# SoURCE CODE
# (Securing our Underlying Resource in Cyber Environments)
# Proposers' Day

Kristopher W. Reese, PhD | Program Manager | Oct. 5, 2023

Intelligence Advanced Research Projects Activity

# IARPA

Creating Advantage through Research and Technology

Thank you for your interest in this program and participating in this event

To assure a clear broadcast stream, audio and video are disabled for meeting participants

Comments and questions can be submitted to the IARPA team via the WebEx chat tool submission or via index cards for in-person attendees

- Please direct questions to "All Panelists" in the chat if you are virtual

Questions submitted to the alias (dni-SoURCE-CODE-proposers-day@iarpa.gov) prior to this meeting and during this presentation, and corresponding answers, may be posted in writing online

- This presentation is provided solely for information and planning purposes

- The Proposers' Day does not constitute a formal solicitation for proposals or proposal abstracts

- Nothing said at Proposers' Day changes the requirements set forth in a BAA

- **The BAA language supersedes anything presented or said by IARPA at the Proposers' Day**

- This meeting is being recorded and will be posted for public viewing

- For those viewing the recording, email aliases and POCs may be dated, please refer to IARPA.gov for updated information.

1. Familiarize participants with IARPA's interest in the SoURCE CODE program and solicit questions and feedback

2. Foster discussion of complementary capabilities among potential program participants, i.e., TEAMING

   - Teaming information can be found at the following address: https://www.iarpa.gov/research-programs/source-code

   - An attendance list, with contact information of participants who approved of sharing will be distributed soon

   - The chat feature is enabled for participants to plan future discussions associated with teaming

   - Teaming interests, capability summaries, and lightning talk slides will be posted publicly on the IARPA SoURCE CODE webpage until the BAA submission period closes

Please ask questions and provide feedback, this is your chance to alter the course of events.
Please talk with others, find great team members.

- Participants are encouraged to find partners and collaborators . . . someone might have a missing piece of your puzzle.

- Lightning talks will take place following the Program presentations.

- Collaborating and capability summaries will be accepted, with minimal review for appropriateness, and made available to the public.
  - Teaming documents and summaries can be submitted until the BAA closes, submit to [dni-SoURCE-CODE-proposers-day@iarpa.gov](mailto:dni-SoURCE-CODE-proposers-day@iarpa.gov).
  - If you would prefer your information not be shared (any recorded videos cannot be modified or removed) email dni-iarpa-source-code-proposersday@iarpa.gov.

- Questions can be submitted <u>until 11:00am ET</u>.

- There will be a break after the contracting presentation at 11:00am ET.

- Responses to selected questions will be broadcast at 12:30pm ET, so please don't log out or close your WebEx connection.

  - All programmatic and contractual questions will be captured but will not be answered in this session

- Feedback (but not questions) about the draft technical section may be submitted to the IARPA team email at [dni-SoURCE-CODE-proposers-day@iarpa.gov](mailto:dni-SoURCE-CODE-proposers-day@iarpa.gov).

  - A new alias will be established when the full BAA is released

- After this Proposers' Day, IARPA will review all the feedback received for a final BAA to be posted on SAM.gov.

# Agenda

| Time | Topic | Speaker |
|------|-------|---------|
| 9:30am-9:40am | Welcome, Logistics, Proposers' Day Goals | Kristopher W. Reese, Program Manager |
| 9:40am-9:50am | IARPA Overview | Robert Rahmer, Director Office of Analysis Research, IARPA |
| 9:50am-10:40am | SoURCE CODE Program Overview | Kristopher W. Reese |
| 10:40am-11:00am | Contracting Overview | TBD |
| 11:00am-12:30pm | Break (Submit questions in chat or drop boxes before 11:00am) | |
| 12:30pm-1:30pm | Answers to Selected Technical Questions | Kristopher W. Reese |
| 1:30pm-1:35pm | Introductions to Lightning Talks | Kristopher W. Reese |
| 1:35pm-4:25pm (est.) | Lightning Talks | Potential Performers |
| 4:25pm-5:30pm | Informal Teaming Discussions* | In-Person Participants |

*The Government will not attend these events

# LIGHTNING TALKS AGENDA

| Time | Speaker | Institution | In person |
|------|---------|-------------|-----------|
| 1:35pm-1:40pm | Xiangyu Zhang / Lin Tan | Purdue | Yes |
| 1:40pm-1:45pm | Kexin Pei | University of Chicago | Yes |
| 1:45pm-1:50pm | Michael V Le | IBM | Yes |
| 1:50pm-1:55pm | Mike Murphy | SimSpace | Yes |
| 1:55pm-2:00pm | Shiqing Ma | UMASS Amherst | Yes |
| 2:00pm-2:05pm | William Liu | CACI | Yes |
| 2:05pm-2:10pm | Sheikh Rabiul Islam | Rutgers University - Camden | Yes |
| 2:10pm-2:15pm | Andrew Hendela | Karambit.AI | Yes |
| 2:15pm-2:20pm | Aleksey Nogin | Red Balloon Security | Virtual |

| Time | Speaker | Institution | In person |
|------|---------|-------------|-----------|
| 2:20pm-2:25pm | Thomas Wahl | GrammaTech | Yes |
| 2:25pm-2:30pm | Dan Thomsen | SIFT | Virtual |
| 2:30pm-2:35pm | Nathan Clark | Noblis | Yes |
| 2:35pm-2:40pm | Chris Taylor | Tactical Computing Labs | Yes |
| 2:40pm-2:45pm | Tomas Pevny | Czech Technical University in Prague | Virtual |

Break and Informal Teaming Discussion at end of talks.

IARPA envisions and leads *high-risk, high-payoff research* that delivers innovative technology *for future overwhelming intelligence advantage*

- Our problems are complex and multidisciplinary
- We emphasize technical excellence & technical truth

- **Bring the best minds to bear on our problems**
  - Full and open competition to the greatest possible extent
  - World-class, term-limited Program Managers

- **Define and execute research programs that:**
  - Have goals that are clear, ambitious, credible and measurable
  - Run from three to five years
  - Publish peer-reviewed results and data, to the greatest possible extent
  - Employ independent and rigorous Test & Evaluation
  - Involve IC partners from start to finish
  - Transition new capabilities to intelligence community partners

- Technical <u>and</u> programmatic excellence are required

- Each program has a clearly defined and measurable end-goal

  - Intermediate milestones to measure progress are also required

  - Every program has a beginning and an end

- This approach, coupled with term-limited PM positions, ensures

  - IARPA does not "institutionalize" programs

  - Fresh ideas and perspectives are always coming in

  - Status quo is always questioned

  - Only the best ideas are pursued, and only the best performers are funded

IARPA's research portfolio is diverse, including math, physics, chemistry, biology, microelectronics, neuroscience, linguistics, political science, cognitive psychology, and more.

- 70% of completed research transitions to U.S. Government partners

- 3,000+ journal articles published

- IARPA funded researchers have been awarded the Nobel Prize in Physics for quantum computing research, a MacArthur Fellowship, and a Bell prize

- IARPA serves on National Science and Technology Council (NSTC) committees and actively engages with the White House BRAIN Initiative, National Strategic Computing Initiative, and the NSTC Select Committee on Artificial Intelligence, the NSTC Subcommittee on Quantum Information Science (SCQIS), and NSTC Subcommittee on Economic and Security Implications of Quantum Science (ESIX)

## ENGAGE WITH US

Throughout our website you can learn more about engaging with us on our highly innovative work that is having a positive impact in the Intelligence Community and society in general. Click on any of the below links to learn more.

### iarpa.gov | 301-243-1995

dni-iarpa-info@iarpa.gov

- Reach out to our Program Managers.
- Schedule a visit if you are in the DC area or invite us to visit you

**Open BAAs**
Broad Agency Announcements (BAAs) solicit research proposals for specific programs. Learn more about current BAA opportunities and ways to get involved...

**Requests For Information**
Requests for Information (RFIs) are designed to gather more information on an idea in an area in which our program managers are not fully informed...

**Seedlings**
Seedlings are typically 9 – 12 month research efforts that are less than $1M in cost. They are intended to address highly innovative ideas and concepts within...

# SoURCE CODE Overview

Kristopher W. Reese, PhD | Program Manager | Oct. 5, 2023
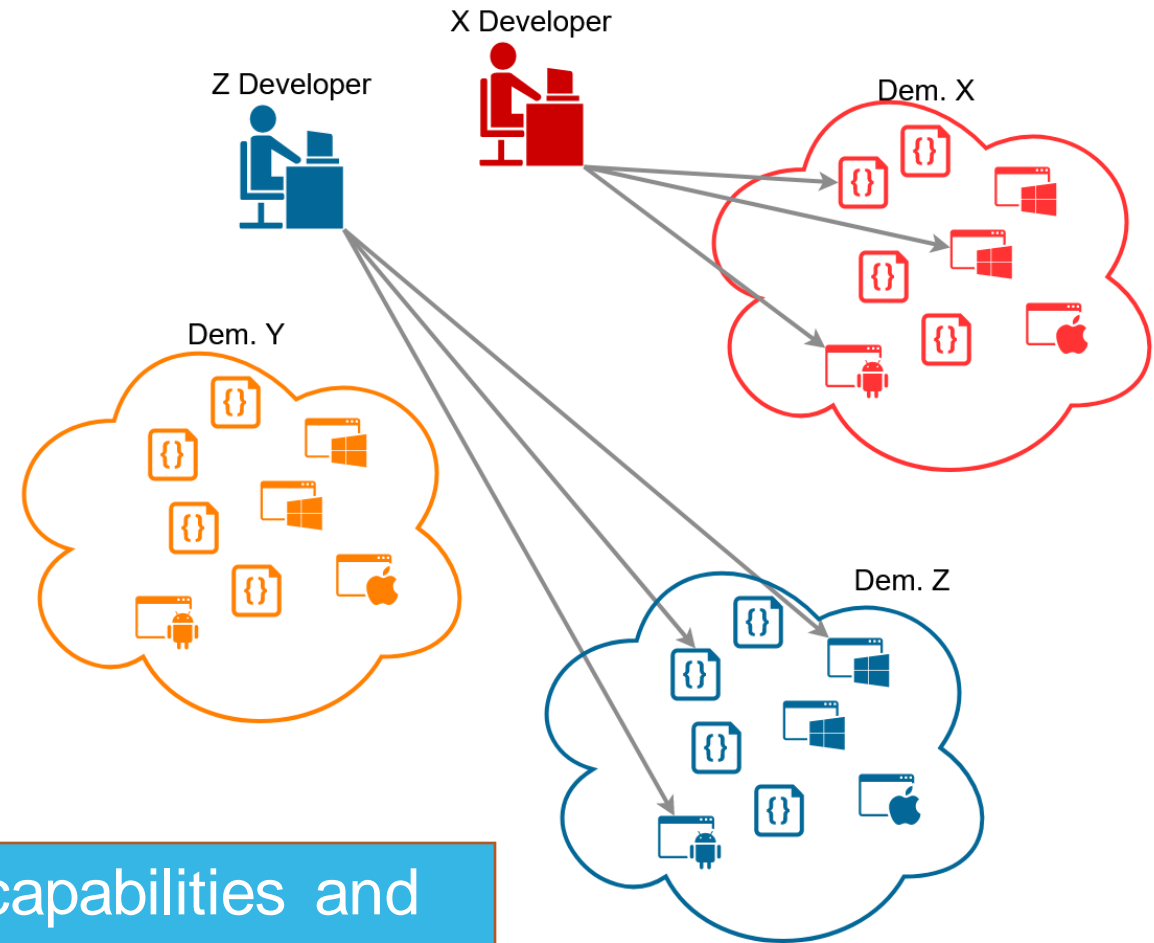
Intelligence Advanced Research Projects Activity

IARPA

Creating Advantage through Research and Technology

SoURCE CODE seeks to create automated, scientifically validated forensic similarity and demographic analytic technologies

- Measure similarity of code and binaries and identify components that may analyze hidden demographic information.



(U) SoURCE CODE will improve forensic capabilities and speed up threat intelligence analytics!

Image is UNCLASSIFIED

- Current methods are highly manual, requiring substantial human expertise, training, and time to conduct forensic analysis of code.
  - Prior automated attempts leverage a small subset of static analysis features Lexical and syntactic.

- Executable binaries and corresponding source code include numerous other features that can measure similarity, and thereby assist in attribution and demographic analytics.

- SoURCE CODE will develop capabilities to use the full feature space to measure similarity, especially between source code and binaries

(U) The full feature set of source code and binary features are an untapped resource for forensic capabilities.

[2]

[3]

[5]
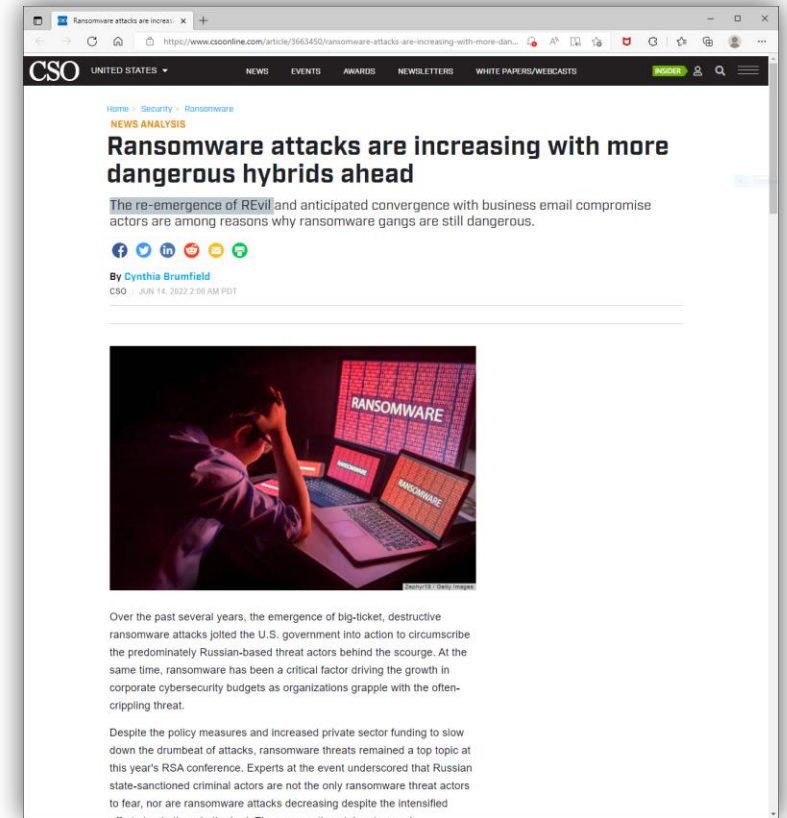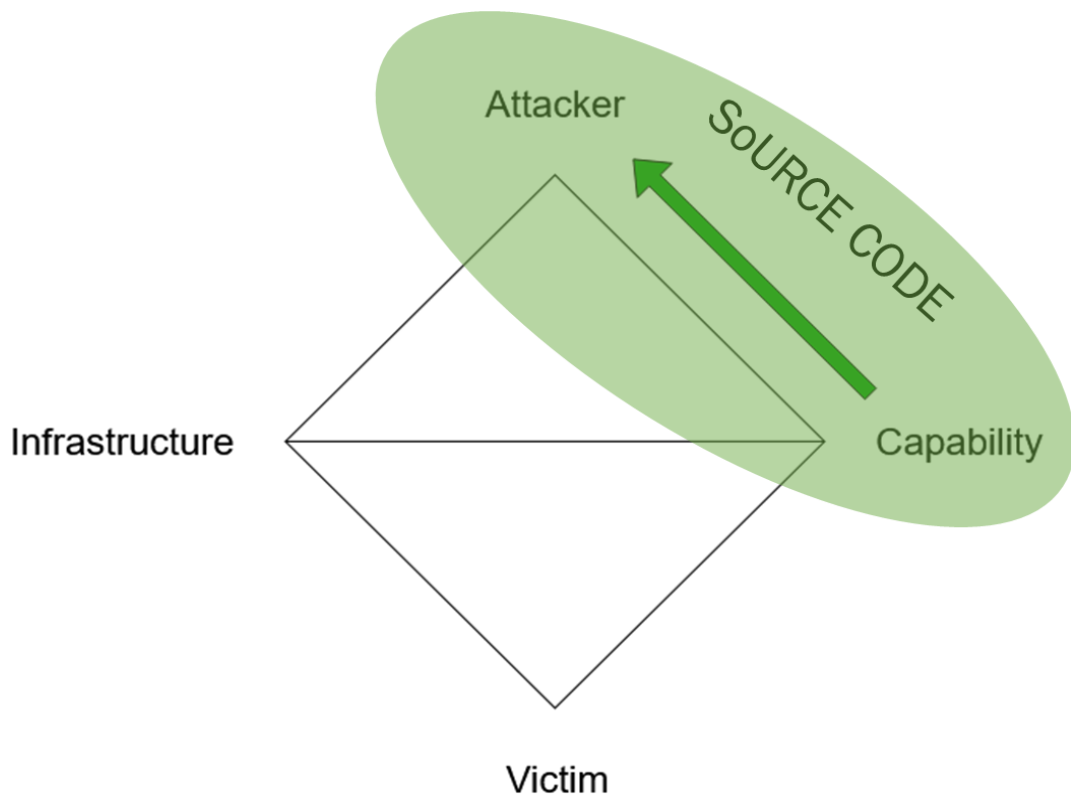
Image is UNCLASSIFIED

Similarity and Demographic traits can play a role in attributing the increasing number of attacks around the world.

SoURCE CODE will create similarity and demographic analytics to automate the "Attacker-Capability" edge.

Image modeled after original research paper [38]
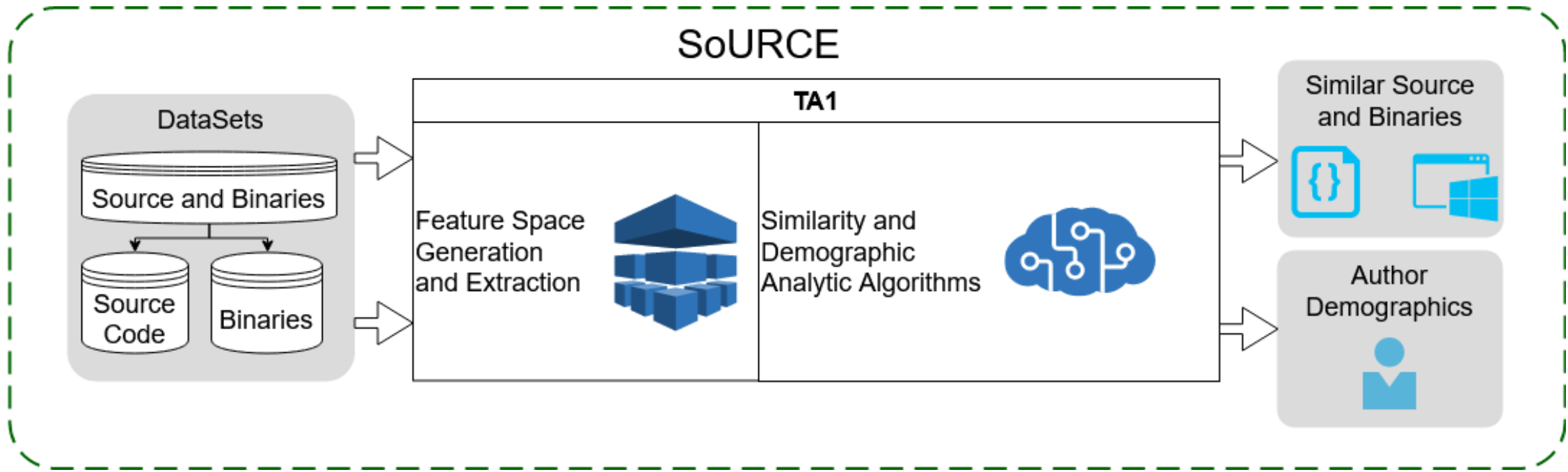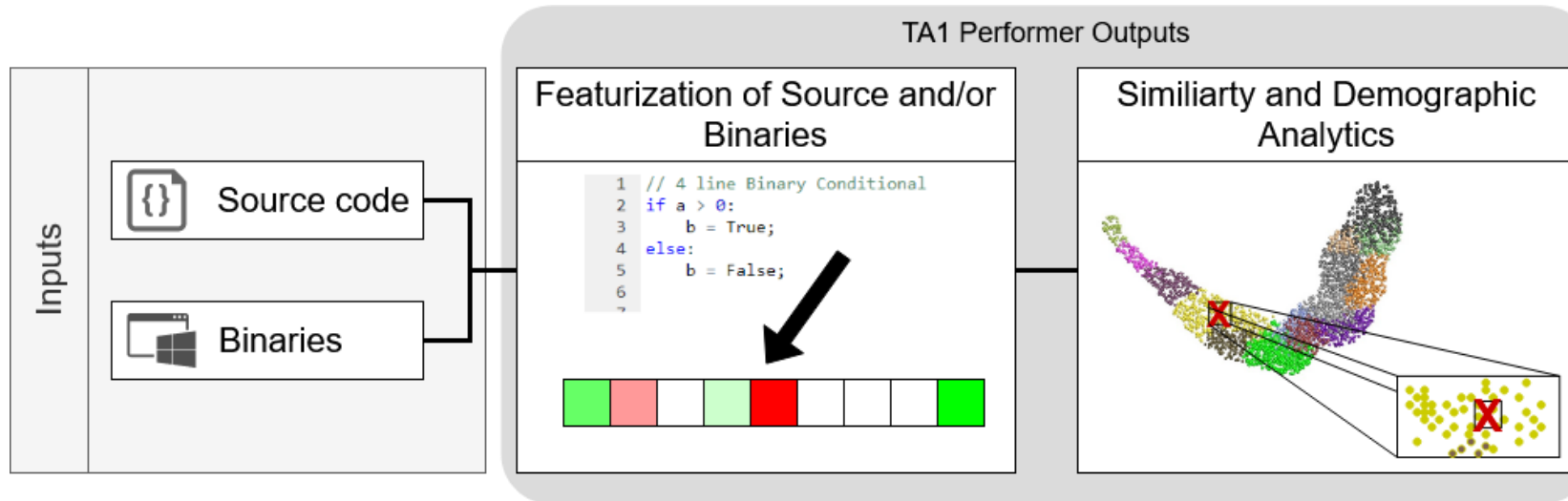
Image is UNCLASSIFIED

# Technical Areas

Image is UNCLASSIFIED

SoURCE CODE seeks to utilize the full feature set (focus area) to measure source code and binary similarity and demographic attributes (e.g. Country, Group, individual).

## INNOVATIONS REQUIRED

1. Identify and map salient features in both binaries and/or source code that capture author style.

2. Identify and implement algorithms, or ensembles, that can effectively utilize authorship features to measure similarity and identify most likely author(s)
   a) **PHASE 1:** binary to binary, source code to source code
   b) **PHASE 2:** binary to source code, source code to binary

3. Explain similarity score / decision process to forensic experts to assist in making final attribution decisions.
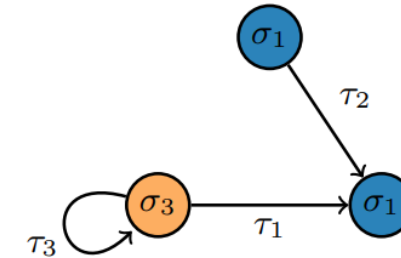
| | Code Properties | Examples |
|---|---|---|
| **Disassembled** | **Instruction** | Byte level n-grams<br>Idioms<br>Instruction Summary Graphlets<br>SuperGraphlets |
| | **Control Flow** | Instruction Summary Graphlets<br>SuperGraphlets<br>Call Graphlets |
| | **External** | Call Graphlets<br>Library Calls |
| **Decompiled** | **Lexical** | Word Unigrams:<br>• Integer types<br>• Names of library functions<br>• Names of Internal functions (when symbol information is available) |
| | **Syntactic** | Fuzzy Abstract Syntax Trees<br>• AST n-grams<br>• Labeled AST edges<br>• AST Node TF-IDF<br>• AST Node Avg. depth |

$$u_1 = (\text{push ebp} \mid * \mid \text{mov esp,ebp})$$

*Example 1: Code Idioms*

```
        cpuid
        jmp L2
        ...
L1:
        cmp ecx,edx
        jle L1
L2:
        mov eax, 0x5
        sysenter
```

*Example 2: Graphlets*

Images screen captured from original research paper [21]

All images are UNCLASSIFIED

Decompiled features are similar to source code features

There are limited features in this space, and a complete exploration of possible features needs to be conducted – including learned features!

Harder to utilize at scale

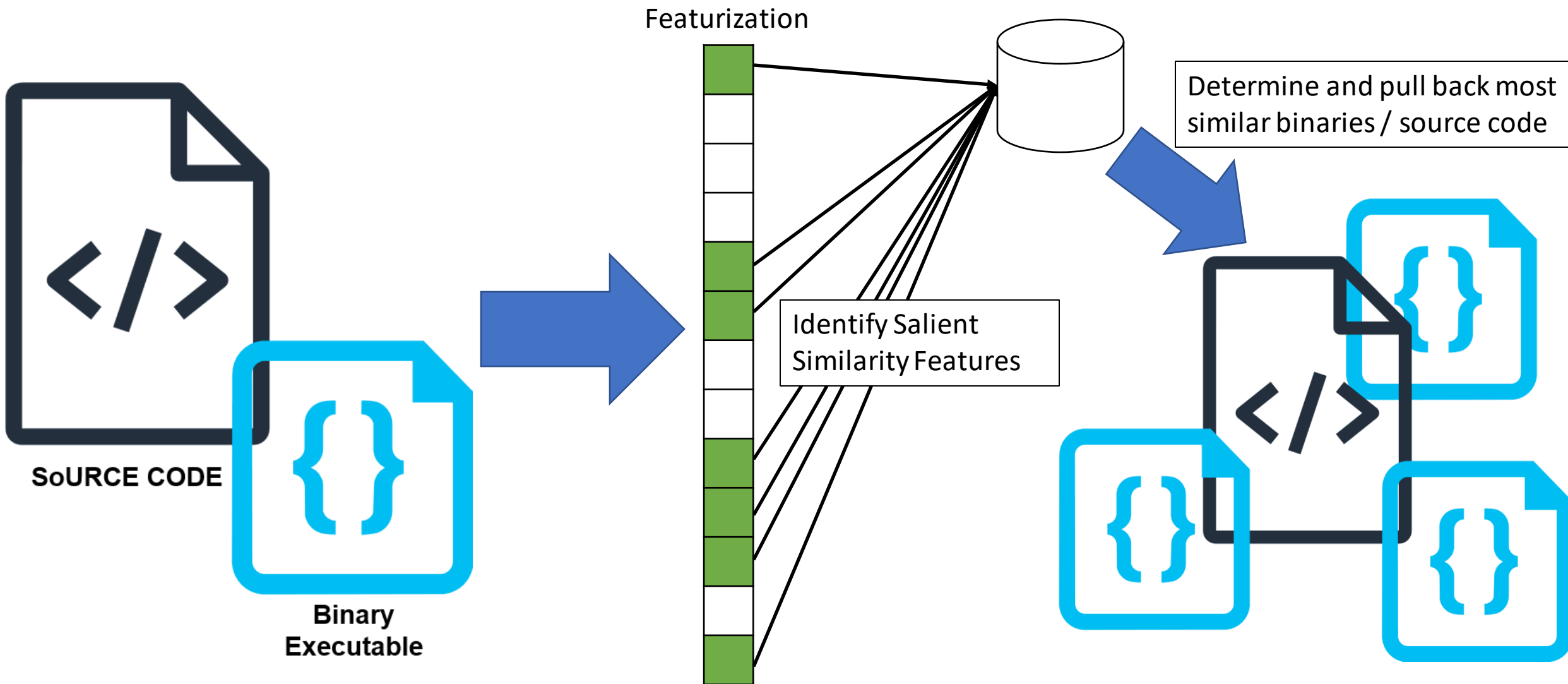| Feature Type | Examples |
|---|---|
| **Lexical** | Lines of Code |
| | Operands |
| | Variables |
| | Spaces |
| | Word n-grams |
| | Char n-grams |
| | Function names |
| **Syntactic** | Average function size |
| | Special Macros |
| | Data Structure choice |
| | Control Structure choice |
| | Input Statements |
| | Conditional Statements |
| | Assignment Statements |
| **Semantic** | Loops |
| | Dataflow analysis |
| | Control flow analysis |
| | Algorithms implemented |
| | Procedure-dependent analysis |
| **Behavioral** | System calls |
| | Files accessed |
| | Created mutex |
| | Visited URLs |
| | Dynamic values |
| | Network connections |
| **App-dependent** | Log file strings |
| | Error message file strings |
| | Property file strings |

Current source code Authorship systems largely exploit Lexical and Syntactic structures, limiting forensic applications to unobfuscated / de-linted code
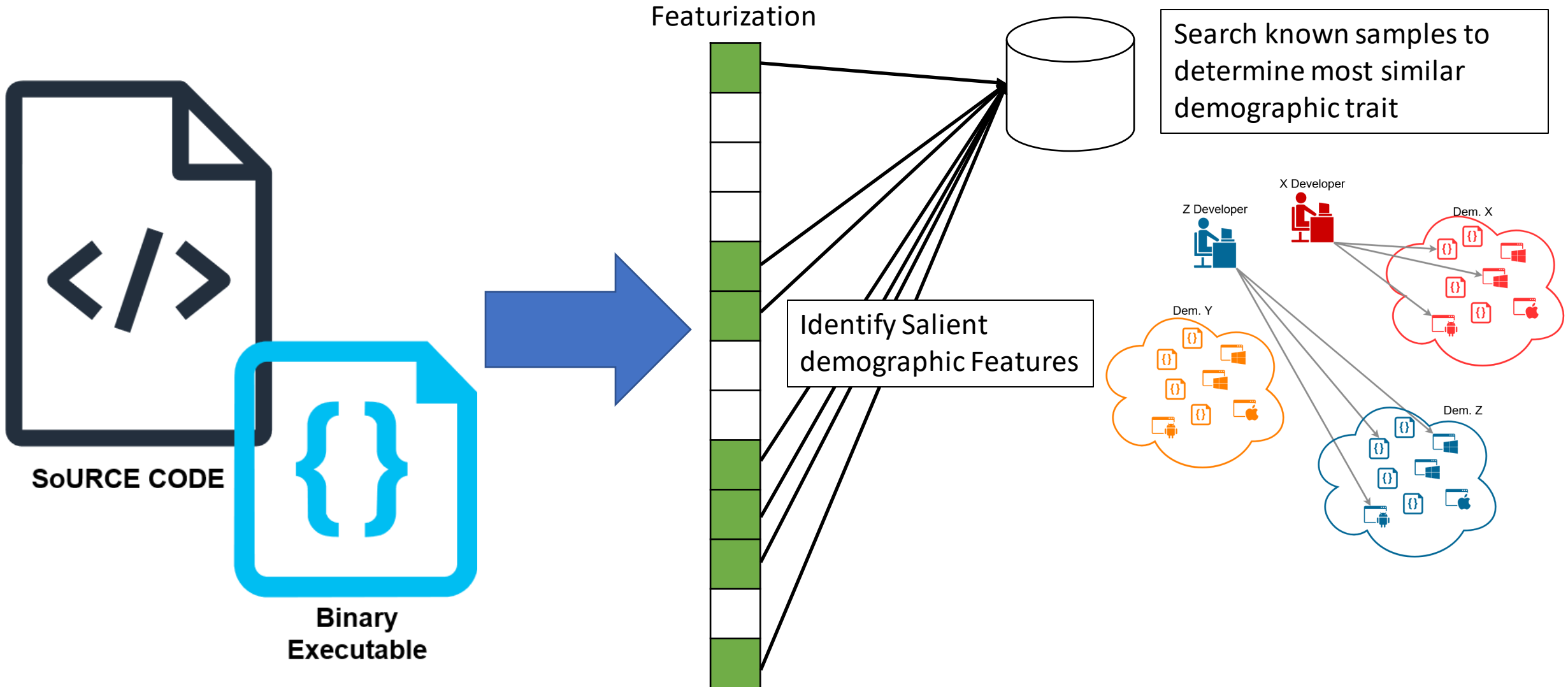
SoURCE CODE seeks to exploit the full feature set!

All images are UNCLASSIFIED

Featurization

Determine and pull back most similar binaries / source code

Identify Salient Similarity Features

SoURCE CODE

Binary Executable

Featurization

Search known samples to determine most similar demographic trait

Identify Salient demographic Features

SoURCE CODE

Binary Executable

- The following scientific gaps in authorship attribution of source code / binaries have been identified as being underexplored in literature, and may impact features identified making authorship determination more difficult:
  - Codebases vs a single Individual's code
  - Understanding the impact and influences of Project Domains: e.g. Android vs. iOS vs. Windows vs. Linux
  - Understanding the impact and influences of specific development tools
    - Integrated Development Environments
    - Version Control Systems
    - Compiler, Build Environment, and Deployment tools
  - Impact of project naming conventions, company style guides, etc.

These gaps highlights the need for understanding the impact of standard coding practices over educational coding practices!

# SoURCE CODE Program Phases

| Task | # Mo. | Phase 1 | | | | | | | | | | | | | | | | | | Phase 2 | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| Kickoff | - | ★ | | | | | | | | | | | | | | | | | | ★ | | | | | | | | | | | |
| TA1: Feature Space | 30 | | | | ★ | | | | ★ | | | | ★ | | | | ★ | | | | | | ★ | | | | ★ | | | | ★ |
| TA1: Similarity | 30 | | | | | | | | ★ | | | | ★ | | | | ★ | | | | | | ★ | | | | ★ | | | | ★ |
| PI Virtual Meetings / Calls | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| PI Workshop (in-person) | - | | | | | | | | ★ | | | | | | | | ★ | | | | | | | | | | ★ | | | | |
| Site Meetings | - | | | | ★ | | | | | | | | ★ | | | | | | | | | | ★ | | | | | | | | |
| Program Closeout | - | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ★ |

★ Milestone Deliverables

Image is UNCLASSIFIED

**Phase 1**

Binary

C++

Rust

**Phase 2**

Binary

C++

Rust

- Phase 1 programming languages pending – subject to data availability.

- Phase 2 programming languages subject to data availability

- Potential Proposers' can offer potential datasets
  - Proprietary Datasets cannot be approved unless they can be shared with all potential performers on the overall program or purchased for research purposes.

- Google Code Jam:
  - Datasets from 2008-2020 competitions
  - Contains Author data as ground truth
  - No binaries – Can compile with different build environments

- Other Coding Competitions:
  - Codeforces.com (Russian-based; Coding Competitions)
  - Topcoders (US-based; Coding Competitions)

- Other datasets identified from RFI will be checked by T&E

(U) These competition datasets act as surrogates for Malware, but they do not represent Malware source code!

- APTClass [31]
  - One of the Largest ground-truth datasets (15,000 samples)
  - Need to Request Access to download (UK University)
  - Data Sources / Ground Truth:

| Source Name | Last Updated |
|---|---|
| MISP [32] | Oct. 2020 |
| APT Operation Tracker [33] | Oct. 2020 |
| MITRE ATT&CK [34] | Oct. 2020 |
| sapphirex00 [35] | Nov. 2018 |
| Thailand CERT [36] | Oct. 2020 |
| Council on Foreign Relations [37] | Oct. 2020 |

  - Other datasets identified from RFI will be checked by T&E

# Evaluation and Metrics

Robust, independent test and evaluation is a crucial part of every IARPA program

- For SoURCE CODE, T&E will be responsible for providing data and product evaluation.
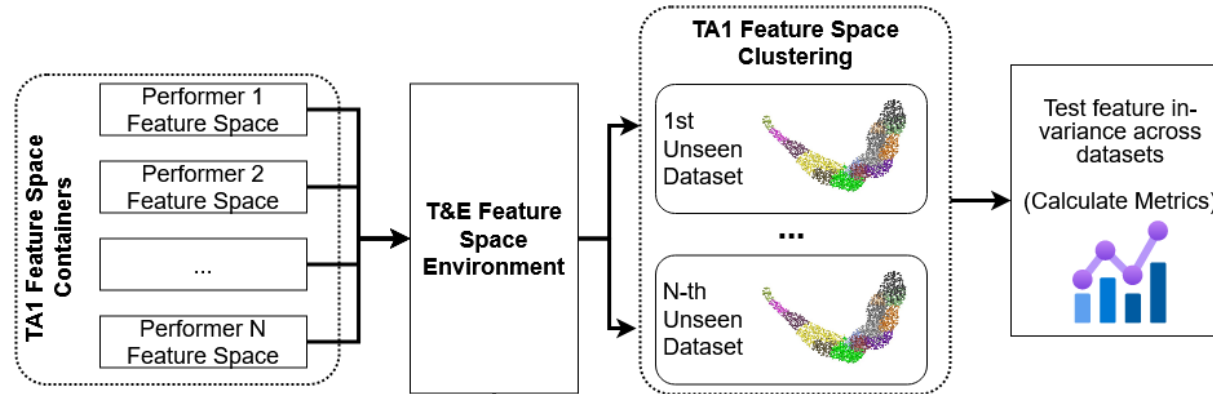
Performer systems will be executed by T&E on stand alone systems / networks

- Specifications of the SoURCE CODE Test System will be provided at a later date
- SoURCE CODE anticipates using multiple T&E teams for various aspects of the overall program.
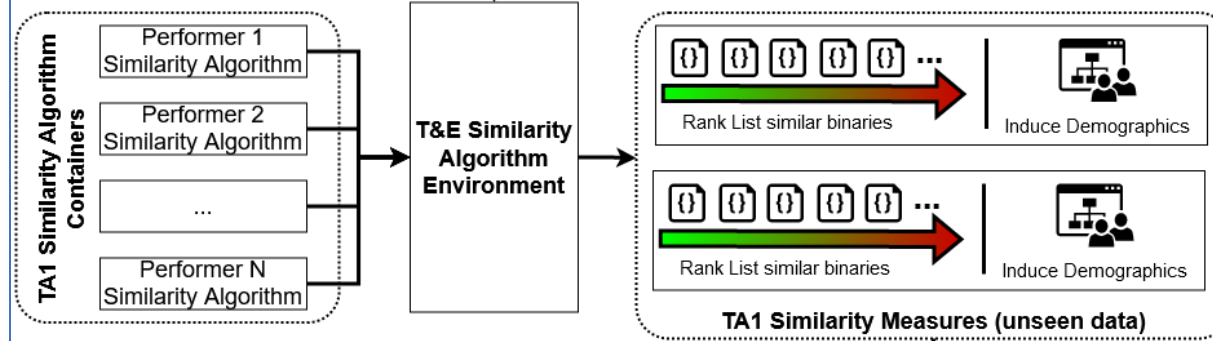
1) T&E will provide interface for automatically running experiments in the feature space and similarity algorithms.

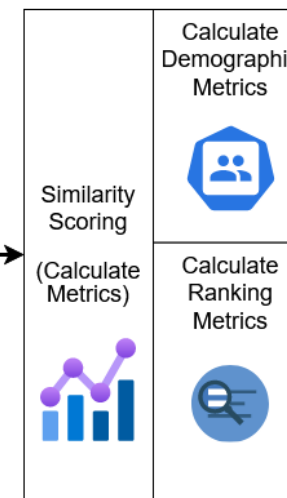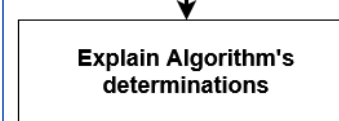2) Performers will submit Feature Space findings and extractions to T&E.

3) T&E will evaluate feature space for invariance across datasets to ensure features are salient.

4) Performers will submit similarity and demographic algorithms to T&E for evaluation and experimentation

5) T&E will compare algorithms to baseline systems and calculate metrics in similarity scoring and demographic analytics

6) T&E will conduct analysis of explanation algorithms for generalizability



**TA1 Feature Space Clustering**

TA1 Feature Space Containers
- Performer 1 Feature Space
- Performer 2 Feature Space
- ...
- Performer N Feature Space

T&E Feature Space Environment

1st Unseen Dataset
...
N-th Unseen Dataset

Test feature invariance across datasets (Calculate Metrics)

TA1 Similarity Algorithm Containers
- Performer 1 Similarity Algorithm
- Performer 2 Similarity Algorithm
- ...
- Performer N Similarity Algorithm

T&E Similarity Algorithm Environment

Rank List similar binaries — Induce Demographics

Rank List similar binaries — Induce Demographics

**TA1 Similarity Measures (unseen data)**

Explain Algorithm's determinations

Similarity Scoring (Calculate Metrics)

Calculate Demographic Metrics

Calculate Ranking Metrics

| | Phase I | Phase II |
|---|---|---|
| Top-1 Accuracy [100/10,000] | 85% [95/75] | 90% [98/90] |
| Top-10 Accuracy [100/10,000] | 95% [99/85] | 96% [99/95] |
| EER Average | 30% | 20% |
| d' (sensitivity index) | Measured | increase 30% |
| Attribution methods (data permitting) | Source → Source, Binary →Binary | Source → Binary, Binary →Source |

Metrics shown are for 1000 unique authors and should scale to match cardinality of unique users – additional numbers shown for 100 / 10,000 authors in brackets.

Additional metrics will also be measured to better understand the efficacy of the algorithms:
- Precision
- Specificity
- FAR/TAR
- FRR/TRR

Detection Error Tradeoff / Receiver Operator Curve

TABLE is UNCLASSIFIED

|  | Phase I | Phase II |
|---|---|---|
| Top-1 Accuracy [100/10,000] | 85% [95/75] | 90% [98/90] |
| Top-10 Accuracy [100/10,000] | 95% [99/85] | 96% [99/95] |
| EER Average | 20% | 10% |
| # Groups/Demographics (data permitting) | 50 | 70+ |
| Type of Set | Closed | Open |

Additional metrics will also be measured to better understand the efficacy of the algorithms:
- Precision
- Specificity
- FAR/TAR
- FRR/TRR

→ Detection Error Tradeoff / Receiver Operator Curve

TABLE is UNCLASSIFIED

|          |          | Actual |          |
|----------|----------|--------|----------|
|          |          | Positive | Negative |
| Predicted | Positive | True Positive (TAR) | False Positive (FAR) |
|          | Negative | False Negative (FRR) | True Negative (TRR) |

Binary classification confusion matrix

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Recall = \frac{TP}{TP + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Specificity = \frac{TN}{TN + FP}$$

*(Recall is best when we don't care about the misclassifications of negative samples).*

TP (TAR) – Code sample from Author A matches to Suspect Author (Author A)

FP (FAR) – Code sample from Author B matches to Suspect Author (Author A)

FN (FRR) – Code sample from Author A does not match Suspect Author (Author A)

TN (TRR) – Code sample from Author B does not match Suspect Author (Author A)
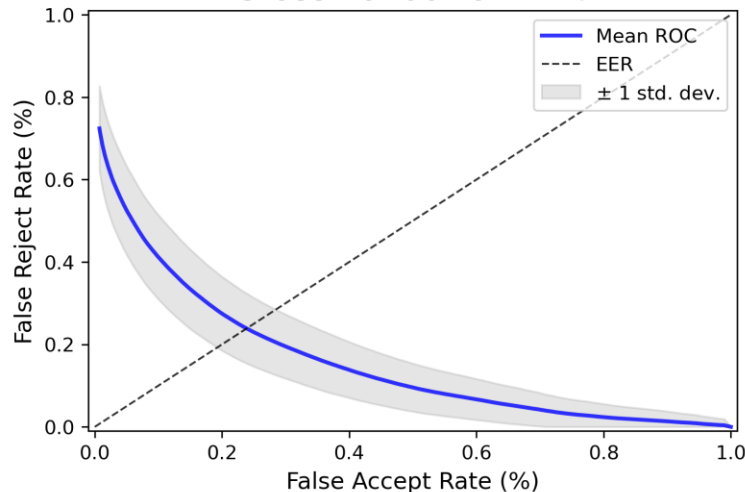
Table and Equations are UNCLASSIFIED

Cross-Validation ROC


Cross-Validation DET

**Equal Error Rate (EER)** is the point at which the proportion of the False Acceptance and False Rejection Rates are equal.

*At what point do we crossover where code samples match the INCORRECT authors and code samples are rejected from CORRECT authors?*

**Receiving Operating Characteristic (ROC)** shows the probability of detection against the probability of false alarm and helps to identify characteristics of the attribution system.

*If we fix rate at which code samples match INCORRECT authors to a specific percentage, what rate will we achieve with correct code samples?*

**Detection Error Tradeoff (DET)** curves map the probability of false alarms against the probability of falsely rejecting a valid author.

*If we fix rate at which code samples math INCORRECT authors to a specific percentage, what rejection rate will we achieve with correct code samples?*
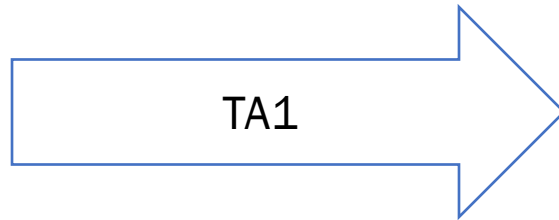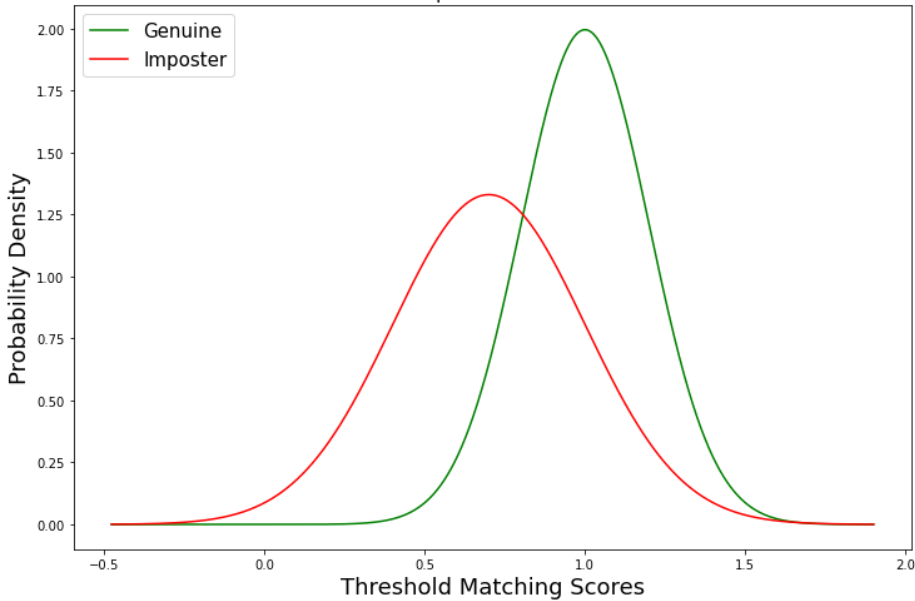
Images are UNCLASSIFIED

(U) *d' is a dimensionless metric indicating the discriminability between two signals – in this case, genuine-imposter scoring from an algorithm's attribution matching.*
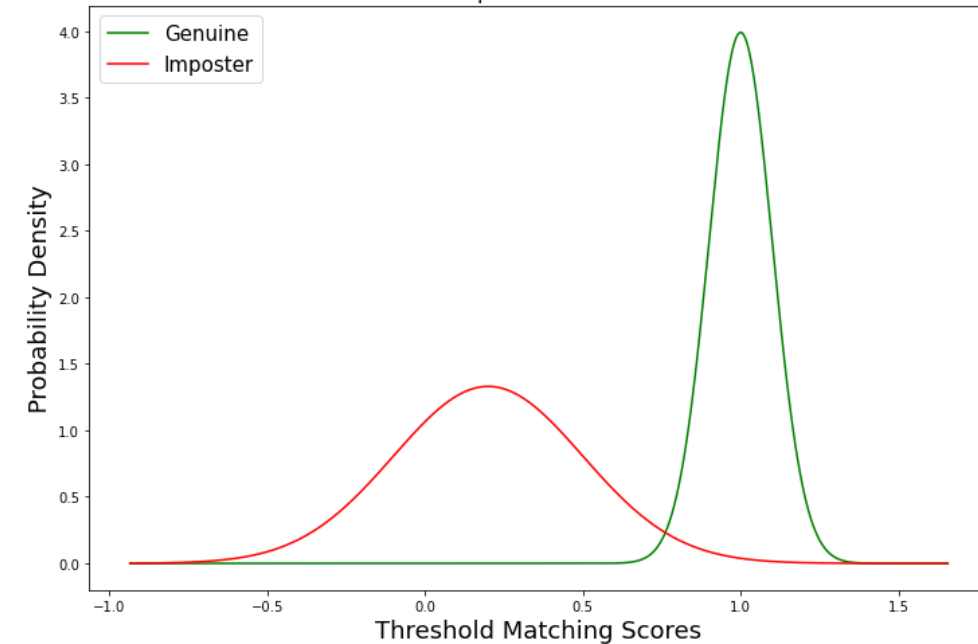


**TA1** – Can we improve the separation between signal distributions of code samples matching correctly (genuine) and code samples matching incorrectly (imposter)

Image is UNCLASSIFIED

Feedback, thoughts and comments:

- SoURCE CODE Team Alias: dni-SoURCE-CODE-proposers-day@iarpa.gov

Additional information:

- SoURCE CODE website: https://www.iarpa.gov/research-programs/source-code.