data extraction, harmonization, & curation.**
   *Activities: Collect, extract, organize, curate, harmonize, and submit deep clinical and phenotypic data; annotate variants with pathogenicity, ClinGen scores.*

3. **Collaborative validation and discovery.**
   *Activities: Building synthetic cohorts; identify structural variants; test pipelines.*
   *Engage trainees in data analysis projects**Bring users to the platform*

4. **Integration and interoperability of external pediatric datasets.**
   *Activities: Using DRC workflow and best practices to harmonize external pediatric datasets; Building tools that can operate across multiple spaces*

5. **Consent and data sharing.**
   *Activities:  Re-consenting cohorts in line with our data sharing expectations*

6. **Validation with model organisms.**
   *Activities: validating KF findings/variants, deep phenotyping of animal models*

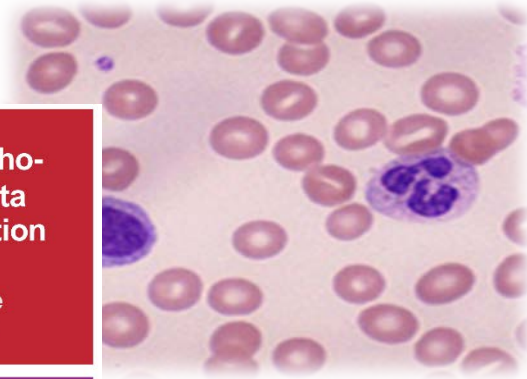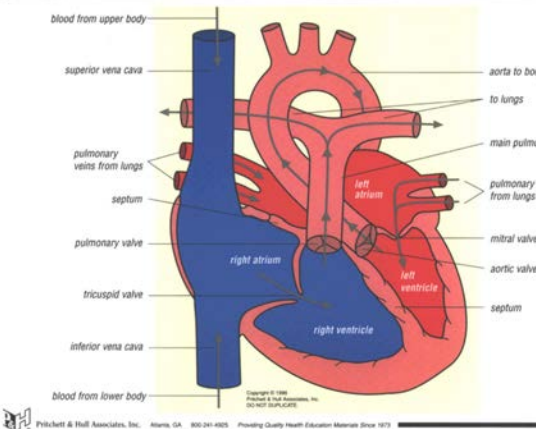7. **Continue WGS & data generation**, invest in long-read, consider other – omics. Reissues of: https://grants.nih.gov/grants/guide/pa-files/PAR-19-104.html

# *Innovation across the Phenotypic Translational Divide Webinar*

Information: https://monarch-initiative.github.io/phenomics/pages/clin-phen-webinar.html

Curation with ontologies that support heterogenous data types in Kids First

# *Innovation across the Phenotypic Translational Divide Webinar*



Webinar Information: https://monarch-initiative.github.io/phenomics/pages/clin-phen-webinar.html

# New ontology search and visualisation tools

*Sunburst: hierarchical view of participant counts per HPO term at all levels*

# 7 Consensus Recommendation Themes Emerged

1. **Innovation: Resource, infrastructure, or tool development.**
   *Activities: Data Visualization tools; other tools for clinical/phenotypic data*

2. **Clinical/phenotypic data extraction, harmonization, & curation.**
   *Activities: Collect, extract, organize, curate, harmonize, and submit deep clinical and phenotypic data; annotate variants with pathogenicity, ClinGen scores.*

3. **Collaborative validation and discovery.**
   *Activities: Building synthetic cohorts; identify structural variants; test pipelines.*
   *Engage trainees in data analysis projects**Bring users to the platform*

4. **Integration and interoperability of external pediatric datasets.**
   *Activities: Using DRC workflow and best practices to harmonize external pediatric datasets; Building tools that can operate across multiple spaces*

5. **Consent and data sharing.**
   *Activities: Re-consenting cohorts in line with our data sharing expectations*

6. **Validation with model organisms.**
   *Activities: validating KF findings/variants, deep phenotyping of animal models*
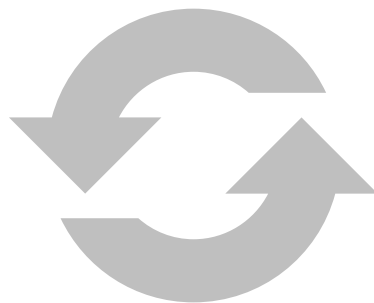
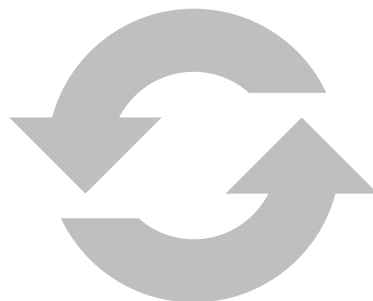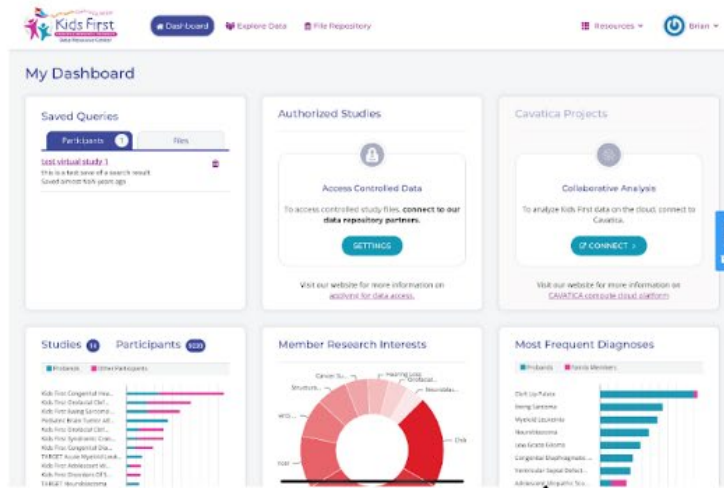7. **Continue WGS & data generation**, invest in long-read, consider other –omics. Reissues of: https://grants.nih.gov/grants/guide/pa-files/PAR-19-104.html

# NIH Cloud Based Platforms Interoperability (NCPI)

# NIH Cloud Based Platforms Interoperability efforts



TOPMed/PCGC

Undiagnosed Diseases Network, Centers for Mendelian Genomics, GTEx

TARGET

Kids First Data Resource Portal

Cavatica Workspace

**Any Portal to Any Workspace**

Terra Workspace

Other Portal

Other Workspace

**Goal: Empower end-user analyses across platforms through federation and interoperability**

*Borrowed from Brian O'Connor / Jack DiGiovanna*

# Layers of Interoperability

| Challenge | NCPI Activities |
|---|---|
| Operational barriers to trans-platform data sharing | Establish principles for promoting interoperability across multiple platforms. |
| Inability to search & access data across platforms | Test & implement technical standards for data exchange (e.g. GA4GH APIs) based on key use cases |
| Teach researchers to use the cloud | Create public "knowledge base" with training materials and cloud cost guide. |
| Lack of standards for clinical data exchange | Pilot and assess FHIR resources to model and share complex clinical and phenotypic data |

# NIH Researcher Auth Services (RAS)

Simplify researcher access to NIH data through
federated **authentication** (linking user identity account; "passport")
and **authorization** (claim to access specific studies/datasets; "visa")

https://datascience.nih.gov/data-infrastructure/researcher-auth-service



*Adapted from Susan Gregurick, ODSS*

# Collabsorations



*INvestigation of Co-occurring conditions across the Lifespan to Understand Down syndrome (INCLUDE)*

# Q & A

- Use the Q&A bar (lower right of your screen) to send your questions to "**All Panelists**".  We will read your questions out loud and answer them.

- You can ask also use the "chat" 💬 service to send private messages to the host or presenters.

**What funding opportunities are available?**

**How can I get involved?**

**How do I access data?**

# *What funding opportunities are available?*

See: FAQs for Funding Opportunities Announcements (FOAs) to Support Data Analyses of Kids First Datasets (https://commonfund.nih.gov/kidsfirst/FAQ)



- **Kids First cohort sequencing opportunity (X01):**
  - 1 more reissue of PAR-19-390 **for 2021**
- **Analyze Kids First data with support from:**
  - **"Kids First R03 PAR":** PAR-19-375
  - **NIH "Parent" R01**:PA-19-056
  - **NIH Parent R03**: PA-19-052
- **Validate variants with support from**:
  - ORIP's Development of Animal Models and Related Biological Materials for Research (R21): https://grants.nih.gov/grants/guide/pa-files/PA-16-141.html
  - Mechanistic Studies of Gene-Environment Interplay in Dental, Oral, Craniofacial, and Other Diseases and Conditions (R01) (PAR-19-292).
  - Development of Novel and Robust Systems for Mechanistic Studies of Gene-Environment Interplay in Dental, Oral, Craniofacial, and Other Diseases and Conditions (R21) (PAR-19-293).
  - To pursue collaborations with the Knockout Mouse Phenotyping Program (KOMP2), contact: KidsFirstKOMP@nih.gov

- **To receive updates about future Kids First opportunities, sign up for the listserv:**
  - https://commonfund.nih.gov/kidsfirst/register

# *How can I get involved?*

- **Connect with and provide <u>feedback</u> to the DRC:** [support@kidsfirstdrc.org](mailto:support@kidsfirstdrc.org)

- **Contact the program for questions or <u>feedback</u>**: [kidsfirst@od.nih.gov](mailto:kidsfirst@od.nih.gov)

- **Learn more about the program & DRC:** [https://commonfund.nih.gov/kidsfirst](https://commonfund.nih.gov/kidsfirst) & [https://kidsfirstdrc.org/](https://kidsfirstdrc.org/)

- **Search data available through the Kids First Data Resource Portal:** [https://portal.kidsfirstdrc.org/](https://portal.kidsfirstdrc.org/)

Community Members

Healthcare Professionals

Patients & Family Members

Researchers

# *How do I access data?*



**Anyone can register & login to the portal to filter, search, visualize datasets**

Submit dbGaP Data Access Requests (DARs) for **individual-level sequence data**

**Push approved sequence data to Cavatica from the portal:**
https://kidsfirstdrc.org/support/analyze-data/

**NIH Kids First Data Access Committee**

# Individual-level sequence data

- To learn more about submitting dbGaP Data Access Requests (DARs) watch:
https://www.youtube.com/watch?v=39cba0gF2tw&index=3&t=503s&list=PLoXwgZflAe4aMwWpVQU_WVeWHzyhI3BCu



**Also see:**
https://dbgap.ncbi.nlm.nih.gov/aa/dbgap_request_process.pdf

Submitting an Approvable
dbGaP Data Access Request
Vivian Ota Wang, Ph.D
Office of Data Sharing
NCI

# How are sequences released by the Kids First DRC?



**X01 Research Groups**
*Patient Phenotypic Data*
including diagnoses, demographics, pathology reports, radiology images…

**Sequencing Centers**
*Sequencing Data*
including whole genomes, whole exomes, transcriptomes…

**Kids First DRC**
*Harmonizes the Data*
Pairs sequences with the correct patient and double-checks everything

**Data Released**

**Kids First DRC**
*Releases the Data*

For each sample…
- Complete genome sequence (in .cram format)
- List of unique differences (in VCF and gVCF formats)

56

# *How can I interact with other community members?*

# *What community resources are available?*

# *How can I interact with other community members?*

# *What community resources are available?*

# *How can I interact with other community members?*

# *What community resources are available?*

# Q & A

- Use the Q&A bar (lower right of your screen) to send your questions to "**All Panelists**". We will read your questions out loud and answer them.

- You can ask also use the "chat"  service to send private messages to the host or presenters.

**What funding opportunities are available?**

**How can I get involved?**

**How do I access data?**

# Thank You!

Email Additional Questions and Comments to the Kids First Mailbox: **kidsfirst@od.nih.gov**