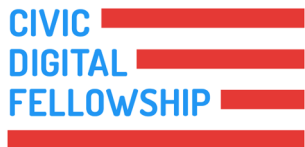


# AUTOMATING HIV/AIDS GRANT ENCODING

Office of AIDS Research

Robert Cregg, Lead Data Scientist  
Danny Murphy, Public Health Analyst



WILLIAM HUANG  
UCLA  
Electrical Engineering

YUYANG ZHONG  
UC Berkeley  
Psychology &  
Data Science

# OBJECTIVE

Automate and cross-verify the manual encoding of SIC codes, objective codes, and areas of emphasis for HIV/AIDS-related projects.

# PROBLEM STATEMENT

## Current Manual Coding:

- Grants receiving HIV/AIDS dollars – SIC codes
- Objective codes/areas of emphasis

## This project:

- Natural Language Processing (NLP) models with titles, abstracts and specific aims
- Facilitate accurate & efficient grant coding (SIC/Obj/AoE)

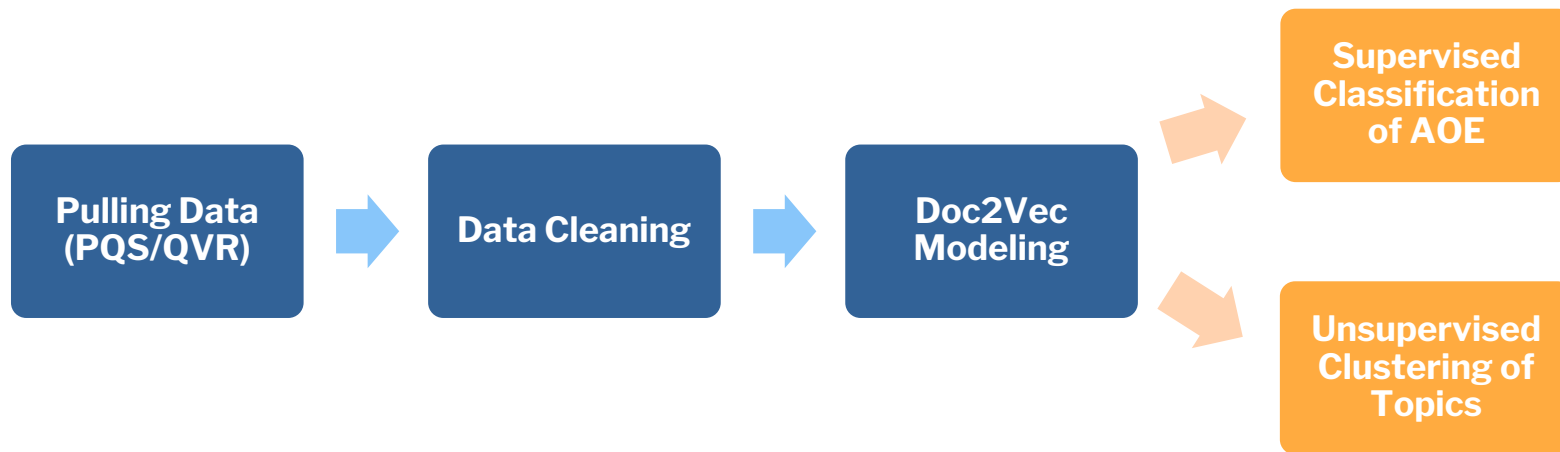
# PROCESS ITERATIONS

- Classification on SIC Codes
- Classification on Objective Codes
- Term Frequency Approach (TF-IDF)
- Deep Learning Implementation
- Classification on Area of Emphasis
- Unsupervised Clustering

# DATASET

- **6,884 Grants**
  - FY 2011-2020
  - PQS data merged with text data from QVR
  - Focus on Areas of Emphasis (AOEs)
  - Research Project Grants (RPG)

# METHODS



# BASIC TEXT CLEANING

## BEFORE CLEANING

Project summary/Abstract Accumulating evidence indicates that the incidence of HPV-associated neoplasia in HIV-positive individuals is substantially higher than in HIV-negative individuals despite effective antiretroviral therapy. These data strongly suggest that HIV may play a critical role in development of HPV-associated neoplasia of the anus, cervix and oropharyngeal cavity. However the mechanisms by which it does so are poorly understood. Our published work and preliminary data show that HIV may interact with oral and anal epithelia creating a tissue microenvironment where epithelial cells lose tight and adherence junctions.

## AFTER CLEANING

accumulating evidence indicates incidence hpv-associated neoplasia hiv-positive individual substantially higher hiv-negative individual despite effective antiretroviral therapy data strongly suggest hiv may play critical role development hpv-associated neoplasia anus cervix oropharyngeal cavity however mechanism poorly understood published work preliminary data show hiv may interact oral anal epithelium creating tissue microenvironment epithelial cell lose tight adherence junction

# RCDC CONCEPT REPLACEMENT

## BEFORE CLEANING

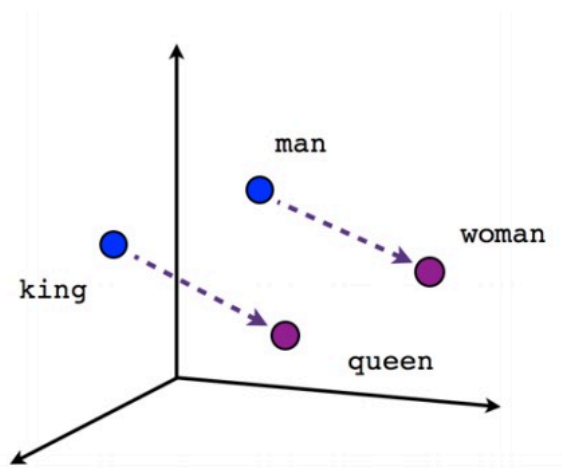
Project summary/Abstract Accumulating evidence indicates that the incidence of HPV-associated neoplasia in HIV-positive individuals is substantially higher than in HIV-negative individuals despite effective antiretroviral therapy. These data strongly suggest that HIV may play a critical role in development of HPV-associated neoplasia of the anus, cervix and oropharyngeal cavity. However the mechanisms by which it does so are poorly understood. Our published work and preliminary data show that HIV may interact with oral and anal epithelia creating a tissue microenvironment where epithelial cells lose tight and adherence junctions.

## AFTER CLEANING

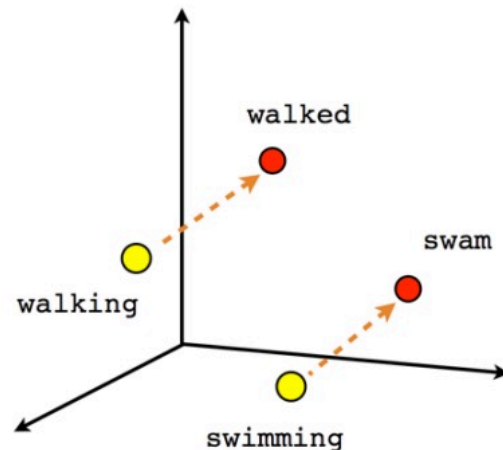
accumulating evidence concept\_incidence  
concept\_human\_papillomavirus-associated concept\_neoplasms  
concept\_hiv-positive individual concept\_hiv-negative individual  
effective concept\_antiretroviral\_therapy concept\_data  
concept\_hiv concept\_play critical concept\_role  
concept\_development concept\_human\_papillomavirus-  
associated concept\_neoplasms concept\_anus  
concept\_cervix\_uteri concept\_oropharyngeal cavity mechanism  
understood published concept\_work preliminary concept\_data  
concept\_hiv interact concept\_oral concept\_anus epithelium  
creating tissue microenvironment concept\_epithelial\_cells lose  
tight concept\_adherence junction



# WORD2VEC OVERVIEW



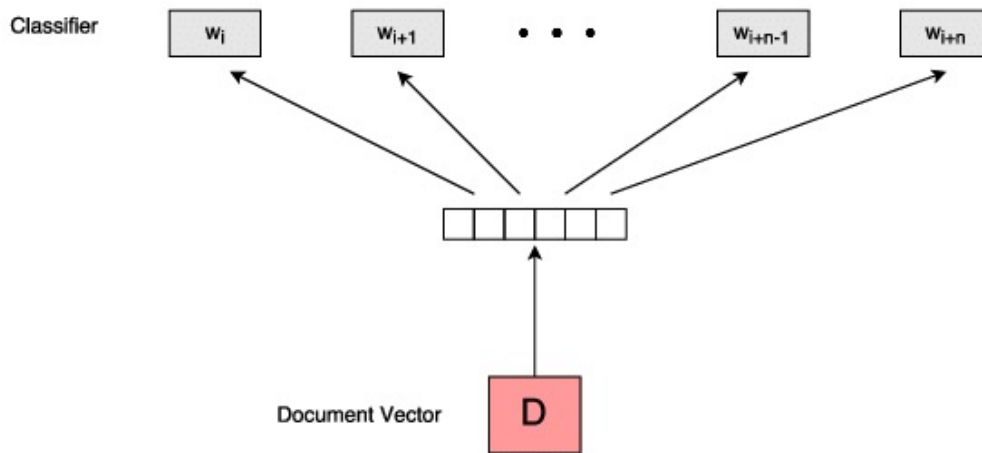
Male-Female



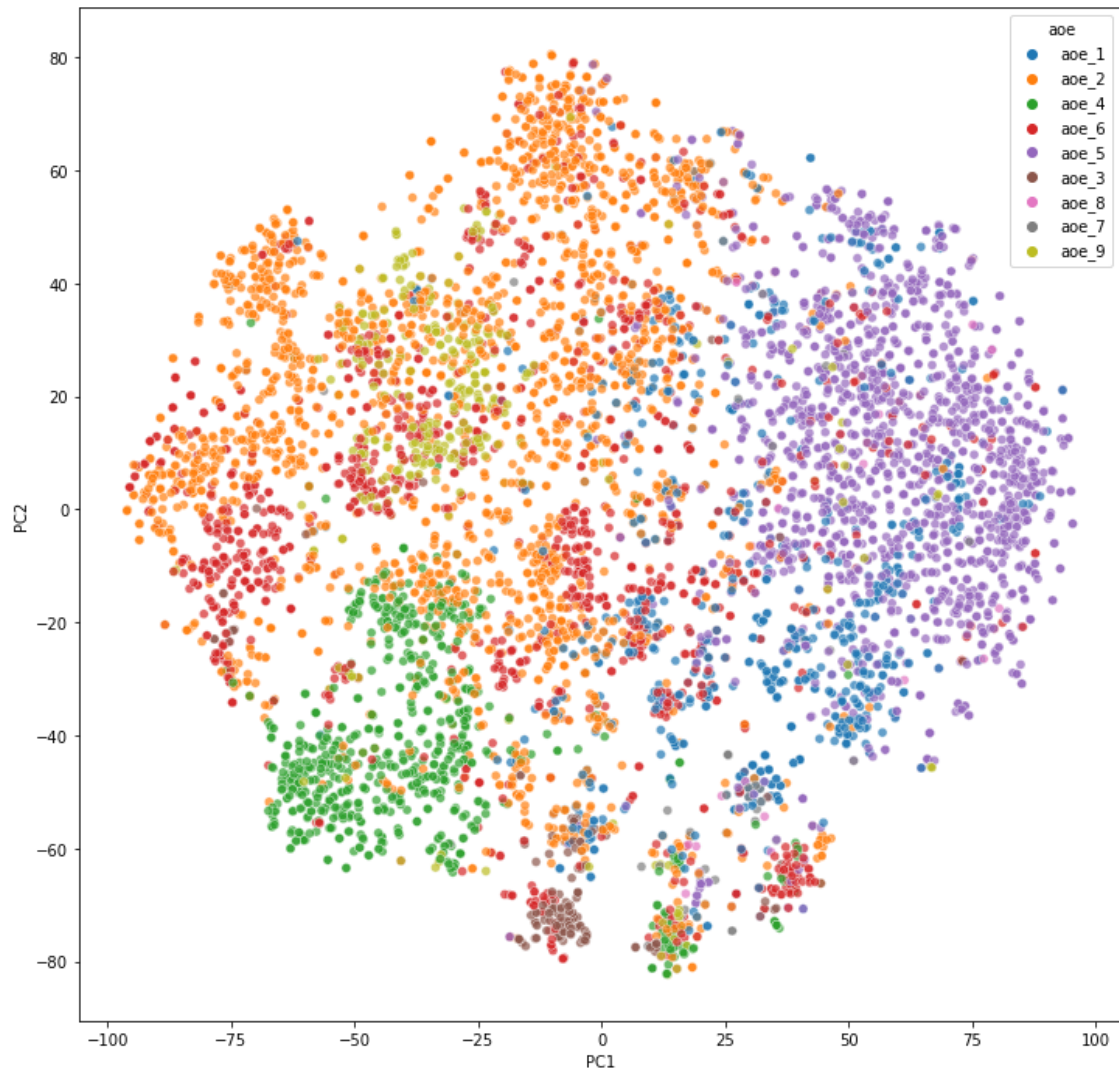
Verb tense

# DOC2VEC OVERVIEW

- Distributed Bag of Words (DBOW)
- Use the center word to predict words around it -> context-based



# Visualizing Doc2Vec via t-SNE

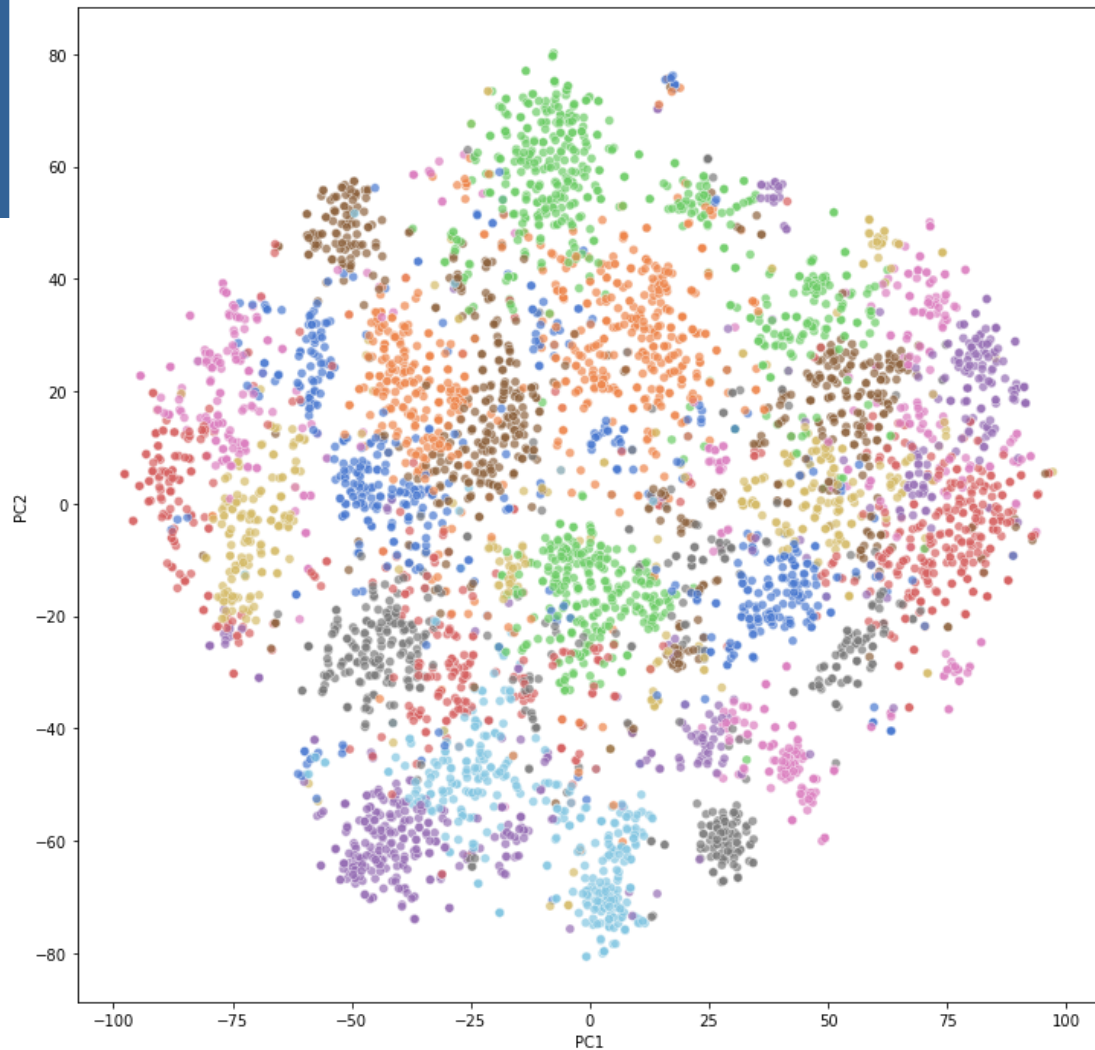


# SUPERVISED CLASSIFICATION

Model	Accuracy
Logistic Regression	72.33%
Neural Network	74.22%
Random Forest	70.15%
Gradient Boosting	57.66%
Voting Classifier (Ensemble)	75.74%

- Partitioned 80% Training & 20% Testing
- The max allocated AOE is used for prediction
  - Check for multiple AOE situation: If predicted AOE is one of the lesser allocated AOE's
- Final Test Accuracy after Correction: **79.59%**

# Unsupervised Clustering By KMeans



# UNSUPERVISED CLUSTERING CONT.

## Select Cluster Topics

Social Stigma, AIDS Prevention, Young Men who have Sex with Men, Counseling, Black Men who have Sex with Men

CD4, Effector, Control, Immunity, GC, CD8B1 Gene, Signal Transduction, B Lymphocyte

Alcohol or other Drugs Use, Imprisonment, Risk Behaviors, User, Services, offender, Community, Injecting Drug User, Drug

- Computationally generate topic labels for different KMeans clusters
  - Select high frequency words that do not appear frequently in other clusters
- Intended to be human readable and descriptive

# TAKEAWAY & NEXT STEPS

- **Automated Grant Categorization**

- Automated method to help in the assignment of evaluators to grant proposals

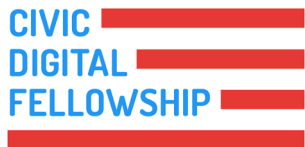
- **Cross-Verification**

- Check existing categorizations and compare them to what a computer might generate
- Use D2V's clustering ability to find closely linked grants

# RCDC CATEGORY OVERLAP VISUALIZATION VIA TABLEAU

Office of AIDS Research

Robert Cregg, Lead Data Scientist  
Danny Murphy, Public Health Analyst



National Institutes of Health  
*Office of AIDS Research*

WILLIAM HUANG  
UCLA  
Electrical Engineering

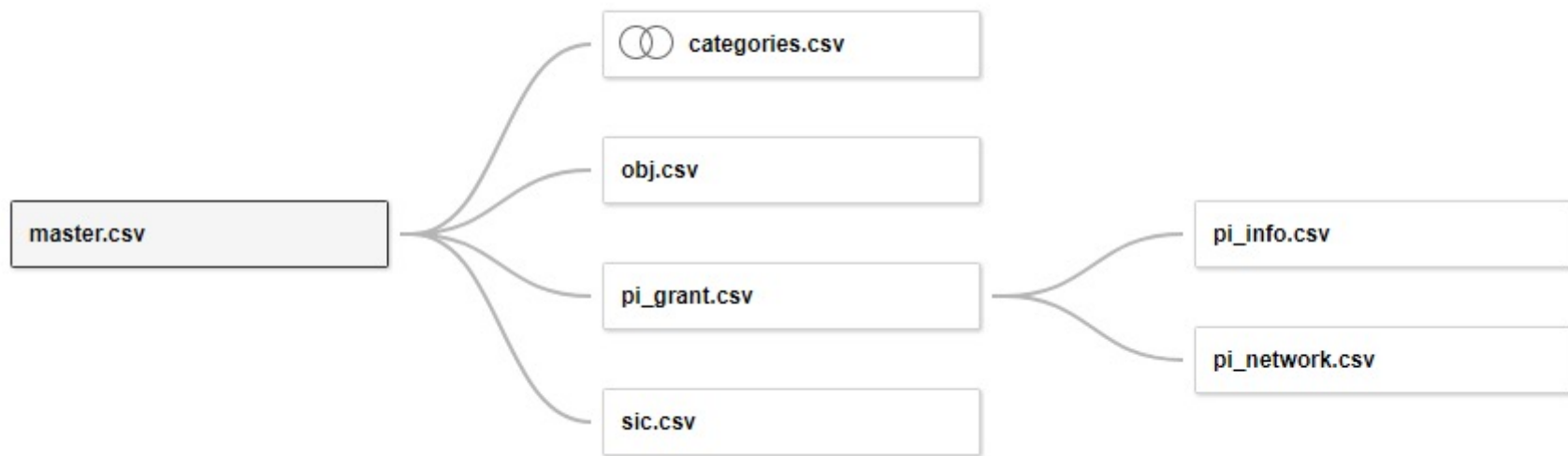
YUYANG ZHONG  
UC Berkeley  
Psychology &  
Data Science



# OBJECTIVE

Visualize funding overlap and other relationships between different major research topics related to HIV/AIDS.

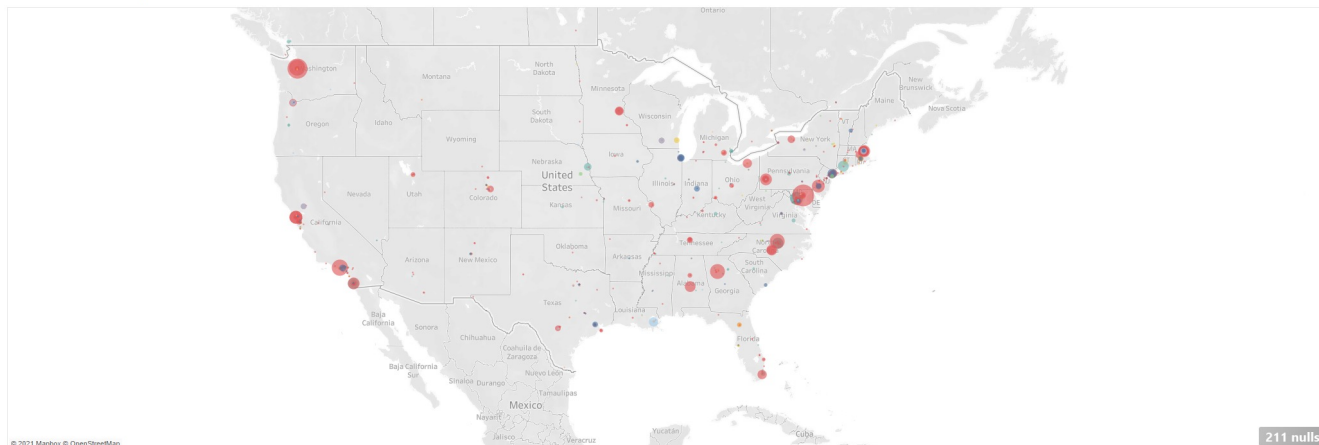
# DATABASE SETUP



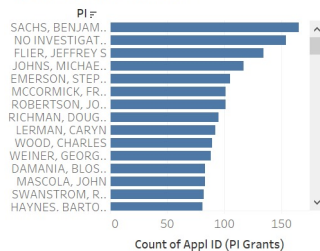
Current database organization from Tableau.

# INSTITUTION VIEW

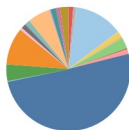
Institution Map



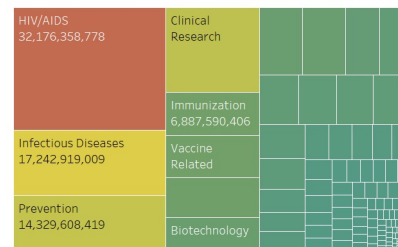
Total Grants Written



IC Funding Distribution

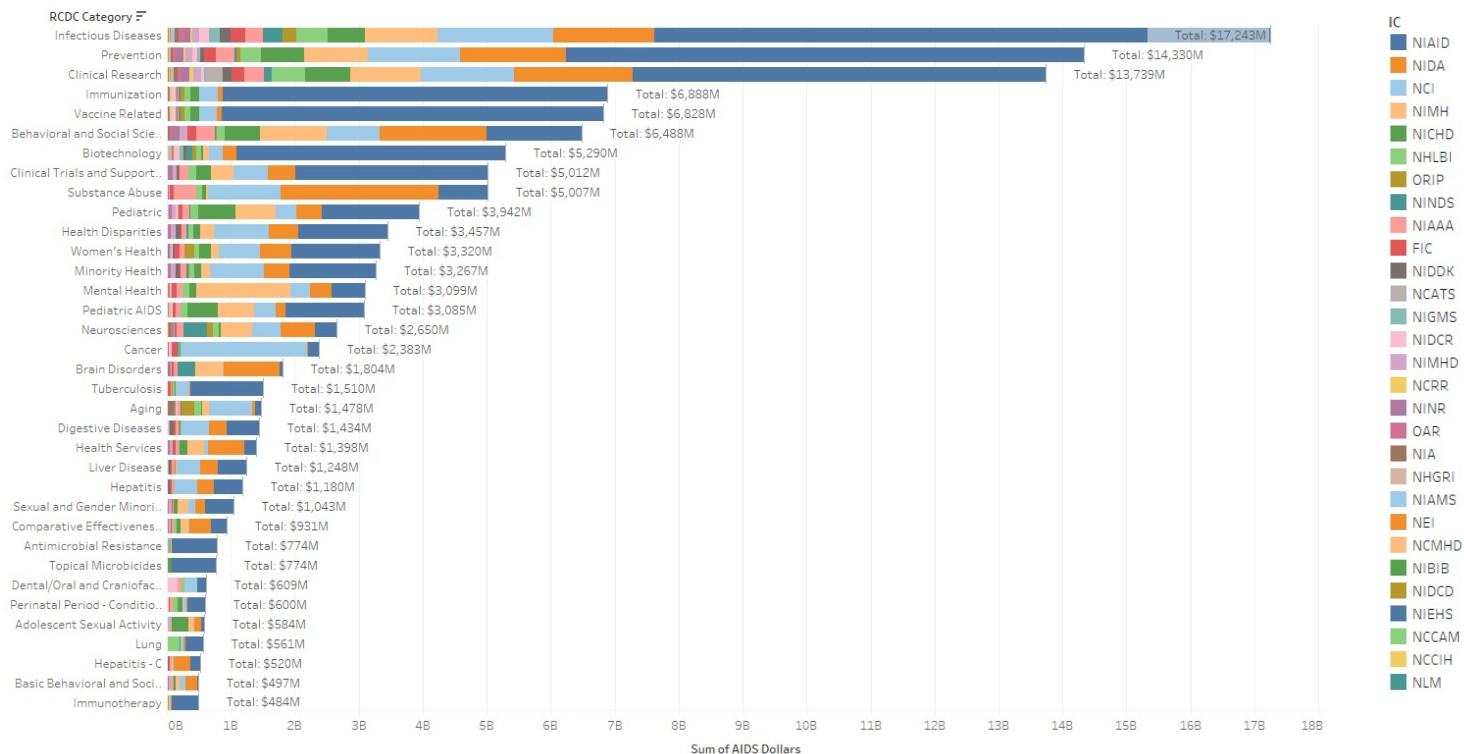


RCDC Funding Distribution (Filter by INST)



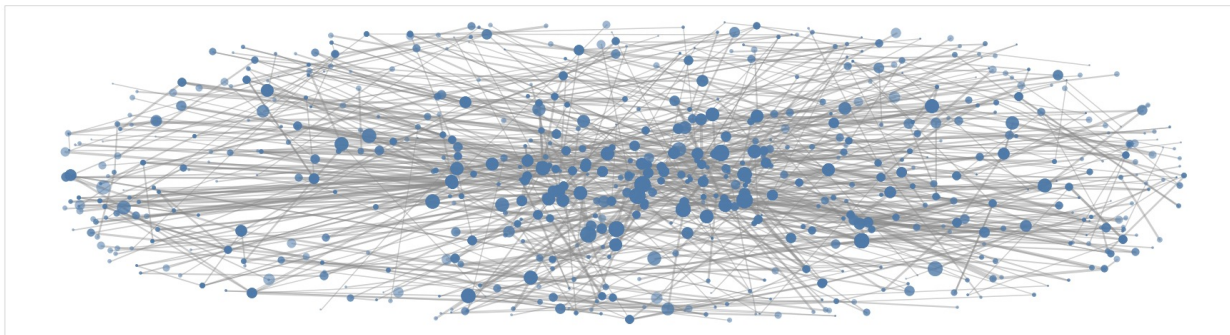
# RCDC FUNDING DISTRIBUTION BY IC

RCDC Funding Distribution (ALL IC) for HIV Projects FY 2009-2020



# PI COLLABORATION NETWORK

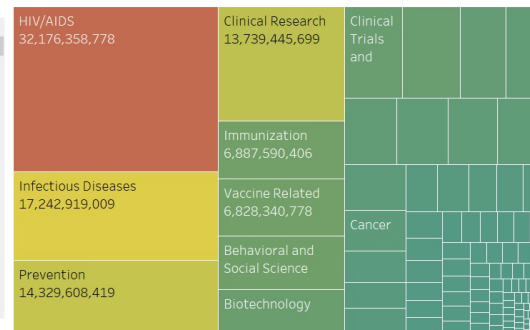
PI Network



PI Grants

Appl ID	Fiscal Y.	Project Title
74911..	2009	A MULTILEVEL HIV-PREVENTION STRATEGY FOR HIGH-RISK YOUTH
74932..	2009	DRUG ABUSE, SUSTANCE P AND HIV
74953..	2009	PLGA/ANTISENSE IL-10: GENE THERAPY FOR COCAINE ABUSERS WITH HAD
74962..	2009	ANALYSIS OF RESTRICTION FACTOR FUNCTION IN SITU
75327..	2009	IMMUNE RESPONSE TO PNEUMOCYSTIS IN SIMIAN MODEL OF AIDS
75334..	2009	MONOCLONAL ANTIBODIES FOR THE STUDY OF P. CARINII
75343..	2009	HIV/STD RISK REDUCTION FOR AFRICAN AMERICAN COUPLES
75343..	2009	ROLE OF ESTROGEN IN CERVICAL CANCER
75347..	2009	INVESTIGATION OF APOBEC3G/VIF INTERACTIONS BY SAXS AND X-RAY CRYSTAL
75348..	2009	MOLECULAR MECHANISMS OF CELLULAR GENE EXPRESSION REPROGRAMMING
75350..	2009	ADHERENCE TO HIV CARE: ADVANCING THE SCIENCE
75351..	2009	STUDIES OF NATURAL SIV INFECTION OF SOOTY MANGABEYS
75352..	2009	HIV-1 INTEGRASE STRUCTURE AND FUNCTION AS A THERAPEUTIC TARGET
75352..	2009	STRUCTURAL STUDIES OF HIV-1 INTEGRATION
75372..	2009	STRUCTURE AND FUNCTION OF THE HEPATITIS C VIRUS GENOME
75383..	2009	NATURAL HISTORY OF HIV INFECTION IN INJECTION DRUG USERS
75399..	2009	INTERFACE OF INNATE AND ADAPTIVE IMMUNITY IN HIV-1 INFECTION
75399..	2009	INTERVENTION FOR HIV+ ADULTS WITH CHILDHOOD SEXUAL ABUSE
75399..	2009	CONDITIONAL CASH TRANSFERS TO PREVENT SEXUALLY TRANSMITTED INFECT
75404..	2009	KILLER CELLS & VIRAL LOAD IN VERTICAL HIV INFECTION
75409..	2009	MUCOSAL REOVIRUS-ADENOVIRUS VACCINES AGAINST HIV-1

RCDC Funding Distribution (Filter by PI)



# TAKEAWAY & NEXT STEPS

- **Easily Extendable**

- Standardized database setup allows future data scientists to build more workbooks and visualizations

- **Accessible**

- Leverage Tableau's different features to develop easily readable visualizations that can be shown to internal OAR members, Congress, or the general public.