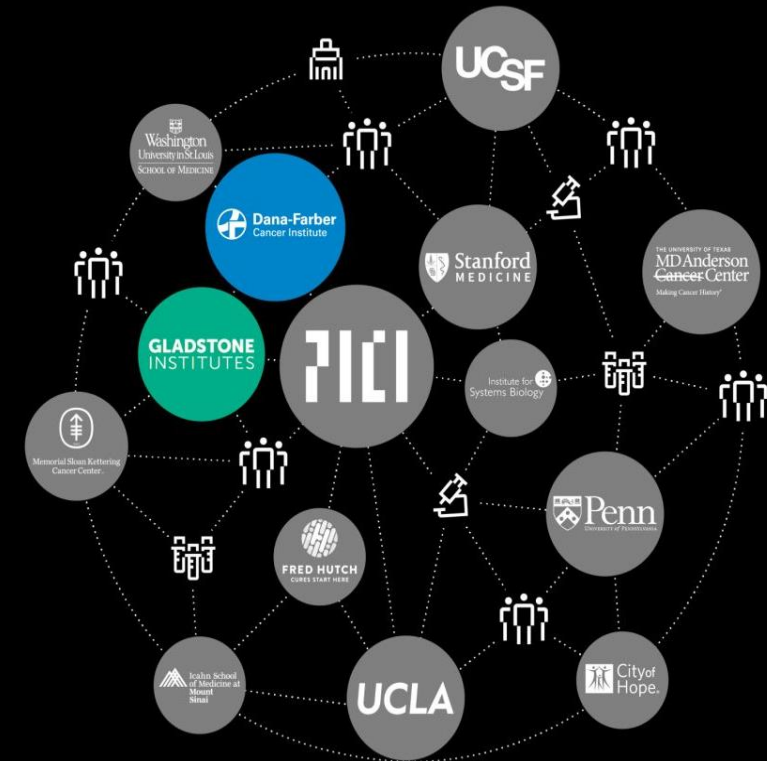# PICI INTERNSHIP 2025

## CATERINA PONTI

Data Science and Bioinformatics Intern

# PARKER INSTITUTE FOR CANCER IMMUNO THERAPY (PICI)

Our mission is to accelerate the development of breakthrough immune therapies to turn all cancers into curable diseases.
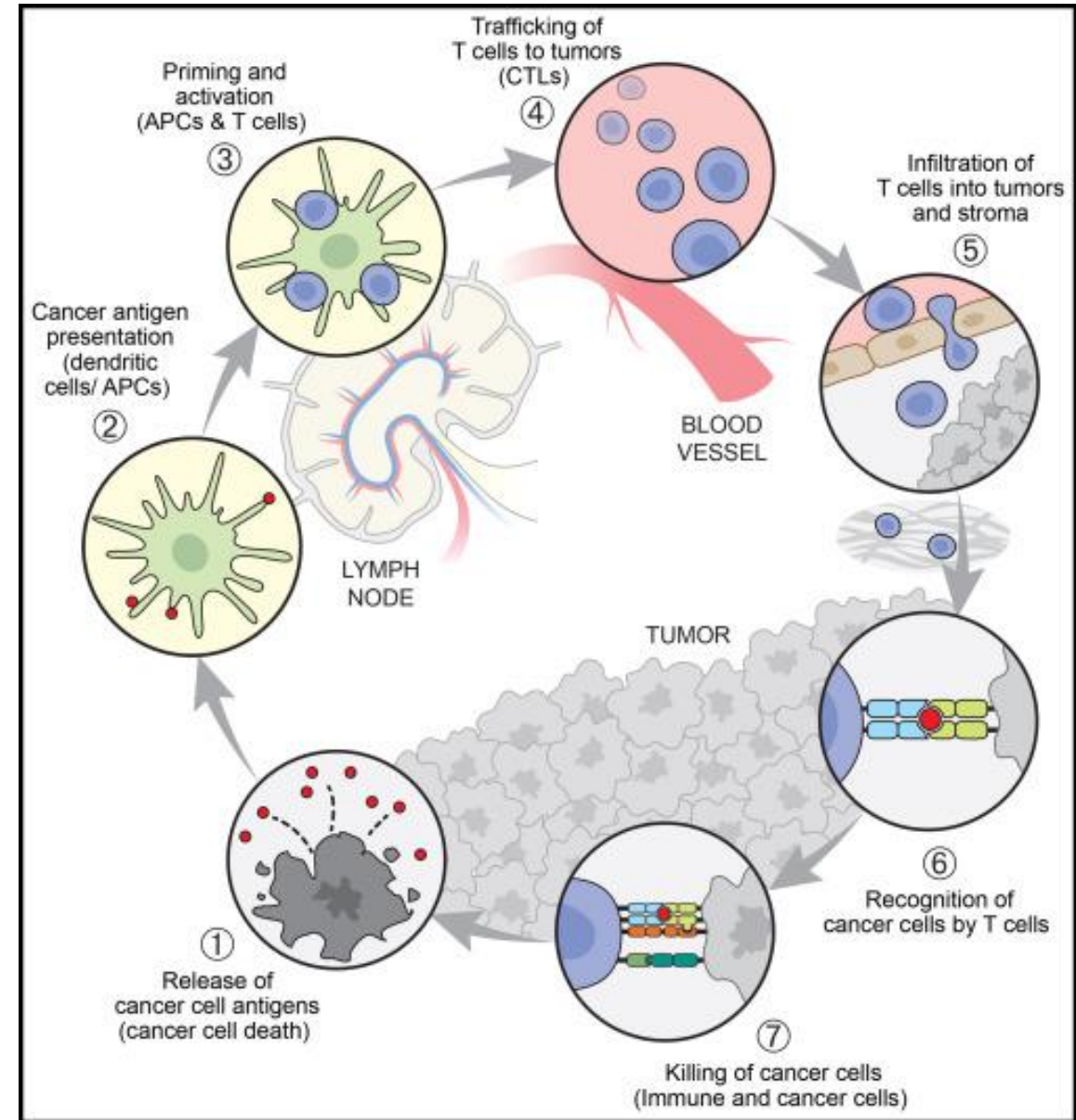
WHAT IS PICI?

- Patient-centric, grant-giving non-profit organization

- Fund & accelerate high-impact research

- Drive collaboration (700+ Cross-Institutional PICI Investigators)

- Enable bold Innovations

# THE CANCER-IMMUNITY CYCLE

1. **Release:** Dying cancer cells release "markers" (Antigens).

2. **Presentation:** Specialized scout cells (Dendritic cells) capture these markers.

3. **Priming:** Scouts travel to lymph nodes to "teach" T-cells what the cancer looks like.

4. **Trafficking:** Newly trained T-cells enter the bloodstream to find the tumor.

5. **Infiltration:** T-cells exit the blood and "invade" the tumor tissue.

6. **Recognition:** T-cells identify the specific cancer cells using the markers from Step 1.

7. **Killing:** T-cells destroy the cancer cells, which releases *more* markers, restarting the cycle at Step 1.

Where the System Breaks Down: Cancer cells can **turn off** T cells by using "brakes" on the immune system called immune checkpoints -> this stops T cells from attacking
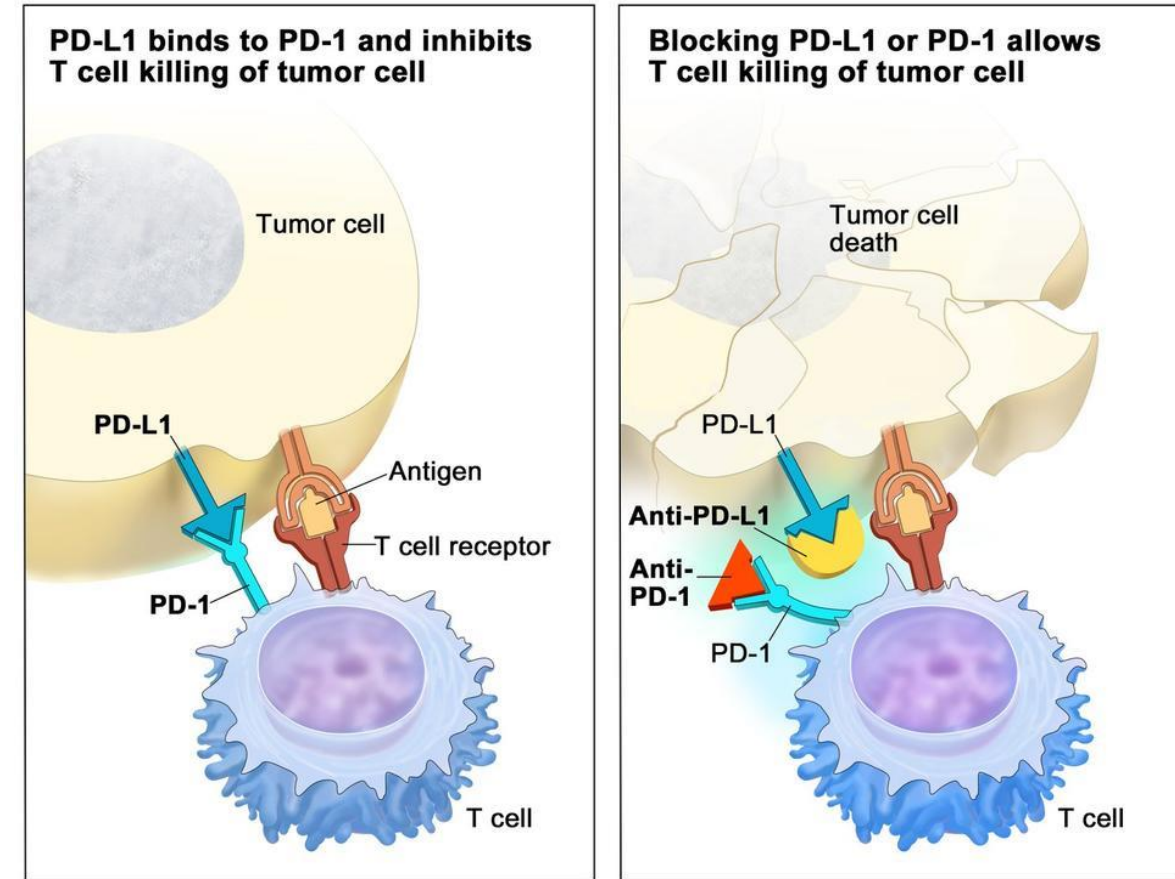
# IMMUNE CHECKPOINT INHIBITOR (ICI) THERAPY

Goal: To restart the Cancer-Immunity Cycle by removing the biochemical "brakes" (checkpoints) that cancer uses to survive.

- Anti–PD-1 / Anti–PD-L1 (e.g., pembrolizumab, atezolizumab)

- Anti–CTLA-4 (e.g., ipilimumab)

Clinical Challenges

- Immune-related Adverse Events (iRAEs) Because we are "releasing the brakes," the immune system can become overactive..

- It may mistakenly attack healthy organs (e.g., lungs, skin, or colon), leading to inflammation.



PD-L1 binds to PD-1 and inhibits T cell killing of tumor cell

Tumor cell

PD-L1

Antigen

T cell receptor

PD-1

T cell

Blocking PD-L1 or PD-1 allows T cell killing of tumor cell

Tumor cell death

PD-L1

Anti-PD-L1

Anti-PD-1

PD-1

T cell

© 2015 Terese Winslow LLC
U.S. Govt. has certain rights

# RADIOHEAD Study Planning

## Study Support and Goals:

- Funded by the Helmsley Charitable Trust, JDRF, and Bristol Myers Squibb.

- Goal: Identify **biomarkers** to prevent or intervene in:
  - Cancer immunotherapy-induced **Type 1 Diabetes**.
  - Severe **Immune-Related Adverse Events (irAEs)**.

## Academic Guidance and Study Direction:

- UCSF: data collection and trial objectives.
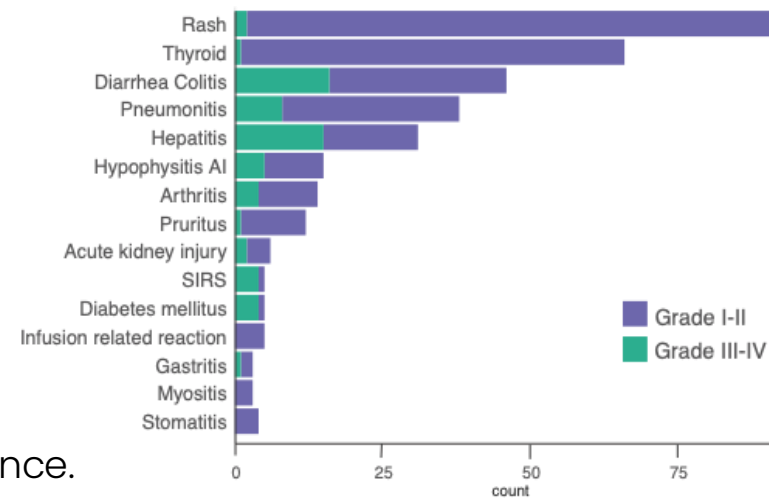- Huntsman Cancer Institute: data cleaning and curation work.

## Key Challenges:

- Recruitment for a Rare Event was difficult: 12 Expected, 3 Recruited.

## Root Cause Analysis:

- Patient Behavior
- COVID Impact
- Technical Issues
  - Difficulty in diagnosing T1D in community centers lacking specialized experience.

**New Focus:** Association between all Immune-Related Adverse Events (irAEs) and Cancer Treatment Response.
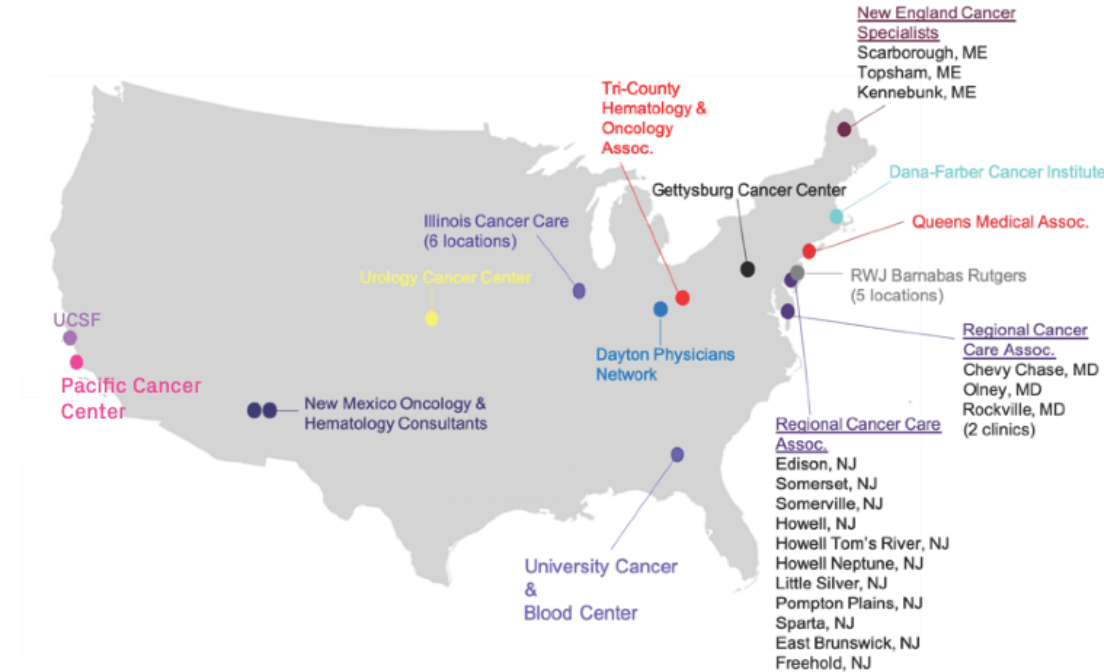
# RADIOHEAD Sample Collection

- **Standardized Processes**: All blood samples shipped overnight to ensure consistency across sites.

- **Central Laboratory**: Usage of a single central facility to eliminate processing variables.

- **Site Selection Strategy**
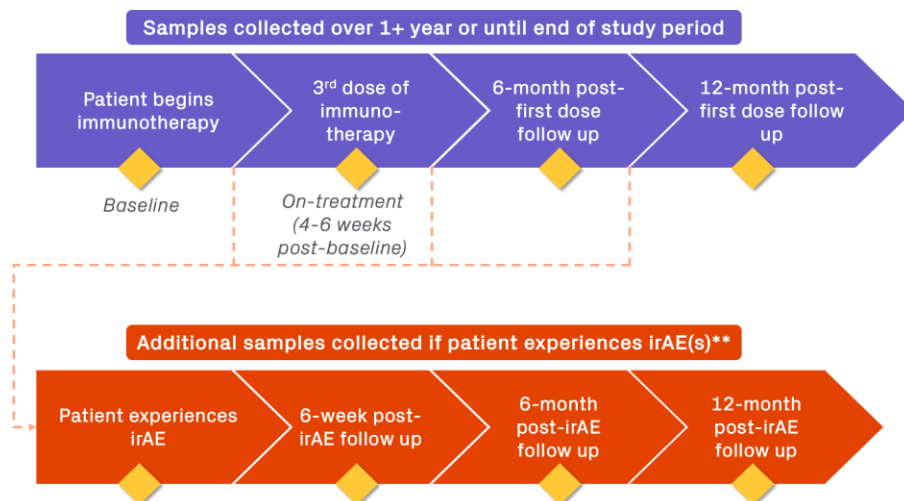
- **Batch Effect Analysis**

## Systemic Data Biases

- Geographic Bias: study privileges expensive, resource-rich urban hubs

- Demographic Gap: underrepresentation of minority

- Economic Barriers: Cost of Access

- IRAE Potentially underreported

- Lack of Longitudinal Data (patient follow up is inconsistent because of pandemic)

- Clinical Limitations: additional disease information was often a free textbox



New England Cancer Specialists
Scarborough, ME
Topsham, ME
Kennebunk, ME

Tri-County Hematology & Oncology Assoc.

Dana-Farber Cancer Institute

Gettysburg Cancer Center

Illinois Cancer Care (6 locations)

Queens Medical Assoc.

Urology Cancer Center

RWJ Barnabas Rutgers (5 locations)

UCSF

Dayton Physicians Network

Regional Cancer Care Assoc.
Chevy Chase, MD
Olney, MD
Rockville, MD
(2 clinics)

Pacific Cancer Center

New Mexico Oncology & Hematology Consultants

Regional Cancer Care Assoc.
Edison, NJ
Somerset, NJ
Somerville, NJ
Howell, NJ
Howell Tom's River, NJ
Howell Neptune, NJ
Little Silver, NJ
Pompton Plains, NJ
Sparta, NJ
East Brunswick, NJ
Freehold, NJ

University Cancer & Blood Center

# What Makes RADIOHEAD *Unique* ?

- Scale and Data:
  - Data collected from 1,070 patients (~2.5 draws per patients).
  - Generated over 70,000 samples across 3,500+ combinations of timepoints.
  - Samples were collected over one year.

- Event-Driven Sampling
  - Additional samples collected in the event of an adverse event (irAE)



| 1070 Patients 3500+ Timepoints | 70K+ Sample Aliquots Banked |
|---|---|
| Patient Information | Whole Blood |
| Treatment Detail | Plasma |
| Outcomes Data | Serum |
| Additional Metadata | PBMC |

# What Makes RADIOHEAD *Unique* ?

Distribution of patient tumor types in RADIOHEAD dataset
Number of patients



| Comprehensive Multiomic Analyses | | |
|---|---|---|
| ctDNA | ~750 genes, completed | GUARDANT |
| Serum Proteomics | 600+ proteins, completed | nomic |
| WES / SNP-panel | completed | Genentech *A Member of the Roche Group* / Teiko.bio |
| HD Flow Cytometry | 80 markers, in progress (complete by mid-2025) | BostonGene |
| Bulk-RNAseq | In progress (complete by mid-2025) | |
| Single cell RNA-seq | Projected to start in April 2025 | immunai |
| Metabolomics | Evaluating partners | |

# BIOREPOSITORY



BIOREPOSITORY 2024 PULL

- The biorepository includes samples from 10 additional PICI-owned clinical trials.

- Limitations: inconsistent naming conventions and duplicate entries.

    - I developed a **data clean-up and standardization script** to support future data pulls.

- RADIOHEAD (PICI-009) is a high-priority study due to its external partnerships and multi-omics integration capabilities.

# RADIOHEAD DATASETS

## Biorepository

Multiple pulls: keeping track of sample shipments

**REDCap**: clinical annotations + samples' barcodes of patients (1,300) collected during clinical trial

**df_clean**: cleaned REDCap with 1,070 patients of all patients with a pretreatment sample



Sample Material Types in PICI 009 Biorepository
(from df_clean patients: 83.3% of all samples)

# BARCODES MATCHING AND OUTLIER DETECTION (RADIOHEAD)

Sample Material Types in Biorepository from the 1070 RADIOHEAD patients selected in df_clean.



- Cleaned and standardized barcodes using custom parsing rules.

- Detected outliers (barcodes assigned to the wrong project name) using pairwise Levenshtein distances (via RapidFuzz) via RapidFuzz to compare alphanumeric barcodes.

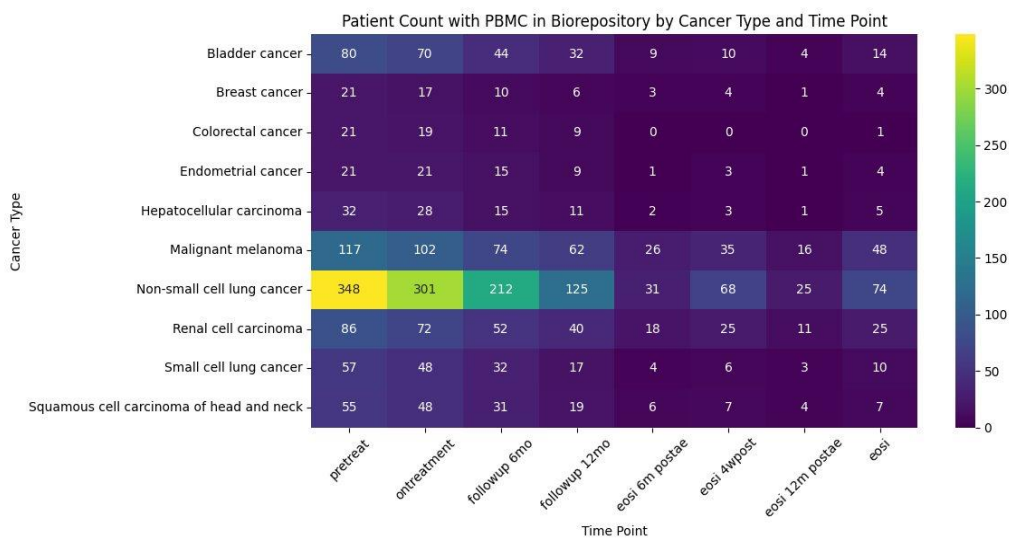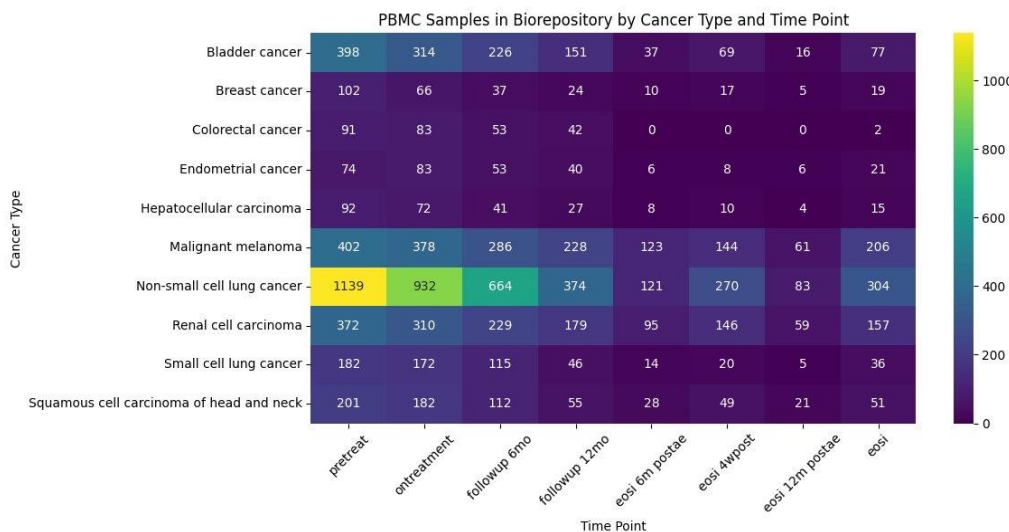- Matched barcodes between biorepository and REDCap using exact and substring matching.

Actual vs Expected Tubes by Phase

# SAMPLE COUNT AND PATIENT BY TIMEPOINT ACTUAL VS. EXPECTED (FROM REDCAP)

- Gap between the expected samples tracked in REDCap and the actual samples available in the biorepository as of 2025

- Percentage of expected tubes collected per patient during clinical trial vs samples available now in the biorepository

PBMC Samples in Biorepository by Cancer Type and Time Point

| Cancer Type | pretreat | ontreatment | followup 6mo | followup 12mo | eosi 6m postae | eosi 4wpost | eosi 12m postae | eosi |
|---|---|---|---|---|---|---|---|---|
| Bladder cancer | 398 | 314 | 226 | 151 | 37 | 69 | 16 | 77 |
| Breast cancer | 102 | 66 | 37 | 24 | 10 | 17 | 5 | 19 |
| Colorectal cancer | 91 | 83 | 53 | 42 | 0 | 0 | 0 | 2 |
| Endometrial cancer | 74 | 83 | 53 | 40 | 6 | 8 | 6 | 21 |
| Hepatocellular carcinoma | 92 | 72 | 41 | 27 | 8 | 10 | 4 | 15 |
| Malignant melanoma | 402 | 378 | 286 | 228 | 123 | 144 | 61 | 206 |
| Non-small cell lung cancer | 1139 | 932 | 664 | 374 | 121 | 270 | 83 | 304 |
| Renal cell carcinoma | 372 | 310 | 229 | 179 | 95 | 146 | 59 | 157 |
| Small cell lung cancer | 182 | 172 | 115 | 46 | 14 | 20 | 5 | 36 |
| Squamous cell carcinoma of head and neck | 201 | 182 | 112 | 55 | 28 | 49 | 21 | 51 |

Patient Count with PBMC in Biorepository by Cancer Type and Time Point

| Cancer Type | pretreat | ontreatment | followup 6mo | followup 12mo | eosi 6m postae | eosi 4wpost | eosi 12m postae | eosi |
|---|---|---|---|---|---|---|---|---|
| Bladder cancer | 80 | 70 | 44 | 32 | 9 | 10 | 4 | 14 |
| Breast cancer | 21 | 17 | 10 | 6 | 3 | 4 | 1 | 4 |
| Colorectal cancer | 21 | 19 | 11 | 9 | 0 | 0 | 0 | 1 |
| Endometrial cancer | 21 | 21 | 15 | 9 | 1 | 3 | 1 | 4 |
| Hepatocellular carcinoma | 32 | 28 | 15 | 11 | 2 | 3 | 1 | 5 |
| Malignant melanoma | 117 | 102 | 74 | 62 | 26 | 35 | 16 | 48 |
| Non-small cell lung cancer | 348 | 301 | 212 | 125 | 31 | 68 | 25 | 74 |
| Renal cell carcinoma | 86 | 72 | 52 | 40 | 18 | 25 | 11 | 25 |
| Small cell lung cancer | 57 | 48 | 32 | 17 | 4 | 6 | 3 | 10 |
| Squamous cell carcinoma of head and neck | 55 | 48 | 31 | 19 | 6 | 7 | 4 | 7 |

# SAMPLES AND PATIENT COUNT BY TIMEPOINT AND CANCER TYPE

1. Counts unique samples and patients remaining in the biorepository, grouped by cancer type and timepoint.

2. Provides a clear view of sample availability across disease groups and visits.

3. Helps scientists quickly identify gaps, trends, and well-represented cohorts when planning analyses or selecting samples.

# WHY MOVING FROM SCIRPTS TO A WEB PLATFORM?

- Internal Efficiency: Automated scripts to replace manual queries
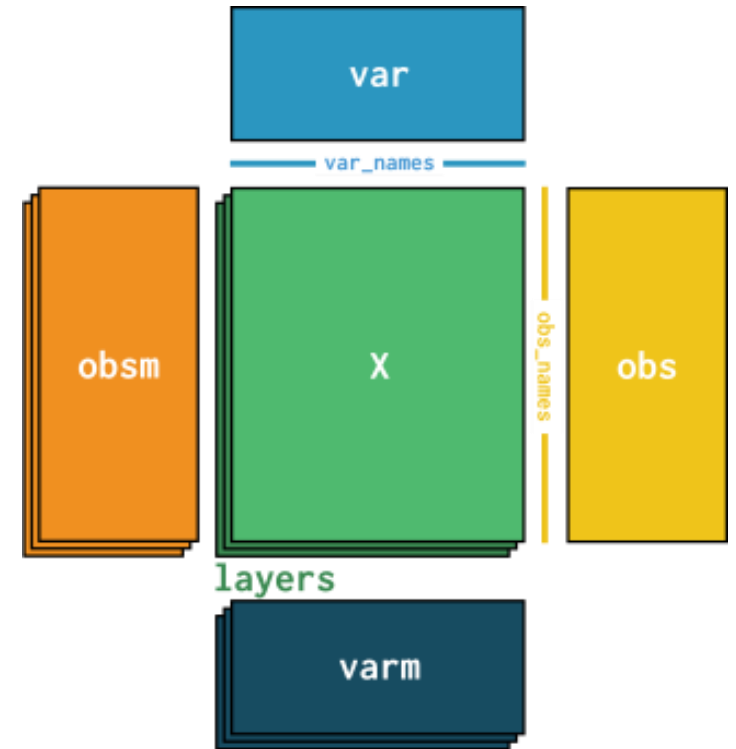- Partner Collaboration
- Real-Time Updates
- User Friendly

# Mass Cytometry Data (CyTOF) from Teiko

## Technology and Overview:
- Uses **antibodies** to quantify up to 50 single-cell biological markers.
- Partnership with **Teiko** completed data acquisition for all melanoma patients.

## Raw Data:
- Each sample ~10,000 cells
- High dimensionality: 50 markers x 10,000 cells x 500 samples
- **AnnData** (Annotated Data) object



| Parse FCS into AnnData | → | Downsample cells per sample | → | Marker-Aware downsampling |

# The Cloud Solution

**Challenges**
- High-dimensional data → long runtimes & heavy compute
- Running on a laptop took ~8 hrs
- Local workflow limits collaboration

**The Cloud Solution**
- What is the Cloud?
- Remote compute & storage instead of local machine

**Why the Cloud**
- Scales analysis beyond local computing limits.
- Enables fast processing of large single-cell datasets.
- Supports reproducible, shareable pipelines.

# Cloud Workflow Architecture

**Cloud Environment Setup**
- Linux Virtual Machine (Ubuntu)

**Data Access and Storage**
- **GCS Bucket** mounted locally to the Virtual Machine for data and workflows access.

**Data Processing and Analysis**
- Running pipelines for parsing (FCS to AnnData) and data analysis.

# Analytical Workflow

**1. Data Processing**

- Preprocessing and normalization of raw data using **Scanpy** and **Anndata** to ensure quality.

**2. Dimensionality Reduction**

- Application of **UMAP** to project high-dimensional marker data into interpretable 2D embeddings.

**3. Unsupervised Clustering**

- Utilizing **Leiden** (graph-based) and **FlowSOM** (SOM-based) algorithms to identify distinct cell populations.

**4. Visualization & Output**

- Generation of marker expression heatmaps, cluster comparison plots, and interactive visualizations.

# CLUSTER VALIDATION: SILHOUETTE ANALYSIS

Clustering: Cells were clustered using unsupervised methods (Leiden / FlowSOM), with K = 5 clusters selected for evaluation.

For every cell (data point) $i$, we calculate the Silhouette Score based on two distances:

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

- **Intra-cluster distance ($a$):** Distance to cells in its **own** cluster (should be small).
- **Inter-cluster distance ($b$):** Distance to cells in the **nearest neighbor** cluster (should be large).

Interpreting the Silhouette Plot
- **Each Line = One Cell:** The plot is made of horizontal bars; every single line represents an individual cell's score.
- **The "Knife" Shape:** You want thick, "knife-shaped" blocks. This indicates that most cells have high scores and are well-assigned.
- **The Threshold:** The vertical dashed line is the **average score**. You want most of your "knives" to extend past this line.



Silhouette analysis for KMeans clustering (n_clusters = 5)

The silhouette plot for the various clusters.

Cluster label

The silhouette coefficient values

The visualization of the clustered data.

UMAP2

UMAP1

# HOW TO DETERMINE THE OPTIMAL NUMBER OF CLUSTERS (K)?

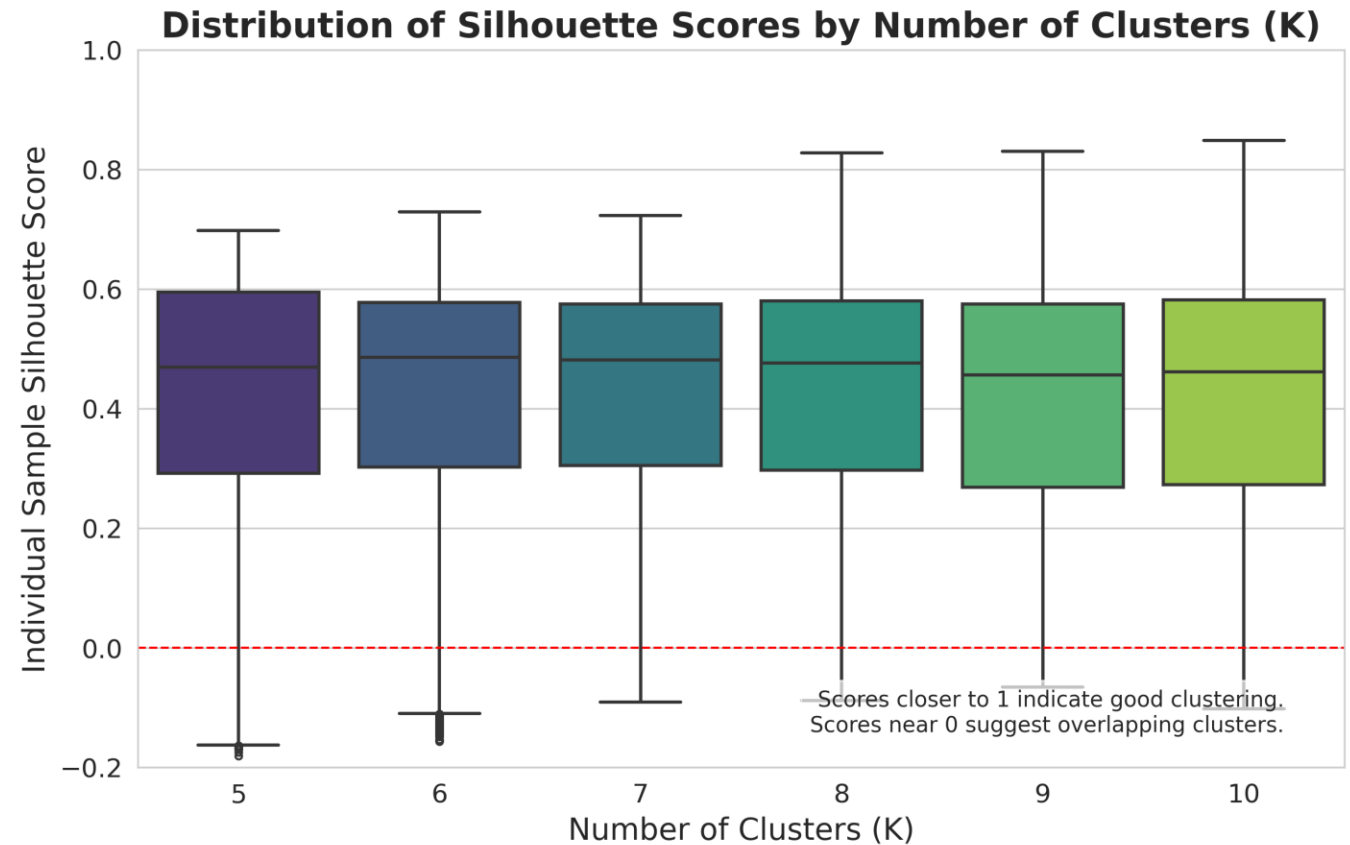## SILHOUETTE PLOTS FOR K MEANS CLUSTERING ANALYSIS

# HOW TO DETERMINE THE OPTIMAL NUMBER OF CLUSTERS (K)?

What is a **Silhouette Scores** Boxplot?

Represents the distribution of the silhouette scores for every data point across a range of cluster counts (K)
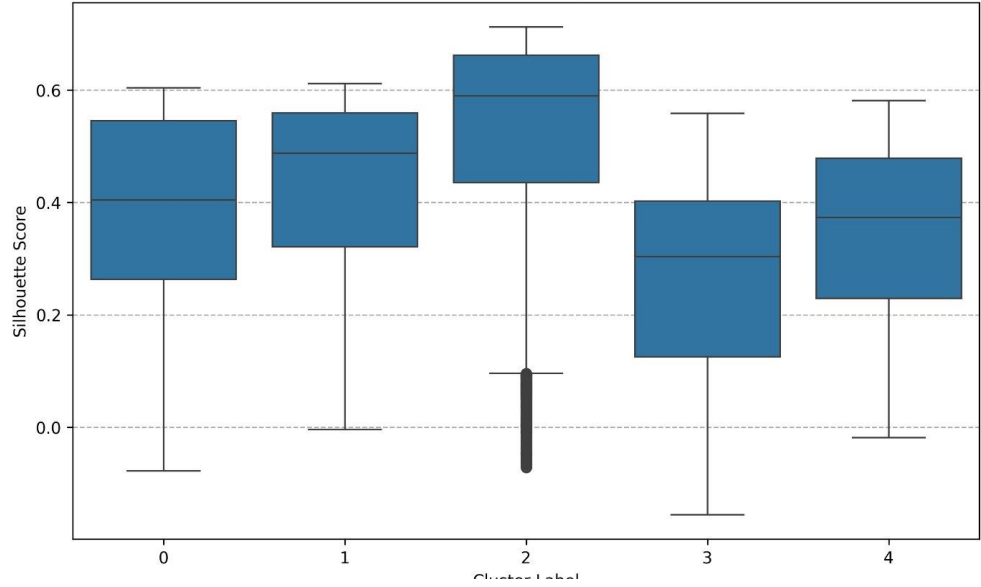
Optimal K:

- **Median of Each Silhouette Score**: higher indicates better separation on average

- **Interpreting Cluster Cohesion** (Variability/IQR): less spread suggests more consistent cluster assignments across samples
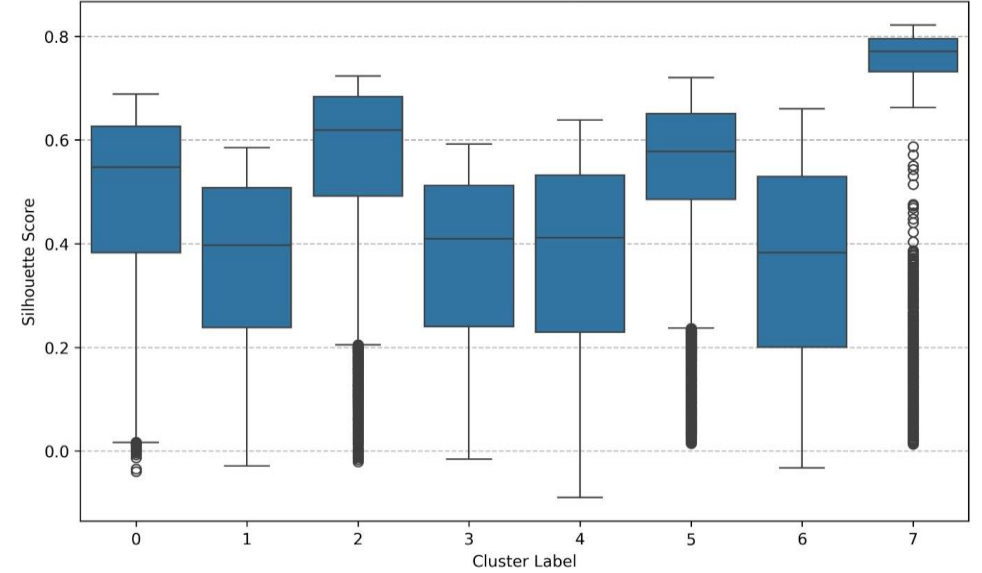


**Distribution of Silhouette Scores by Number of Clusters (K)**

Scores closer to 1 indicate good clustering.
Scores near 0 suggest overlapping clusters.

# HOW GOOD IN AVERAGE ARE MY CLUSTERS?
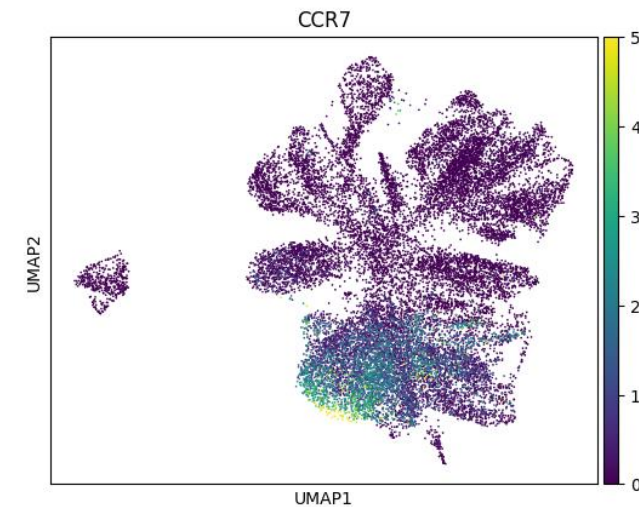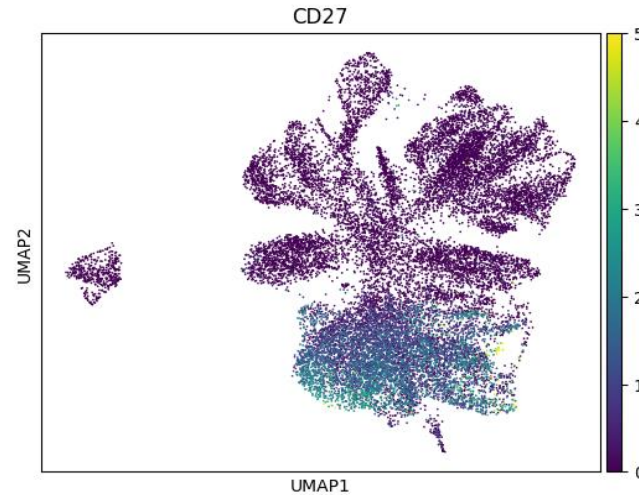
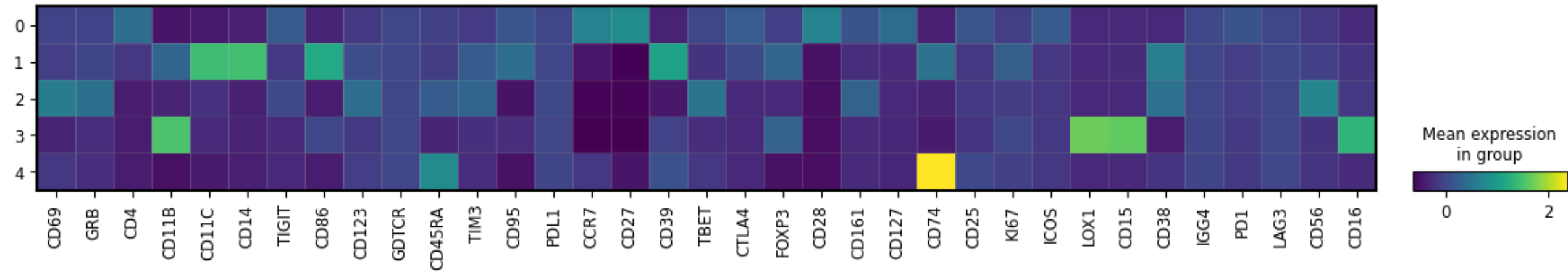# Marker Expressions Across Treatment Stages

## Pre Treatment



## On Treatment



- **CD27**: co-stimulatory receptor on CD8+ T cells, essential for their activation, clonal expansion, survival, and long-term memory formation in anti-tumor and anti-viral immunity

- **CCR7**: Chemokine receptor marking naive and central-memory CD8+ T cells.

HOW CAN A COMPUTATIONAL APPROACH SUPPORT CELL MARKER CLUSTERING? HOW DATA ANALYSIS CAN SUPPORT BIOLOGY RESEARCH?

- Unbiased gating and clustering: can reveal new cell populations and relationships

- Validate biological gating

- Handling High Dimensionality and Automated Clustering
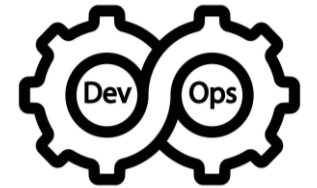
# Tools and Technologies

## Analysis & Parsing

- Pandas
- Anndata
- Scanpy
- Fcsparser
- pytometry

## DevOps & Environment

- Linux (Ubuntu)
- Miniconda
- Pyenv
- Git/GitHub
- Bash Scripting

## Web & Visualization

- Matplotlib
- Altair
- Flask
- SQLite3

# Next Steps

Apply mixed linear models to appropriately account for longitudinal data.

Integrate Circulating Tumor DNA (ctDNA) and mass cytometry together in analysis

Perform clustering analyses on additional cell types

Explore determinants of immune checkpoint responses and iRAEs

Run analyses in parallel across samples, donors, or batches.

# THANK YOU!

CHECK OUT MY CODE:
https://github.com/caterinaponti/radiohead-pici-internship.git

# CONTACT ME

Caterina personal – cponti@dons.usfca.edu

LinkedIn – https://www.linkedin.com/in/caterina-ponti

# SOURCES

- *Associations between Immune Checkpoint Inhibitor Response, Immune-Related Adverse Events, and Steroid Use in Radiohead: A Prospective Pan-Tumor Cohort Study | Journal for Immunotherapy of Cancer*, jitc.bmj.com/content/13/5/e011545.
- *The Cancer-Immunity Cycle: Indication, Genotype, and Immunotype: Immunity*, www.cell.com/immunity/fulltext/S1074-7613(23)00416-8.
- Dyikanov D;Zaitsev A;Vasileva T;Wang I;Sokolov AA;Bolshakov ES;Frank A;Turova P;Golubeva O;Gantseva A;Kamysheva A;Shpudeiko P;Krauz I;Abdou M;Chasse M;Conroy T;Merriam NR;Alesse JE;English N;Shpak B;Shchetsova A;Tikhonov E;Filatov I;Radko A;Bolshakova A;K. "Comprehensive Peripheral Blood Immunoprofiling Reveals Five Immunotypes with Immunotherapy Response Characteristics in Patients with Cancer." *Cancer Cell*, U.S. National Library of Medicine, pubmed.ncbi.nlm.nih.gov/38744245/.
- Liang, Samantha I, et al. "Methylation-Based Ctdna Tumor Fraction Changes Predict Long-Term Clinical Benefit from Immune Checkpoint Inhibitors in Radiohead, a Real-World Pan-Cancer Study." *Cancer Research Communications*, U.S. National Library of Medicine, 1 Aug. 2025, pmc.ncbi.nlm.nih.gov/articles/PMC12365632/%E2%80%8B.
- Mellman, Ira, et al. "The cancer-immunity cycle: Indication, genotype, and immunotype." *Immunity*, vol. 56, no. 10, 2023, pp. 2188–2205, doi:10.1016/j.immuni.2023.09.011.
- Figure. Immune-related adverse events and impact on survival outcomes. Image from Zoe Quandt et al., Associations between immune checkpoint inhibitor response, immune-related adverse events, and steroid use in RADIOHEAD: a prospective pan-tumor cohort study, Journal for ImmunoTherapy of Cancer, vol. 13, no. 5, 12 May 2025, e011545, doi:10.1136/jitc-2025-011545. PubMed, https://pubmed.ncbi.nlm.nih.gov/40355283/
- Figure. Data center security layers, from One percent of Googlers get to visit a data center, but I did, Google Blog, 30 June 2020, blog.google/inside-google/infrastructure/how-data-center-security-works/
- Figure. Mellman, Ira, et al. "The cancer-immunity cycle: Indication, genotype, and immunotype." *Immunity*, vol. 56, no. 10, 2023, pp. 2188–2205, doi:10.1016/j.immuni.2023.09.011. Figure of the cancer-immunity cycle.
- Figure: Immune checkpoint inhibitors illustrating PD-1/PD-L1 interaction and blockade (Credit: © Terese Winslow, National Cancer Institute, cancer.gov). Originally published on Cancer.gov in "Immune Checkpoint Inhibitors." https://www.cancer.gov/about-cancer/treatment/types/immunotherapy/checkpoint-inhibitors.