# Washington University CPTAC3 Genomic Characterization Pipelines

Matthew A. Wyczalkowski, Liang-Bo Wang, Daniel Cui Zhou, Lijun Yao, Ruiyang Liu, Yige Wu, Wen-Wei Liang, Song Cao, Houxiang Zhu, Li Ding

Department of Medicine and McDonnell Genome Institute, Washington University in St. Louis, MO 63108, USA.
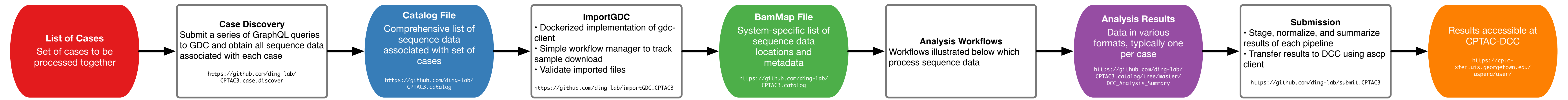
DingLab
dinglab.wustl.edu

## Data Acquisition and Processing

Each batch of cases is processed in a series of steps, beginning with the discovery of sample details at GDC, the import of data to local systems, running of the various workflows, and finally upload to DCC.
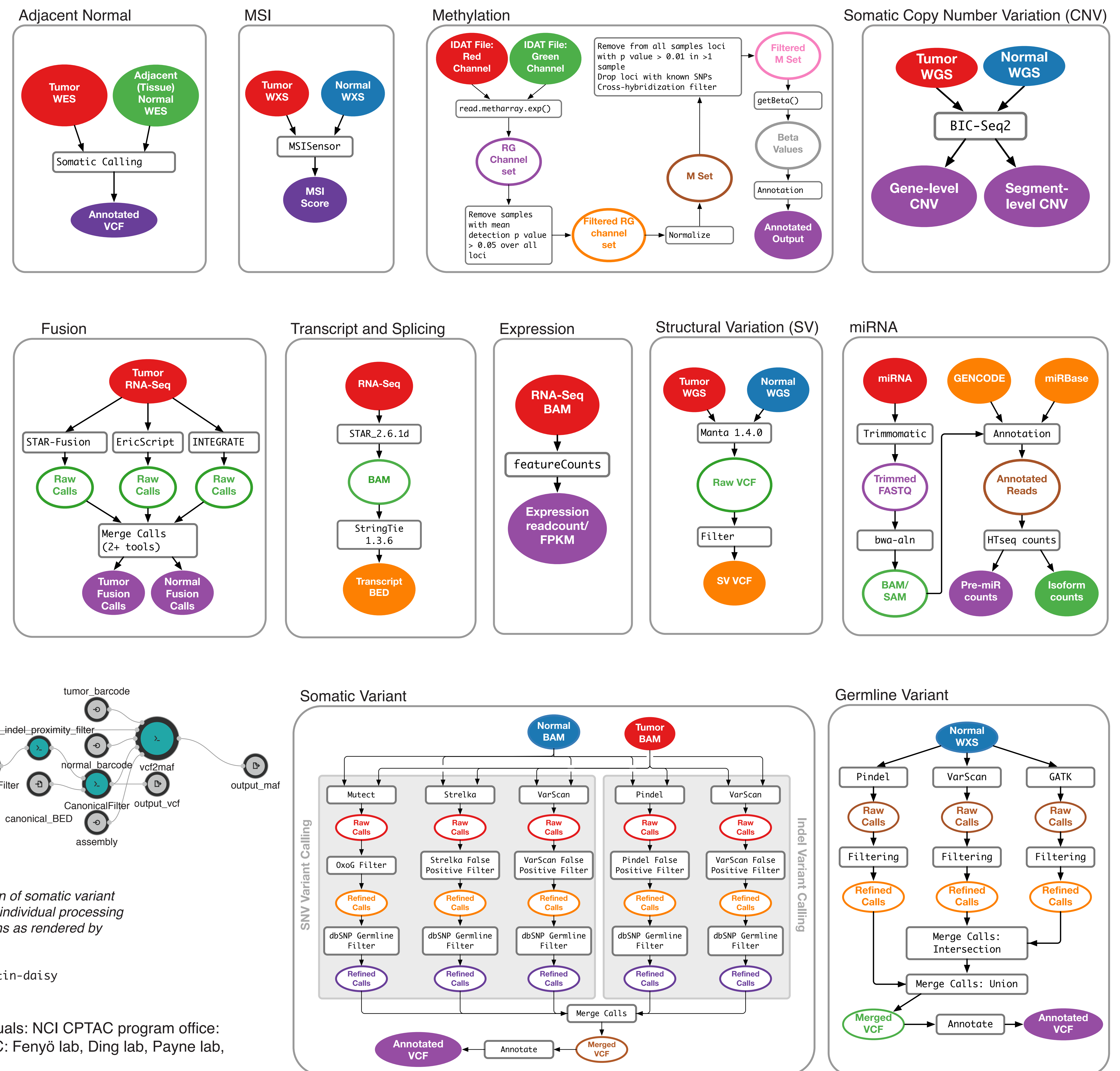


Data flow: List of Cases → Case Discovery → Catalog File → ImportGDC → BamMap File → Analysis Workflows → Analysis Results → Submission → Results accessible at CPTAC-DCC

- **List of Cases** — Set of cases to be processed together
- **Case Discovery** — Submit a series of GraphQL queries to GDC and obtain all sequence data associated with each case. https://github.com/ding-lab/CPTAC3.case.discover
- **Catalog File** — Comprehensive list of sequence data associated with set of cases. https://github.com/ding-lab/CPTAC3.catalog
- **ImportGDC** — Dockerized implementation of gdc-client; Simple workflow manager to track sample download; Validate imported files. https://github.com/ding-lab/importGDC.CPTAC3
- **BamMap File** — System-specific list of sequence data locations and metadata. https://github.com/ding-lab/CPTAC3.catalog
- **Analysis Workflows** — Workflows illustrated below which process sequence data
- **Analysis Results** — Data in various formats, typically one per case. https://github.com/ding-lab/CPTAC3.catalog/tree/master/DCC_Analysis_Summary
- **Submission** — Stage, normalize, and summarize results of each pipeline; Transfer results to DCC using ascp client. https://github.com/ding-lab/submit.CPTAC3
- Results accessible at CPTAC-DCC. https://cptc-xfer.uis.georgetown.edu/ppweb/user/

## Processing Summary

Table below details batches which have been submitted to DCC to date.

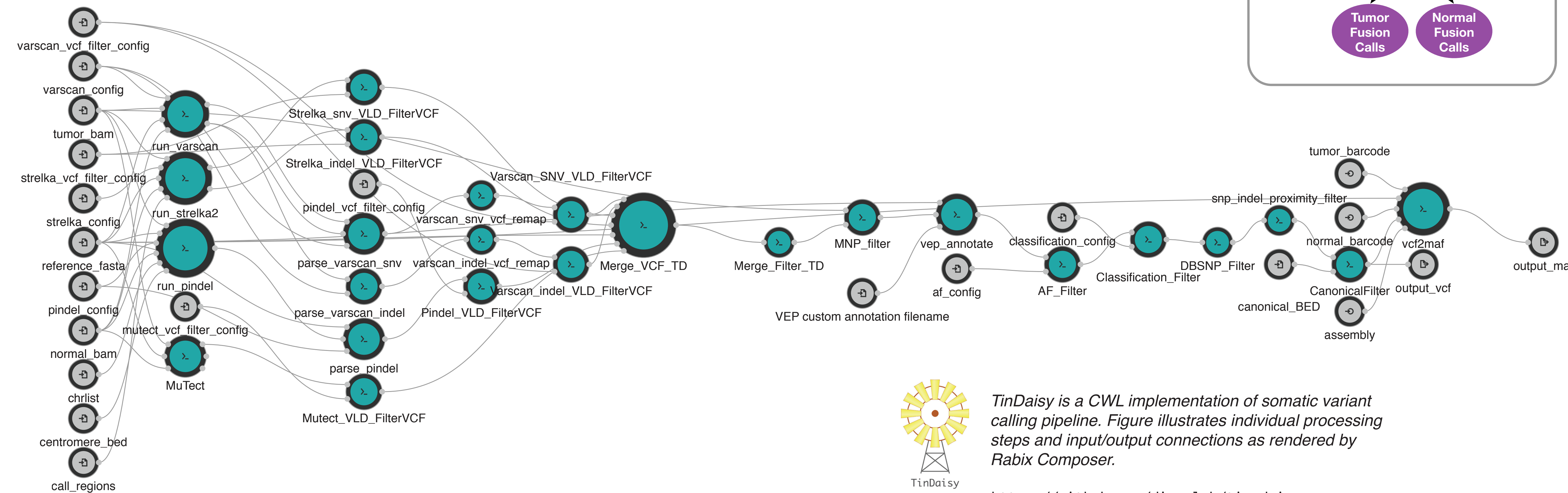| Pipeline | Data and Sample Type | CCRCC | GBM | HNSCC | LSCC | LUAD | PDA | UCEC | Total |
|---|---|---|---|---|---|---|---|---|---|
| Methylation Array | | 222 | 116 | 111 | 113 | 229 | 164 | 246 | 1201 |
| miRNA | | 222 | 114 | 111 | 113 | 229 | 164 | 247 | 1200 |
| Expression | | 225 | 120 | 118 | 117 | 230 | 164 | 251 | 1225 |
| Fusion | | 222 | 119 | 111 | 113 | 164 | 164 | 222 | 1115 |
| Transcript & Splicing | | 222 | 119 | 111 | 113 | 53 | 164 | 222 | 1004 |
| SV | | 222 | 59 | 109 | 113 | 111 | 166 | 217 | 997 |
| Somatic CNV | | 222 | 59 | 109 | 113 | 121 | 166 | 39 | 829 |
| MSI | | 222 | 118 | 111 | 113 | 111 | 166 | 244 | 1085 |

**Data Type** ● WGS   ○ WES   ◆ RNA-Seq   ▲ miRNA   ▦ Methylation array

**Sample Type** ■ Tumor   ■ Blood normal

## Looking Ahead

- Import and storage of very large datasets is the principal bottleneck for high throughput genomic analysis.
- Containerization and workflow definition languages are technologies which enable cloud-based computing, bringing pipelines to the data rather than the other way around.
- We have implemented a number of pipelines in Docker / CWL. Somatic and germline variant calling is currently being processed using somatic calling pipeline on Cromwell workflow engine.



TinDaisy is a CWL implementation of somatic variant calling pipeline. Figure illustrates individual processing steps and input/output connections as rendered by Rabix Composer.

https://github.com/ding-lab/tin-daisy

## Acknowledgments

## Pipeline Details



Pipeline detail diagrams including: Adjacent Normal, MSI, Methylation, Somatic Copy Number Variation (CNV), Fusion, Transcript and Splicing, Expression, Structural Variation (SV), miRNA, Somatic Variant, and Germline Variant.