Year 3 All Hands Meeting
April 20, 2023

# Project Overview

9:30 – 10:30 AM

Jennifer Fowler, PI
Arkansas Economic Development Commission

Jackson Cothren, Co-PI
University of Arkansas, Fayetteville

# The motivation for DART

# Vision

The Arkansas research community - academic, government, and industry - collaborating often and easily on a shared computing platform through a modern cyberinfrastructure which facilitates cutting-edge data science research.

# Mission

The mission of DART is to improve research capability and competitiveness in Arkansas by creating an integrated statewide consortium of researchers and educators working to establish a synergistic, statewide focus on excellence in data analytics research and training.

# The continued (growing) relevance of DART

# How DART is organized



- + Addressing barriers to effective and accepted integration of data science into our lives.

- + Building a shared research cyberinfrastructure in Arkansas

- + Creating a state-wide, latticed data science education program

- + Fostering a collaborative research community in Arkansas

# Project Management

# Significant Challenges



Enhanced cybersecurity footprint
- 🔒 Increased oversight
- 🔒 Off-campus access to research computing more controlled
- 🔒 ScienceDMZ architecture reconfigured
  - 🔒 15-month planning process
  - 🔒 Limited access to cluster and storage
  - 🔒 Limited access to Enterprise GitLab

💡 **Solution: Build enduring relationships between campus IT and research communities**

# Significant Challenges



Research themes identified a need for additional programming support

- several projects need to scale
- deploying code on ARP clusters and commercial cloud
- 1 FTE currently supporting this need, more talent is needed
- Keeping talent, competing with industry, is hard

**Solution: Reallocate program funds to support diverse programming requirements.**

# Significant Accomplishments – Grants & Proposals



**$ 24,389,199.41 Awarded to Date**

# Significant Accomplishments – Grants & Proposals

| Project Team(s) | Amount Awarded |
|---|---|
| Cyberinfrastructure (CI) | $ 884,318 |
| Data Curation (DC) | $ 7,176,848 |
| Learning & Prediction (LP) | $ 478,763 |
| Learning & Prediction (LP), Cyberinfrastructure (CI) | $ 895,000 |
| Learning & Prediction (LP), Data Curation (DC) | $ 155,199 |
| Learning & Prediction (LP), Data Curation (DC), Social Awareness (SA) | $ 1,450,003 |
| Social Awareness (SA) | $ 1,463,789 |
| Social Awareness (SA), Workforce Development (WD) | $ 49,999 |
| Social Media (SM) | $ 11,835,280 |

# Significant Accomplishments – Grants > $500,000

| Title | Investigators | Institutions | Award Amount | Funding Agency |
|-------|---------------|--------------|--------------|----------------|
| Multi-Level Models of Covert Online Information Campaigns | Nitin Agarwal | UALR | $ 4,965,214 | Department of Defense |
| Developing Rapid Response Capabilities to Evaluate Emerging Social Cyber Threats | Nitin Agarwal | UALR | $ 3,773,526 | U.S. Office of Naval Research |
| Fusing Narrative and Social Cyber Forensics to Understand Covert Influence | Nitin Agarwal | UALR | $ 2,500,000 | Department of Defense |
| Targeting heat shock protein 72 to improve renal function after transplantation | Se-Ran Jun | UAMS | $ 2,461,707 | NIDDK |
| Center for studies of host response to cancer therapy | Se-Ran Jun | UAMS | $ 2,280,000 | NIGMS |
| RII Track-2 FEC: Artificial Intelligence on Sustainable Energy Infrastructure Network (AI SUSTEIN) and Beyond towards Industries of the Future | Haitao Liao, Xintao Wu, Xiao Liu | UARK | $ 1,450,003 | NSF |
| Assessment of antibiotic resistance in fresh vegetables from farm to fork | Se-Ran Jun | UAMS | $ 1,000,000 | USDA |
| Photogrammetry Services, Task Order | Jackson Cothren, Chase Rainwater | UARK | $ 800,000 | Department of Energy |
| FAI: A novel paradigm for fairness-aware deep learning models on data streams | Xintao Wu | UARK | $ 628,789 | NSF |
| IUCRC Phase I The University of Arkansas: Center for Infrastructure Trustworthiness in Energy Systems (CITES) | Qinghua Li | UARK | $ 525,000 | NSF |

# Significant Accomplishments - Publications

| Publication Status | Project Total |
|---|---|
| **Accepted** | **42** |
| **Published** | **140** |
| **Submitted** | **25** |

| Project Team | Total Published | Journal Articles | Books | Conference Proceedings |
|---|---|---|---|---|
| **Cyberinfrastructure (CI)** | **2** | **2** | **0** | **0** |
| **Cyberinfrastructure (CI), Data Curation (DC)** | **2** | **1** | **1** | **0** |
| **Data Curation (DC)** | **29** | **19** | **2** | **6** |
| **Education (ED)** | **1** | **0** | **0** | **0** |
| **Learning & Prediction (LP)** | **37** | **18** | **1** | **18** |
| **Learning & Prediction (LP), Cyberinfrastructure (CI)** | **1** | **1** | **0** | **0** |
| **Social Awareness (SA)** | **30** | **5** | **5** | **25** |
| **Social Media (SM)** | **23** | **11** | **0** | **11** |

**DART**

Home | About ▾ | Calendar | Events ▾ | Resources ▾ | FAQ

Subscribe | Contact

Home / About DART /

# Project Publications

Most recent publication is listed first; browse, filter, and review abstracts on Zotero.

1. Yang X, Zhang N, Schrader P (2022) A study of brain networks for autism spectrum disorder classification using resting-state functional connectivity. Machine Learning with Applications 8:100290. https://doi.org/10.1016/j.mlwa.2022.100290
2. Tran M, Ly L, Hua B-S, Le N (2022) SS-3DCapsNet: Self-supervised 3D Capsule Networks for Medical Segmentation on Less Labeled Data. arXiv:220105905 [cs, eess]
3. Vo-Ho V-K, Yamazaki K, Hoang H, Tran M-T, Le N (2022) Meta-Learning of NAS for Few-shot Learning in Medical Image Applications. arXiv:220308951 [cs, eess]
4. Tran M, Vo-Ho V-K, Quinn K, Nguyen H, Luu K, Le N (2022) CapsNet for Medical Image Segmentation. arXiv:220308948 [cs, eess]
5. Spann B, Mead E, Maleki M, Agarwal N, Williams T (2022) Applying diffusion of innovations theory to social networks to understand the stages of adoption in connective action campaigns. Online Social Networks and Media 28:100201. https://doi.org/10.1016/j.osnem.2022.100201
6. Bellis ES, von Münchow CS, Odero CO, Kronberger A, Kelly E, Xia T, Huang X, Wicke S, Runo SM, dePamphilis CW, Lasky JR (2022) Genomic signatures of host-specific selection in a parasitic plant. Evolutionary Biology
7. Moon SH, Udaondo Z, Abram KZ, Li X, Yang X, DiCaprio EL, Jun S-R, Huang E (2022) Isolation of AmpC- and extended spectrum β-lactamase-producing Enterobacterales from fresh vegetables in the United States. Food Control 132:108559. https://doi.org/10.1016/j.foodcont.2021.108559
8. Hu Y, Zhang L (2022) Achieving Long-Term Fairness in Sequential Decision Making. In: Proceedings of the First MiniCon Conference
9. Wassenaar TM, Wanchai V, Buzard G, Ussery DW (2022) The first three waves of the Covid-19 pandemic hint at a limited genetic repertoire for SARS-CoV-2. FEMS Microbiology Reviews fuac003. https://doi.org/10.1093/femsre/fuac003
10. McNerney HW, Spann B, Mead EL, Kready J, Marcoux T, Agarwal N (2022) Assessing the influence and reach of digital activity amongst far-right actors: A comparative evaluation of mainstream and 'free speech' social media platforms. For(e)Dialogue. https://doi.org/10.21428/e3990ae6.60c47409
11. Maleki M, Arani M, Mead E, Kready J, Agarwal N (2022) Applying an Epidemiological Model to Evaluate the Propagation of Toxicity related to COVID-19 on Twitter
12. Tanu SS, Zhang L, Gauri D, Sha Z (2022) An Exploratory Study on Fairness-Aware Design Decision-Making
13. Daneshjou R, Barata C, Betz-Stablein B, Celebi ME, Codella N, Combalia M, Guitera P, Gutman D, Halpern A, Helba B, Kittler H, Kose K, Liopyris K, Malvehy J, Seog HS, Soyer HP, Tkaczyk ER, Tschandl P, Rotemberg V (2022) Checklist for Evaluation of Image-Based Artificial Intelligence Reports in Dermatology: CLEAR Derm Consensus Guidelines From the International Skin Imaging Collaboration Artificial Intelligence Working Group. JAMA Dermatol 158:90. https://doi.org/10.1001/jamadermatol.2021.4915
14. Quach KG, Le N, Duong CN, Jalata I, Roy K, Luu K (2022) Non-volume preserving-based fusion to group-level emotion recognition on crowd videos. Pattern Recognition 128:108646. https://doi.org/10.1016/j.patcog.2022.108646
15. Khaund T, Kirdemir B, Agarwal N, Liu H, Morstatter F (2022) Social Bots and Their Coordination During Online Campaigns: A Survey. IEEE Trans Comput Soc Syst 9:530–545. https://doi.org/10.1109/TCSS.2021.3103515
16. Le N, Rathour VS, Yamazaki K, Luu K, Savvides M (2022) Deep reinforcement learning in computer vision: a comprehensive survey. Artif Intell Rev 55:2733–2819. https://doi.org/10.1007/s10462-021-10061-9
17. Alassad M, Hussain MN, Agarwal N (2022) Comprehensive decomposition optimization method for locating key sets of commenters spreading conspiracy theory in complex social networks. Cent Eur J Oper Res 30:367–394. https://doi.org/10.1007/s10100-021-00738-5
18. Mead E, Agarwal N (2022) Surface Web vs Deep Web vs Dark Web. In: Schintler LA, McNeely CL (eds) Encyclopedia of Big Data. Springer International Publishing, Cham, pp 901–905
19. Nasiri E, Milanova M, Nasiri A (2022) Masked Face Detection Using Artificial Intelligent Techniques. In: Kountchev R, Mironov R, Nakamatsu K (eds) New Approaches for Multidimensional Signal Processing. Springer, Singapore, pp 3–34
20. Pierce E (2022) A Balanced Scorecard for Maximizing Data Performance. Frontiers in Big Data 5:
21. Hu C, Sheng VS, Wu N, Wu X (2022) Managing Uncertainty in Crowdsourcing with Interval-Valued Labels. In: Rayz J, Raskin V, Dick S, Kreinovich V (eds) Explainable AI and Other Applications of Fuzzy Techniques. Springer International Publishing, Cham, pp 166–178
22. Jiang R, Chazot P, Pavese N, Crookes D, Bouridane A, Celebi ME (2022) Private Facial Prediagnosis as an Edge Service for Parkinson's

---

## zotero

Other Group Libraries
- dart publications

| Title | Creator | Date | |
|---|---|---|---|
| Narrative trends of COVID-19 misinformation | Marcoux and Agarwal | 2021 | |
| Narrow Band Active Contour Attention Model for Medic... | Le et al. | 2021-07-31 | |
| Non-volume preserving-based fusion to group-level e... | Quach et al. | 2022 | |
| Offset Curves Loss for Imbalanced Problem in Medical ... | Le et al. | 2021-01-10 | |
| PairFlow: Enhancing Portable Chest X-Ray By Flow-Bas... | Le et al. | 2021-09-19 | |
| Performance Analysis of Tor Website Fingerprinting ov... | Oh et al. | 2020-12 | |
| Plasma Metabolomics in a Nonhuman Primate Model of... | Jun et al. | 2021-08-13 | |
| Point-Unet: A Context-Aware Point-Based Neural Netw... | Ho et al. | 2021 | |
| Poisoning Attacks on Fair Machine Learning | Van et al. | 2021-10-17 | |
| Predicting Stance Polarity and Intensity in Cyber Argu... | Sirrianni et al. | 2021 | |
| Private Facial Prediagnosis as an Edge Service for Park... | Jiang et al. | 2022 | |
| ProdMX: Rapid query and analysis of protein functional... | Wanchai et al. | 2020 | |
| Proposing a Broader Scope of Predictive Features for ... | Mead et al. | 2021-04-19 | |
| PTMViz: a tool for analyzing and visualizing histone po... | Chappell et al. | 2021 | |
| Quantitative Modeling and Analysis of Argumentation P... | Sirrianni et al. | 2021 | |
| Removing Disparate Impact on Model Accuracy in Diffe... | Xu et al. | 2021-08-14 | |
| Roughness Index and Roughness Distance for Benchm... | Rathour et al. | 2021 | |
| Skin Melanoma Detection in Microscopic Images Using... | Rastghalam et al. | 2021 | |
| Social Bots and Their Coordination During Online Cam... | Khaund et al. | 2022 | |
| SS-3DCapsNet: Self-supervised 3D Capsule Networks ... | Tran et al. | 2022-03-28 | |
| Studying the Role of Social Bots During Cyber Flash M... | Al-khateeb et al. | 2021 | |

Info | Notes | Tags | Attachments | Related | Show Empty Fields

| | |
|---|---|
| Item Type | Journal Article |
| Title | PTMViz: a tool for analyzing and visualizing histone post translational modification data |
| Author | Chappell, Kevin |
| Author | Graw, Stefan |
| Author | Washam, Charity L. |
| Author | Storey, Aaron J. |
| Author | Bolden, Chris |
| Author | Peterson, Eric C. |
| Author | Byrum, Stephanie D. |
| Publication | BMC Bioinformatics |
| Volume | 22 |
| Issue | 1 |
| Pages | 275 |
| Date | 12/2021 |
| Journal Abbr | BMC Bioinformatics |
| Language | en |
| DOI | 10.1186/s12859-021-04166-9 |
| ISSN | 1471-2105 |
| Short Title | PTMViz |
| URL | https://bmcbioinformatics.biomedcentral.co... |
| Accessed | 4/1/2022, 4:43:43 PM |
| Library Catalog | DOI.org (Crossref) |

Tags: AmpC, Analytical models, ASD, Black Box, blogs, Blogs, BlogTracker, bots, Brain networks, causal discovery, Causality, Collective action, Complexity Theory, Computer Science - Artificial Intelligence, Computer Science - Computation and Langu..., Computer Science - Computer Vision and Pat..., Computer Science - Computers and Society, Computer Science - Cryptography and Securi..., Computer Science - Databases, Computer Science - Machine Learning, Computer Science - Social and Information N..., Computer Science - Sound, Computer Vision and Pattern Recognition (cs...

Filter Tags

**Abstract**

Abstract

Background
Histone post-translational modifications (PTMs) play an important role in our system by regulating the structure of chromatin and therefore contribute to the regulation of gene and protein expression. Irregularities in histone PTMs can lead to a variety of different diseases including various forms of cancer. Histone modifications are analyzed using high resolution mass spectrometry, which generate large amounts of data that requires sophisticated bioinformatics tools for analysis and visualization. PTMViz is designed for downstream differential abundance analysis and visualization of both protein and/or histone modifications.

Results
PTMViz provides users with data tables and visualization plots of significantly differentiated proteins and histone PTMs between two sample groups. All the data is packaged into interactive data tables and graphs using the Shiny platform to help the user explore the data in a fast and efficient manner to assess if changes in the system are due to protein abundance changes or epigenetic changes. In the example data provided, we identified several proteins differentially regulated in the dopaminergic pathway between mice treated with methamphetamine compared to a saline control. We also identified histone post-translational modifications including histone
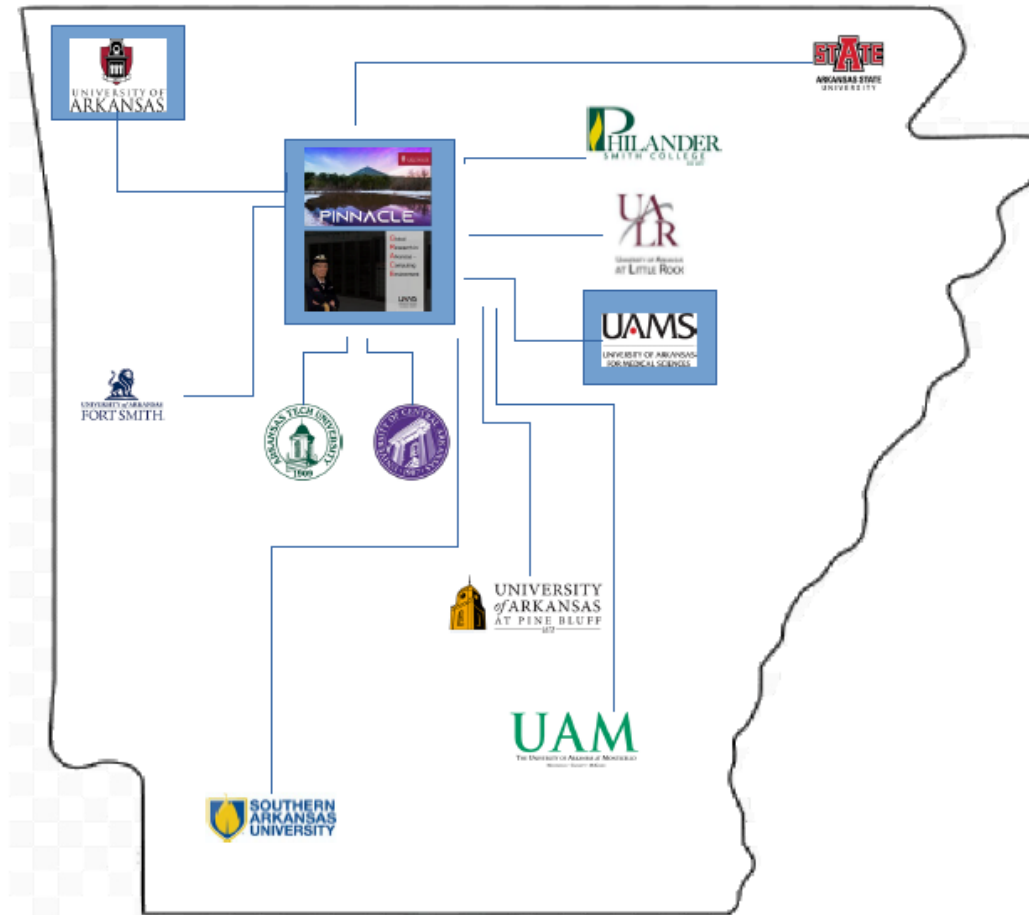
# Arkansas Research Platform: Computational Resources
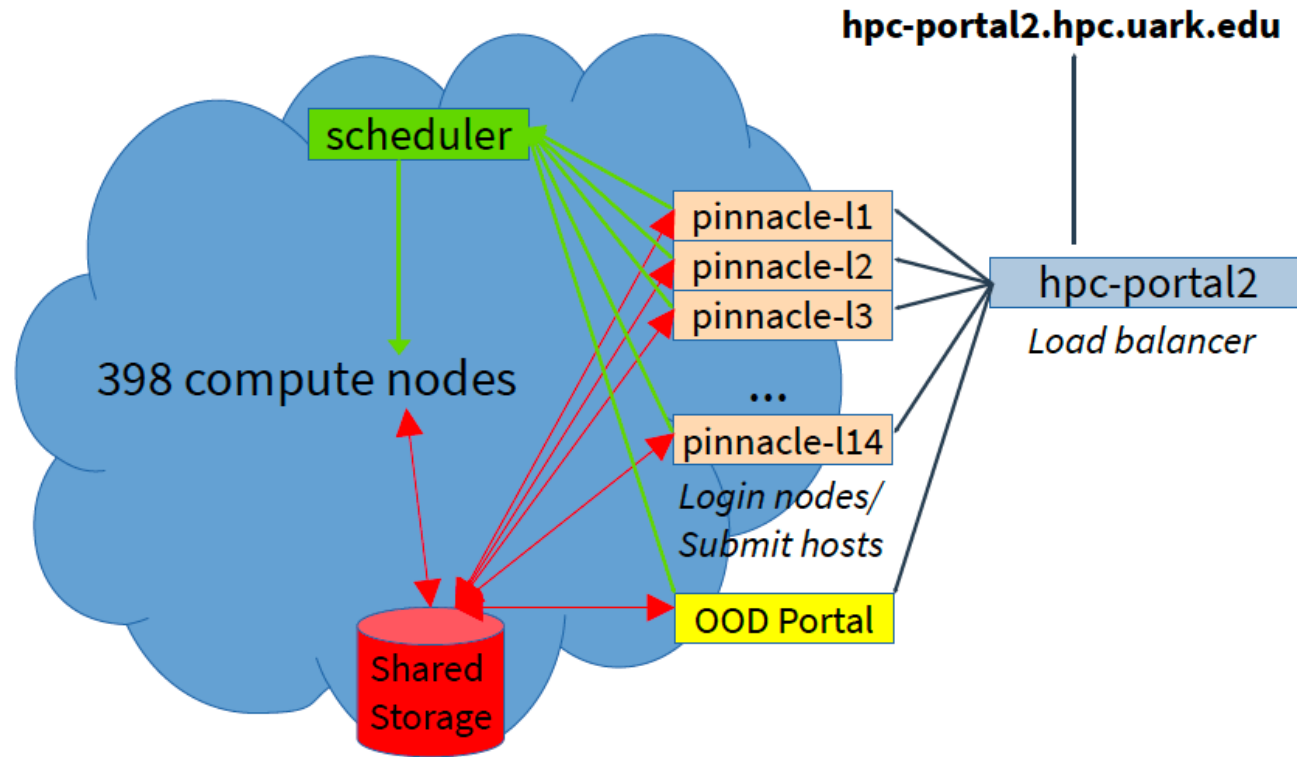
# Arkansas Research Platform



Arkansas Research Platform is a collaboration of higher education institutions within the state of Arkansas which provides computing resources and data storage to all its members (free of charge to students, faculty and staff). The core components are the **Pinnacle cluster** managed by the Arkansas High Performance Computing Center (UAF, Fayetteville) and the **Grace cluster** managed by UAMS High Performance Computing Center (UAMS, Little Rock).

supported by UAMS High Performance Computing  AHPCC Arkansas High Performance Computing Center

ARKANSAS RESEARCH PLATFORM ARP

# ARP Clusters – functional diagrams



## Pinnacle

hpc-portal2.hpc.uark.edu

scheduler

pinnacle-l1
pinnacle-l2
pinnacle-l3
...
pinnacle-l14

*Login nodes/ Submit hosts*

398 compute nodes

hpc-portal2
*Load balancer*

OOD Portal

Shared Storage

## Grace

login.hpc.uams.edu

scheduler
scheduler

login1

*Login nodes/ Submit hosts*

196 compute nodes

OOD Portal

Shared Storage

portal.hpc.uams.edu

*scheduler* – coordinates all of the nodes in the queue (SLURM queueing system)
*submit hosts* – nodes from which users submit jobs to the queues
*login hosts* – nodes accessible from the external network
*compute nodes* – nodes that assigned to queues and run jobs

# Arkansas Research Platform - ARP

Deployed a point-to-point virtual research network between UAMS and UAF, extensible to other ARE-ON members
- Forms the backbone of the ARP

Installed 100Gbps connection between North Little Rock and Fayetteville ARE-ON backbone locations to enable higher throughput between UAMS and UAF.

Installed dedicated research 100Gbps connection to the UAMS campus for the UAMS Science DMZ

# ARP Workshops

Began in Fall 2022

- UAMS
- UALR

Offered again Fall 2023



## Getting Started with the
## Arkansas Research Platform

**UAMS Campus**
**September 29-30**

*Day 1*

supported by **UAMS** High Performance Computing **AHPCC** Arkansas High Performance Computing Center

ARKANSAS RESEARCH PLATFORM **ARP**

# Day 1 Schedule

**8:30 am**    **Registration & Set-Up**
        Check with each participant to ensure ability to connect to internet
        Make sure all required programs have been downloaded and installed

**9:00 am**    **Welcome & Introductions**
        Brief introductions: Name, institution, 1-sentence research interests

**9:15 am**    **What is Arkansas Research Platform (ARP)?**
        Brief description of ARP hardware resources
        When to use HPC resources
        Programs that can leverage HPC resources
        Getting familiar with the vocabulary
        Cluster Functional diagrams

**10:15 am**    **The Essentials**
        Getting an account on Pinnacle / Grace
        Logging into Open OnDemand portal
        OOD Overview

**10:45 am**    **Running Jobs – Open Ondemand**
        Batch job (run, modify, rerun)
        Interactive jobs (Jupyter notebook, VMD)

**11:45 am**    **LUNCH**

**12:45 pm**    **Moving Data**
        Open Ondemand File Upload/Download
        Terminal access in Open OnDemand
        wget, scp, rsync, FileZilla (GUI)

**2:00 pm**    Globus
        Log in/create Globus ID account
        Globus endpoints
        Transfer data between Pinnacle and Public endpoint
        Install Globus Connect Personal
        Transfer data between Pinnacle and personal device

**4:00 pm**    **Reflection and Closing Remarks** for Day 1
        Where are you stuck?

supported by UAMS High Performance Computing · AHPCC Arkansas High Performance Computing Center

ARKANSAS RESEARCH PLATFORM · ARP

## Day 2 Schedule

**9:00 am**  **SSH access**
SSH clients (Windows/MacOS/Linux)
Logging into Pinnacle and Grace
Network Accessibility
Pinnacle's 2-step login on the DMZ network

**9:30 am**  **Intro to Command Line Interface (CLI)**
Bash shell
Environment
Basic linux commands
File systems
/home and /scratch

**10:00 am**  **Queueing System**
Queues
Running a sample job
Interactive job
Batch job
Job arrays

**10:45 am**  **Software modules**
Modules

**11:00 am**  **Moving and Storing Data, part 2**
Basic object storage principles; What is ROSS?
Request access to ROSS – identities and secret keys
ecs-sync

**12:00 pm**  **LUNCH**

**1:00 pm**  **Individual project help**
Questions, Tell me more, and Bring your own project (optional for attendees)

# And, in the coming months…



# …federated identity access

# Cyberinfrastructure Plan

Leveraged two NSF CC* funded projects:

GPN CyberTeam (NSF #1925681):

- CI working group formed, led by Great Plains Network executive director and CyberTeam PI's

UAF SHARP CCI (2021) (NSF #2126108):

- ARP working group composed of all DART campus CIO's and research leads to plan for sustainable ARP

- State-wide gap analysis underway

- Reallocated funds from Globus Data Management service to fund federated identify solutions

- Developed MISP and SSP templates for managing CUI

- HIPAA-compliant storage at UAMS

- Engaged with TrustedCI and EPOC teams at Indiana University

# Significant Accomplishments:
CI, LP, DC, SM, SA

# Significant Accomplishments



- Significant progress implementing the Arkansas Research Platform & Computing Collaborative for scaling algorithms and applications

- Integrated research and IT staff statewide working groups consolidating resources and expanding access

- Positive data control concepts and implementations being tested by an industrial partner

- Developed an autoencoder method that improved unsupervised and self-supervised deep learning methods' ability to manage and label much larger datasets.

- Much more sophisticated, coordinated approach to managing controlled or restricted data

# Auto-annotation of multimedia data



**Developing novel methods and algorithms**

- for image and video analysis
  - Multimedia data characteristics towards target applications have been identified.
  - Applications have been defined around the team's research into better informing disaster response with social media (SM4).

- vehicle speed estimation which has tested for road condition assessment

- for rain drop removal using GANs network to improve accuracy in using visual data captured in rainy conditions

Dr. Serhan Dagtas, UALR,

GA: Sharafat Hossain, UALR

# Parameter Discovery Process (PDP) (DC1)

- Make the **Data Washing Machine (DWM)** truly "unsupervised"

- Automatically sets 14 DWM input parameters

- Given a dataset with redundant (duplicate) records, the PDP process
  - Reads the dataset and generates a set of token statistics
  - Compares the generated token statistics to the statistics from benchmark datasets with known optimal parameters
  - Starts with the nearest known optimal parameters
  - Iterates over the dataset to refine the parameter settings to be optimal for the input dataset



PDP System Sequence Diagram (SSD)

# Hadoop-based Data Washing Machine (HDWM)

- The current design of the DWM is single-threaded:
  - Unable to go beyond a few thousands of records
    - Python version – up to 100, 000 records
    - Java version – up to 1, 000, 000 records
  - Goal – cluster more than a million (hundreds of million records)

- HDWM is a highly scalable version of the current Data Washing Machine designed with Hadoop MapReduce

# Hadoop-based Data Washing Machine (HDWM)

It solves the linear-based processing of DWM by:

- Using MapReduce programming model
- Ability to scale and cluster billions of records using HDFS

HDWM handles the "out-of-memory" problem caused by the creation of shared memory tables/dictionaries

- Carries intrinsic metadata alongside tokens
- E.g. "token: {refID, token, position, frequency}"

# Significant Accomplishments

- Developed Master Information Security and System Security Plans for ARP resources to allow for eventual work with CUI data (NIST 800.171)

- Demonstrated a foundation technology for detecting in real time websites disseminating illegal contents

- A positive data curation prototype demonstrated the ability to control both access and metadata reporting for data operations in the Hadoop environment.

- DART researchers continue to develop and refine a suite of novel algorithms (differential privacy preserving multi-party learning, fair and robust learning under sample selection bias or attacks, uncertainty award crowdsourcing, fraud and hate detection in cyberspace, user-centric data sharing in cyberspace, and privacy-preserving analytics in health and genomics).

Educated Public and Workforce

Big Data Management

Model Interpretation

Privacy and Security

# Role of Multimedia in Social Movement Mobilization

Studies of online Brazil and Peru political protests show that visual content on Twitter exhibits a greater degree of user engagement than text-only posts.



Diffusion of innovations theory

Diffusion of Innovations (DOI) theory (Rogers, 1962) is a social science theory that explains how new ideas, products, and technologies spread through society over time.

In this research, we apply DOI theory to online networks to evaluate the emerging adoption of information campaigns influenced by text, images, and video.

# Toxicity and Community Health

- We generate conversation trees for threads with > 50% toxicity and use machine learning to predict the leaf nodes' toxicity.
- The model predicts if the next reply will be toxic or not based on the structural characteristics of the conversation tree.
- Toxic conversation threads tend to have wider, deeper, and larger conversation branches.
- Those conversations are more likely to *end* with toxic comments.

Subreddit Post with comments-replies
Red = Toxic
Blue = non-Toxic





Leaf Node Class Distribution for Top 100 Trees with Most User Engagement



Plotting the Model Accuracies for 100 Conversation trees with Most Comments

| Model | Accuracy |
| --- | --- |
| LogisticRegression | 0.824 |
| KNeighborsClassifier | 0.804 |
| SVC | 0.823 |
| DecisionTreeClassifier | 0.739 |
| RandomForestClassifier | 0.8 |
| GradientBoostingClassifier | 0.823 |
| XGBClassifier | 0.814 |
| AdaBoostClassifier | 0.824 |
| MultinomialNaiveBayes | 0.823 |

COSMOS
Collaboratorium for Social Media and Online
Behavioral Studies

# Deep Learning for Preventing Discrimination and Hate Speech on Social Media

**Design and implementation of robust hate speech detection models via mitigating spurious correlations.**

- Automatic filtrations of hateful content have been deployed to prevent hate speech, but intentionally modifying hate words may bypass auto-detection systems.
- The research team developed a robust hate speech detection model by formulating a causal structure to represent the causal relationship among different variables and integrating the causal strength into a regularized cross-entropy loss for removing the spurious correlation.

(Lu Zhang)

# Cryptography-Assisted Secure and Privacy-Preserving Learning

**Privacy-Preserving Face Recognition Access Control**

- Goal: the backend server never sees the original face image or features

- Key idea: transform facial embeddings to a different space in a partially order-preserving way

- Transformed facial embeddings of the same user will still stay clustered together

- Transformed facial embeddings of different users will be relatively far away from each other

- Face matching model trained and queried based on transformed facial embeddings





Figure 2: Facial Reconstruction vs Proposed Method

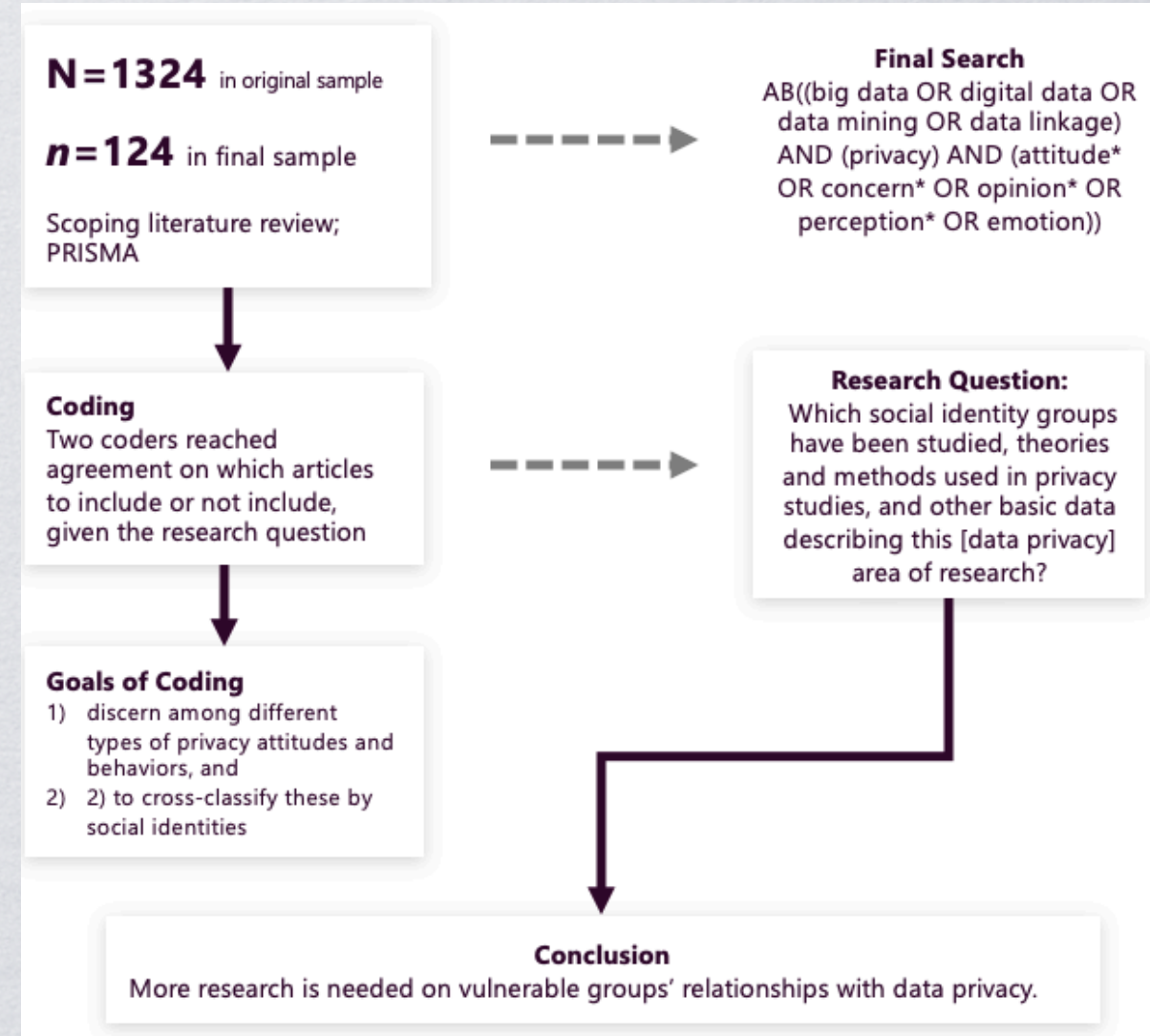| Transformation Parameters | True Verification Accuracy (%) | |
|---|---|---|
| | Training Data | Testing Data |
| No Transformation (Raw Embeddings) | 99.7% | 96% |
| 100 Partitions & 10e17 Partition Size | 99.7% | 96% |
| 100 Partitions & 10e6 Partition Size | 99.9% | 96.7% |
| 10 Partitions & 10e17 Partition Size | 99.9% | 96% |

# The interface between data privacy and user's social identities

- Systematic literature review of empirical research examining the intersection of social identities and privacy concerns within the big data context.

- Our goal is to better understand how users' social identities inform privacy attitudes and behaviors concerning digital technologies and big data.

- Context matters for privacy concerns, attitudes, and behaviors. Many of these studies did not report how privacy concerns relate to social identities. Although age and gender were frequently examined, many of these studies lacked a theoretical framework to generate a detailed explanation of their findings.



**N = 1324** in original sample

**n = 124** in final sample

Scoping literature review; PRISMA

**Final Search**
AB((big data OR digital data OR data mining OR data linkage) AND (privacy) AND (attitude* OR concern* OR opinion* OR perception* OR emotion))

**Coding**
Two coders reached agreement on which articles to include or not include, given the research question

**Research Question:**
Which social identity groups have been studied, theories and methods used in privacy studies, and other basic data describing this [data privacy] area of research?

**Goals of Coding**
1) discern among different types of privacy attitudes and behaviors, and
2) 2) to cross-classify these by social identities

**Conclusion**
More research is needed on vulnerable groups' relationships with data privacy.

# Mining cyberargumentation data for opinion evolution

**Deploy new version of cyberargumentation platform**
- Collect postings from 100s of University of Arkansas students debating 5 sociological topics (e.g., guns on campus) for future analysis.

**Develop expertise in transformer-based natural language processing (e.g., BERT) and apply it to two problems:**
- Sentiment analysis (IMDB movie reviews)
  - Publication:  G. Nkhata, U. Anjum and J. Zhan, "Movie Reviews Sentiment Analysis Using BERT," The Fifteenth International Conference on Information, Process, and Knowledge Management (eKNOW 2023), to appear.
- Hate speech detection
  - Publication:  X. Guo, U. Anjum and J. Zhan, "Cyberbully Detection Using BERT with Augmented Texts," 2022 IEEE International Conference on Big Data (Big Data), pp. 1246-1253.

Project leaders: Gauch, Adams

# Sentiment Analysis Architecture



Fig. 2. Fine-tuning of the model

To the best of our knowledge, this is the first work to couple BERT with BiLSTM
Results exceed state of the art accuracy on multiple datasets.

# Cyberbully Detection Architecture



Outperformed state of the art on smaller and/or unbalanced data sets

# Significant Accomplishments



- Successfully integrated CNN approaches on tabular data which incorporates neighborhood effects novel and potentially powerful ways

- Demonstrated a causal inference framework which can explain model predictions using unstructured data

- Developed a Boost-R solution to model stochastic event processes with heterogeneous features

- Bijective Maximum Likelihood image segmentation method achieved state-of-the-art image segmentation performance with no pixel-independent assumption in a tractable and invertible solution.

# Learning & Prediction (LP1): Statistical Learning – Random Forests for Recurrent Event Analytics led by Xiao Liu

- Developed a new, tree-based model selection algorithm which focuses on interaction among regressors outperforms by multiple measures other selection strategies for linear regression and a logistic regression model of repeat events

- Created a Random-Forest-Based algorithm and a Gradient-Boosted-Tree-Based algorithm for learning and prediction of recurrent event processes

- Contributed open-source code/tools for implementing the algorithms (available on GitHub)

# Learning & Prediction (LP3): Deep Learning – Novel Approaches led by Md Karim

- Developed PH-Net, a hybrid image processing model based on Persistent Homology (PH) and Fast Regional Convolutional Neural Network (FRCNN)

- Developed methods for image localization and image preprocessing using topological data analysis

- Performed an empirical study to determine dataset-specific usability of topological features

- Developed a new Self-supervised Domain Adaptation method in crowd counting

- Contributed a new Fairness Domain Adaptation (FREDOM) approach to semantic scene understanding

- Developed a new Self-supervised Spatiotemporal Transformers (SPARTAN) approach to group action recognition

# Learning & Prediction (LP4): Deep Learning – Efficiency and Specification led by Khoa Luu

- Developed self-supervised 3D capsule networks for medical segmentation on less labeled data

- Developed a new self-supervised domain adaptation deep learning method to deal with limited training data

- Developed a new equipollent domain adaptation approach for image deblurring

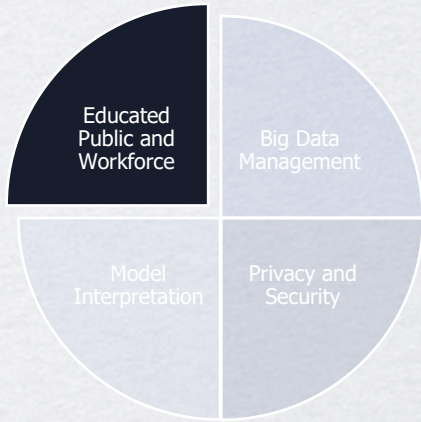- Proposed new fairness metrics to analyze the complexity and bias training data in deep learning

# Significant Accomplishments: ED, Workforce, Broadening Participation

# Significant Accomplishments

In 2019, Arkansas had no 4-year DS degrees

We now have programs implemented with students enrolled at the University of Arkansas, the University of Central Arkansas, and Arkansas State University.

Two-year college students and faculty are using shared computing resources (in ARP) for the first time

Arkansas Summer Research Institute keeps growing

Certificate proposals for 2-year campuses coming soon

Educated Public and Workforce

Big Data Management

Model Interpretation

Privacy and Security

# Education Theme – Year 3 Accomplishments

- **4-year Hubs for Common B.S. Data Science 8-semester Plan**
  - UAF is teaching all four years as of this year
  - UCA is evolving its curriculum to match the 8-semester plan
  - A-State is evolving its curriculum to match the 8-semester plan

- **ACTS for DASC 1003 Intro to Data Science has been approved**
  - DASC 1104 and 1223 will be proposed to be developed next

- **NorthArk and UAF are close to finalizing the model 2+2 plan**
  - This will facilitate +2 with UAF, UCA, and A-State

- **Study Abroad Program with the University of Nicosia is nearly complete**
  - Finalizing details for Year 2 Fall and/or Year 2 Spring
  - Designed to work for any institution following the 2+ or 8-semester plan
  - All courses our students would take are taught in English

# Education Theme – Year 3 Accomplishments (cont.)

- **Participated in many state-wide Data Science Workshops by ADHE**

- **Conference Papers / Presentations Submitted**
  - ACC Annual Fall 2023 Conference: "Data Science Careers Start at Community College"
    - Theme: Workforce Development
    - Presenter: Laura Berry, Interim Dean of Health Professions, North Arkansas Community College
      - Co-authors: Christine Davis, NWACC; Tina Moore, ADHE; Karl Schubert, UAF

  - ASEE 2023 Annual Conference:
    - Expanding & Improving a Multi-College Interdisciplinary B.S. Data Science Program with Concentrations
    - Theme: ASEE MULTI Division
    - Authors: Karl Schubert, UAF; Lee Shoultz, UAF, Shantel Romer, UAF
    - Note: builds on the previous notes on the Ed-Theme project
    - ASEE 2024 Paper submission will include co-leads & those active to cover the entire Ed-Theme at time of submission

# Broadening Participation



**JUNE 1 - 19 @ ONLINE**

Info & Application: tinyurl.com/apply2023asri

**Arkansas Summer Research Institute**
**2023**

**2.5 weeks of interactive data science training for US-based students**

ARKANSAS SCHOOL FOR MATH, SCIENCES, + THE ARTS

NSF

ARKANSAS NSF EPSCoR

**~300 applicants so far from all over the US**

**New sessions, new presenters, and improved flow**

**Extended by additional 3 days**

**Any US residents can apply (targeted for undergrads or beginners)**

# Broadening Participation

**SURE Program-** Summer Undergraduate Research Experience funding available, contact Brittany @ Brittany.Hillyer@arkansasedc.com

**Minigrants up to $5,000** to support professional development, special activities, and other related projects available. Check website for details
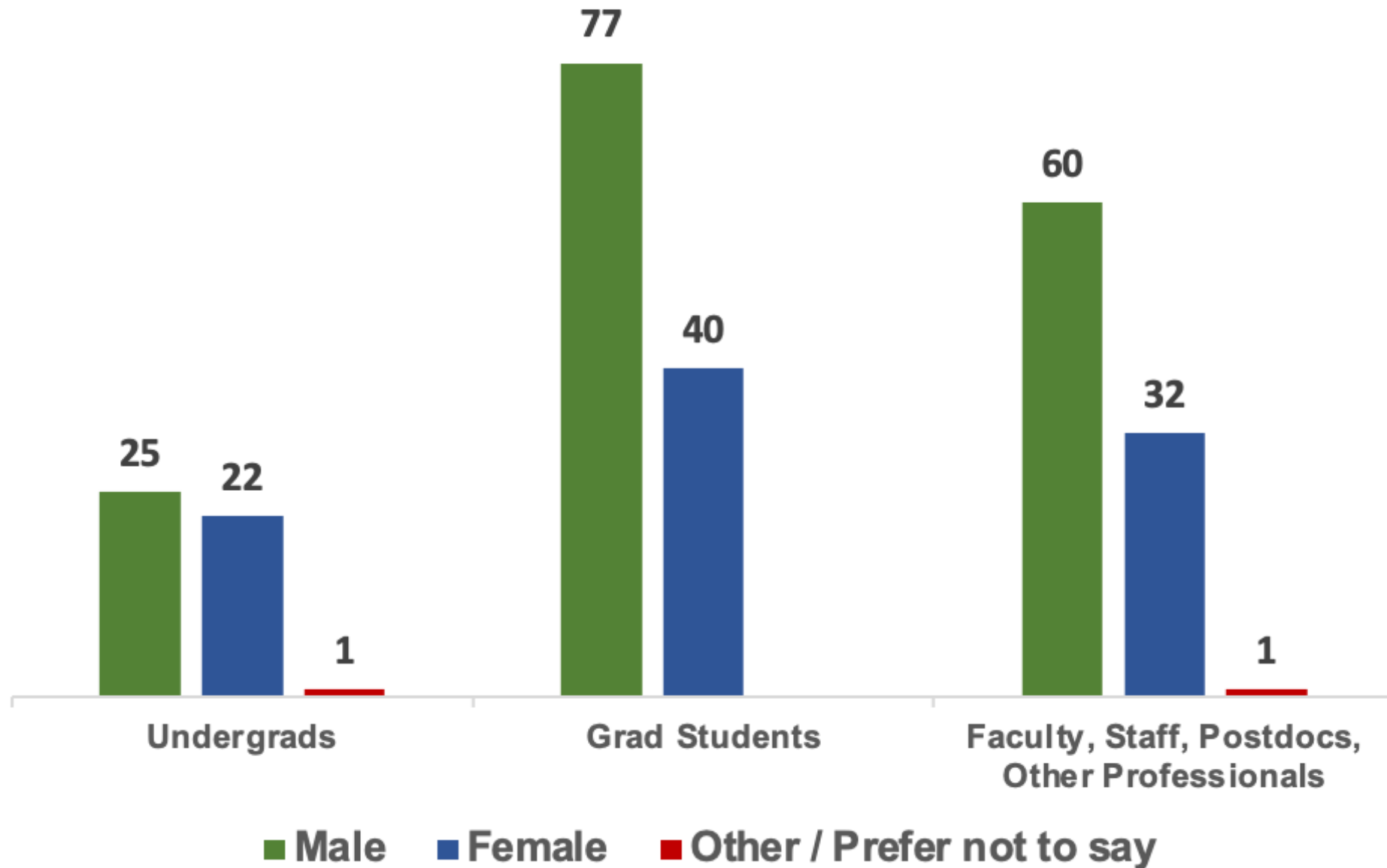
# DART Participants: Year 3

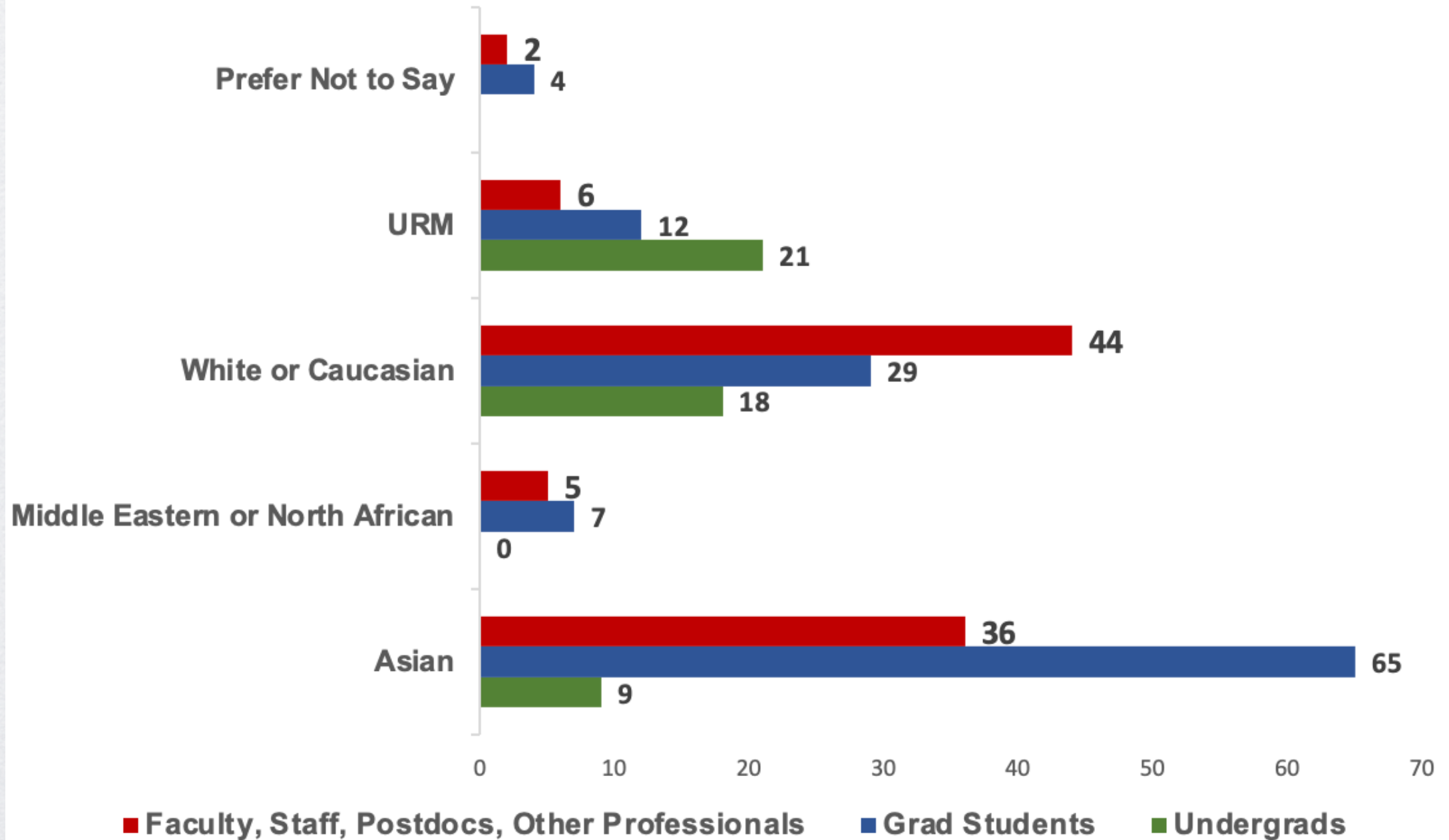| | |
|---|---|
| AEDC | **6** |
| Arkansas State University | **14** |
| Arkansas Tech University | **8** |
| North Arkansas College | **1** |
| Philander Smith College | **12** |
| SAU Tech | **1** |
| Shorter College | **3** |
| Southern Arkansas University | **17** |
| University of Arkansas Division of Agriculture | **1** |
| University of Arkansas Fayetteville | **79** |
| University of Arkansas Little Rock | **79** |
| University of Arkansas Medical Sciences | **21** |
| University of Arkansas Pine Bluff | **7** |
| University of Central Arkansas | **11** |

DART Y3 Participants

DART Y3 Participants

Demographics of DART Participants (Percentage of Whole by Role)

# Summary

- On track to meet strategic plan milestones and objectives
- Submitting a revision to strategic plan soon
- Made significant progress and have impactful results and collaborations
- Even more great things to share soon

DART

**Data Analytics that are Robust and Trusted**