# GDC WGS Variant Calling Workflow Updates

## 31 March 2025

Bill Wysocki, Ph.D – GDC Director of User Services Zhenyu Zhang, Ph.D – GDC Director of Bioinformatics Center for Translational Data Science University of Chicago



### Agenda

- 1. Overview of Whole Genome Variants in the GDC
- 2. Available WGS Data
- 3. WGS Pipeline Overview
- 4. Downloading WGS Data
- 5. Documentation and Repositories
- 6. Questions

# Overview of Whole Genome Variants in the GDC

Somatic Mutations



## WGS Data in the GDC

- Tumor and normal samples are aligned
- The resulting alignment are processed with variant caller workflows of different types:
  - Point mutations/ Indels (SSM)
  - Copy number variants
  - Structural variants



There are two WGS workflow sets -

## Existing Sanger WGS Workflow Set

- 1. Point Mutation: CaVEMan
- 2. Indel: Pindel
- 3. CNV: ascatNGS
- 4. SV: BRASS

## New GDC WGS Workflow Set

- 1. Simple Somatic Mutation (SSM):
  - GATK4 MuTect2
  - VarScan2
  - SvABA Indel
  - Strelka2
- 2. Copy Number Variation (CNV):
  - GATK4 CNV + ABSOLUTE
- 3. Structural Variation (SV):
  - Manta
  - SvABA

## Summary of Major Changes

- 1. Increase SSM calling from 2 callers to 4, which enables proper variant ensemble.
- 2. Increase SV calling from 1 caller to 2 callers.
- **3**. Add ABSOLUTE for better CNV inference.
- 4. Stop Sanger workflow production, and run new WGS workflows across the GDC.

## WGS Data Releases in the GDC

- Data Release 27 (2020):
  - WGS variants were first released from the Sanger workflow set
- Data Release 42 (2025):
  - WGS variants were first released from the new GDC workflow set

#### Transition:

- All data that was processed with the Sanger workflow set will be processed with the new GDC workflow set
- WGS processing is in currently in progress, any pipeline that is completed will be released in the subsequent data release
- For SSMs (point + indels), this means that an annotated VCF is available



## Available WGS Data



# Currently Available in the GDC Data Portal

- 1. More than 15,000 cases with aligned WGS BAMs
- 2. 5,000+ cases with Sanger calls. We no longer use Sanger for production and will gradually release the rest of the Sanger calls that have already been completed.
- 3. Data from <14,000 cases of Manta, GATK4 CNV, GATK4 MuTect2, SvABA Indel callers have been released. We will release more in future releases and hopefully cover all WGS data in the GDC soon.

# WGS Pipeline Overview



## Why Switch Pipelines?

- 1. SSM calling: increased from 2 callers to 4, that enables multi-caller variant ensemble.
- 2. SV calling: increased from 1 caller to 2, that enables users to derive all confidence SV set.
- **3**. CNV calling: added ABSOLUTE for better CNV inference.
- 4. Efficiency: the new workflows, particularly modern tools like Strelka2 and Manta, require significantly fewer compute resources than the existing Sanger workflows.

## WGS Workflows





#### Output: VCF and BedPE

- 1. SvABA SV Caller
  - Generate both VCF and BedPE.
  - Also generate an Indel VCF for SSM ensemble.
- 2. Manta SV Caller
  - Generate both VCF and BedPE.
  - Also generate a candidate SSM VCF. This VCF is then fed into Strelka2 calling to enhance the quality and coverage of Strelka2.



*Output: raw VCF, annotated VCF, aliquot-level MAF, 4-caller ensemble MAF* 

#### 1. GATK4 MuTect2

- Different from the MuTect2 used in WXS calling (GATK3).
- Known issue: Missing variants on chr10 and chr20 in currently released VCFs. We expect to fix in later data releases.
- 2. VarScan2
  - The same caller used in WXS calling.
  - Not released yet.



Output: raw VCF, annotated VCF, aliquot-level MAF, 4caller ensemble MAF

- Implemented as Strelka2 Manta joint calling. Manta candidate VCFs is further filtered by Strelka2 calling to enhance the quality and coverage of Strelka2.
- 4. SvABA Indel
  - Indel output from the SvABA SV calling.



Output: Segmentation, gene-level CNV, purity/ ploidy estimation

- 1. GATK4 CNV caller
  - Non-integer copy number segmentation file (segmean).
  - Per request by genomics analysis groups, we also provide a controlled-access intermediate analysis archive tarball for expert manual review purpose.

#### 2. ABSOLUTE CNV Caller

- Use GATK4 CNV and GATK4 MuTect2 input to derive integerlevel copy number segmentation, model PDF, and gene-level copy numbers.
- Not released yet.



Output: Segmentation, gene-level CNV, purity/ ploidy estimation

- More on ABSOLUTE
  - ABSOLUTE is run in two steps: 1) "model generation" that generates multiple possible CNV models (purity and ploidy combinations); 2) "model extraction" that extracts the segmentations from one selected model.
  - We plan to provide a PDF report that contains all models from a run, and also automatic model extraction of the first model called "ABSOLUTE GATK4 CNV Auto". With curation provided by community experts, we may later release "ABSOLUTE GATK4 CNV Curated" when available.



## Other WGS Data and Workflows

- 1. BAM metrics, including coverage and other information. They are on the AlignedReads node, and available via API query, various facet filtering in GDC portal, and single BAM entity page.
- 2. Microsatellite Instability (MSI) by MSISensor2, at the same place as other BAM metrics, but only available in Tumor BAMs.
- 3. Tumor Purity and Ploidy. They are on the segmentation file node from ascatNGS and ABSOLUTE workflows, available via API.

## Development in Progress

- 1. Tumor Mutational Burden (TMB)
- 2. Tumor Mutational Signature (MuSiCal)

# Downloading WGS Data (Demo)



## **Goal: Download the following:**

Annotated Somatic Mutation VCFs (SNVs + Indels)
Manta BEDPE files

from cases with lung cancer

## GDC 2.0 Workflow



## GDC 2.0 Workflow – Downloading Files





NIH NATIONAL CANCER INSTITUTE GDC Data Portal	Vide	o Guides - 있구 Se	end Feedback 🛛 🖉 Browse Annotatior	ns 🕴 Manage	e Sets 🛛 📜 Cart 🧕 → 🕽 Login	GDC Apps 🔻			
Analysis Center 🎄 Projects	🚱 Cohort Builder 🛛 🥃 Repository				Q e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2				
り Unsaved_Cohort	• • 🔒 🔹 🕂 📋	± ±			✓ 44,736 C	CASES			
Cohort not saved									
	BEATAML1.0	882 (1.97%)	BEATAML1.0-COHORT	826 (1.85%)	acute lymphoblastic leukemia	1,086 (2.43%)			
General Diagnosis	CDDP_EAGLE	50 (0.11%)	BEATAML1.0-CRENOLANIB	56 (0.13%)	🗌 adenomas and adenocarcinom	14,344 (32.06%)			
<b>--</b>	CGCI	645 (1.44%)	CDDP_EAGLE-1	50 (0.11%)	🗌 adnexal and skin appendage ne	22 (0.05%)			
Disease Status and History	СМІ	299 (0.67%)	CGCI-BLGSP	324 (0.72%)	🗌 basal cell neoplasms	22 (0.05%)			
		1,687 (3.77%)	CGCI-HTMCP-CC	212 (0.47%)	blood vessel tumors	1 (0.00%)			
Disease Specific Classifications		😌 16 more		😌 80 more		😌 39 more			
Treatment	Primary Diagnosis	ର 🕼 🕽	Primary Site	୧ 🕼 ୨	Tissue or Organ of Origin	୯ 🕼 🕽			
Exposure	Name 🔺	Cases 💲	Name 🔺	Cases 💲	Name 🔺	Cases 💲			
	🗌 acinar adenocarcinoma	187 (0.42%)	accessory sinuses	1 (0.00%)	🗌 abdomen, nos	236 (0.53%)			
	🗌 acinar cell carcinoma	266 (0.59%)	🗌 adrenal gland	724 (1.62%)	🗌 adrenal gland, nos	608 (1.36%)			
Biospecimen	🗌 acinar cell tumor	31 (0.07%)	(0.07%) 🗌 anus and anal canal		ampulla of vater	4 (0.01%)			
	🗌 acute leukemia, nos	11 (0.02%)	🗌 base of tongue	27 (0.06%)	🗌 anal canal	1 (0.00%)			
Molecular Filters	🗌 acute lymphoblastic leukemia, n	144 (0.32%)	Dladder	767 (1.71%)	🗌 anterior floor of mouth	2 (0.00%)			
	acute lymphocytic leukemia	1,170 (2.62%)	🗌 bones, joints and articular cartila	257 (0.57%)	anterior mediastinum	40 (0.09%)			
Available Data		👴 194 more		🔁 63 more		🕒 194 more			
Custom Filters									
	Case ID	5							
	Upload Cases								

			Experimental Strategy	C [] D
NH) NATIONAL CANCER INSTITUTE	Video Guides	♀ Send Feedback	wgs	×
Analysis Center	le Cohort Builder		Name 🔺	Cases 🗘
J Unsaved_Cohort	•• • • • • •	Ŧ	✓ WGS	15,712 (35.12%)
Cohor Molecular Filters	t not saved Workflow Type २ []	ර Data Format		show less
Available Data	Name ▲     Case       Aliquot Ensemble Somatic Varian     2 (0)	s <b>\$</b> Name ▲		
Custom Filters	Anget NGS   5,444 (12.     BRASS   5,444 (12.     BRASS   5,444 (12.     CaVEMan   5,444 (12.     GATK4 NUTect2   15,271 (34.     CaVEMan   5,444 (12.     GATK4 CNV   13,473 (30.     GATK4 MuTect2   10,941 (24.     GATK4 MuTect2 Annotation   10,941 (24.     GATK4 MuTect2 Tumor-Only   2 (0.     GATK4 MuTect2 Tumor-Only Ann   2 (0.     Manta   5,710 (12.     Pindel   5,444 (12.     SvABA   5,528 (12.     SvABA Indel   1,013 (2.     VCF LiftOver   458 (1.	0.5% badpe   17% bdpe   17% maf   14% tar   17% tsv   12% txt   16% txt   16% 1   16% 1   16% 1   16% 1   16% 1   16% 1   17% 1   16% 1   16% 1   17% 1   16% 1   16% 1   12% 1   12% 1   12% 1	5,444 (12.17%) 13,728 (30.69%) 1 more	

Analysis Center	🎄 Projects 🔞 Cohort Builder 🥃 Repository			Repository Q e.g. B	Q e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2				
5 Unsaved_Cohor	t Cohort not	1 🔻	a   •	+ 🗊 🛨 🛨			✓ 10,945 CASES		
Unsaved_Cohort C	lear All							$\mathbf{\mathbf{x}}$	
EXPERIMENTAL STRA	TEGY ← WGS ×	× WORI	KFLOW TYPE	E ← GATK4 MuTect2 Anno × SvABA Indel Annotat × ×					
Filters									
Collapse All				Download Associated Data  Manifest View Images	🛓 Add Al	Files to Cart	🗧 Remove All Fi	rom Cart	
S + Add a C	ustom Filter	JSON	TSV	TOTAL OF 658,307 FILES \$ 10,945 CASES	<b>6.91</b> PB	<b>Q</b> Search			
へ Experimental Stra	ategy ۹.[)う	Cart A	Access 🌲	File Name 🍦	Cases	Project 🌲	Data Category 🌲	Data For	
Name ▲     Files ↓       ATAC-Seq     351 (0.03%)       Diagnostic Slide     9,088 (0.81%)       Expression Array     602 (0.05%)       Genotyping Array     109,439 (9.76%)       Methylation Array     33,417 (2.98%)       miRNA-Seq     35,064 (3.13%)		Controlled	11c9e0b9-c94c-466f-877a-65325a76dcaf.rna_seq.chimeric.gdc_realn.bam	<mark>Я</mark> <u>1</u>	J TCGA-KIRC	Sequencing Reads	BAM		
		Open	CGA-KIRC.7181379b-1d8a-4d22-818d-b3f7c20a38be.ascat3.allelic_specific.seg.txt	91	🔉 TCGA-KIRC	Copy Number Variation	ТХТ		
	109,439 (9.76%) 33,417 (2.98%)		Controlled	J TCGA-KIRC.997e4ed0-8d01-4901-acb3-50a8f7d15a3c.arriba.rna fusion.bedpe	<mark>9</mark> 1	DICGA-KIRC	Structural Variation	BEDPE	
	35,064 (3.13%)		Controlled	TCGA-KIRC fc76ae2d-e4d0-4ef0-b01e-e4c8e649c24h arriba rna fusion bedne	<b>I</b> 1		Structural Variation	REDPE	

# **Documentation and Repos**



## Documentation – https://docs.gdc.cancer.gov



## **Bioinformatics Pipeline Documentation and Resources:**

## https://github.com/NCI-GDC/gdc-workflow-overview







U.S. Department of Health & Human Services National Institutes of Health | National Cancer Institute

https://www.cancer.gov/

1-800-4-CANCER

Produced April 2024