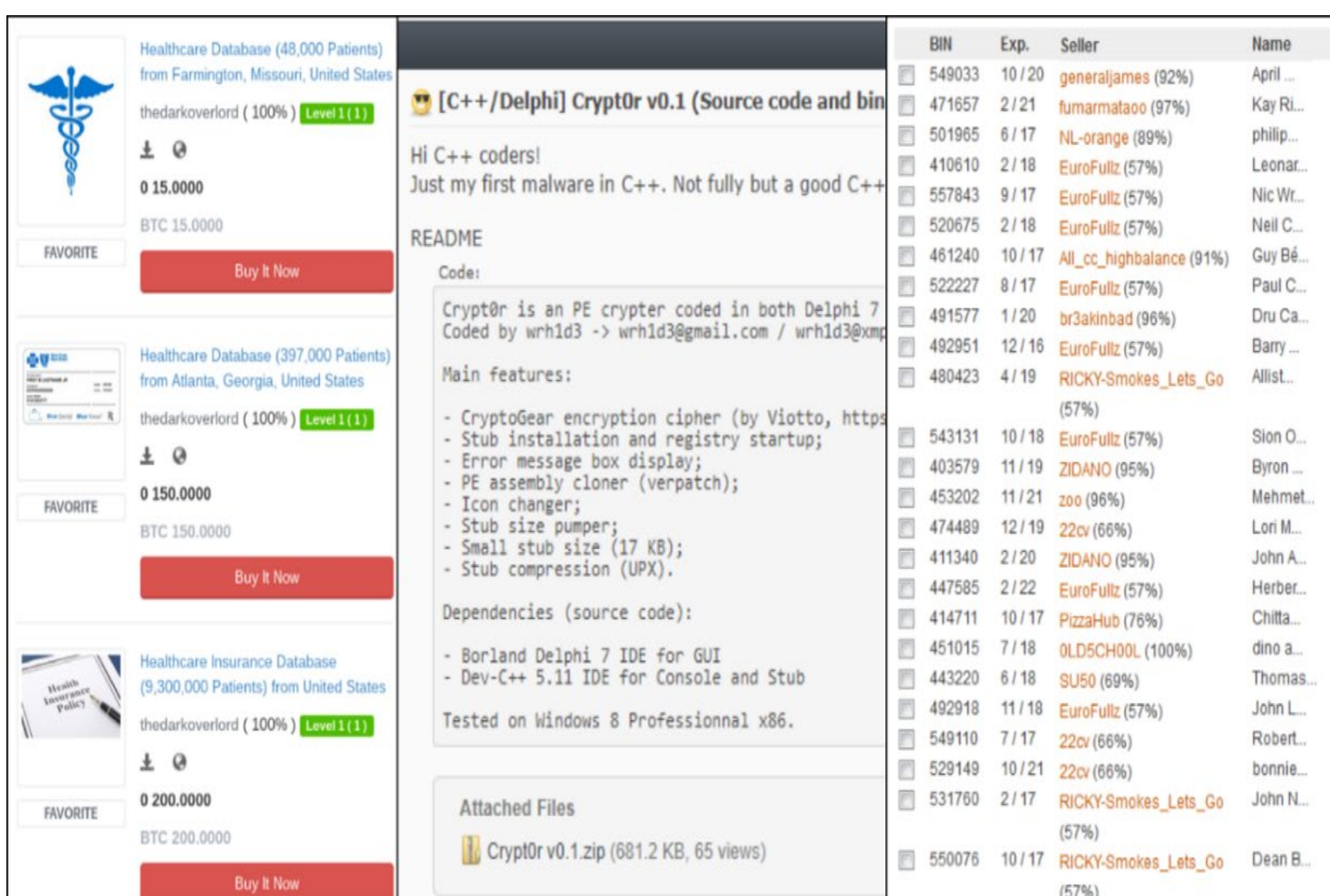# SaTC: CORE: Small: Cybersecurity Big Data Research for Hacker Community: A Topic and Language Modeling Approach

PI: Dr. Hsinchun Chen, Regents' Professor, ACM/IEEE Fellow, University of Arizona, AI Lab Director

Co-PI: Dr. Weifeng Li, University of Georgia

https://eller.arizona.edu/departments-research/centers-labs/artificial-intelligence/research/big-data
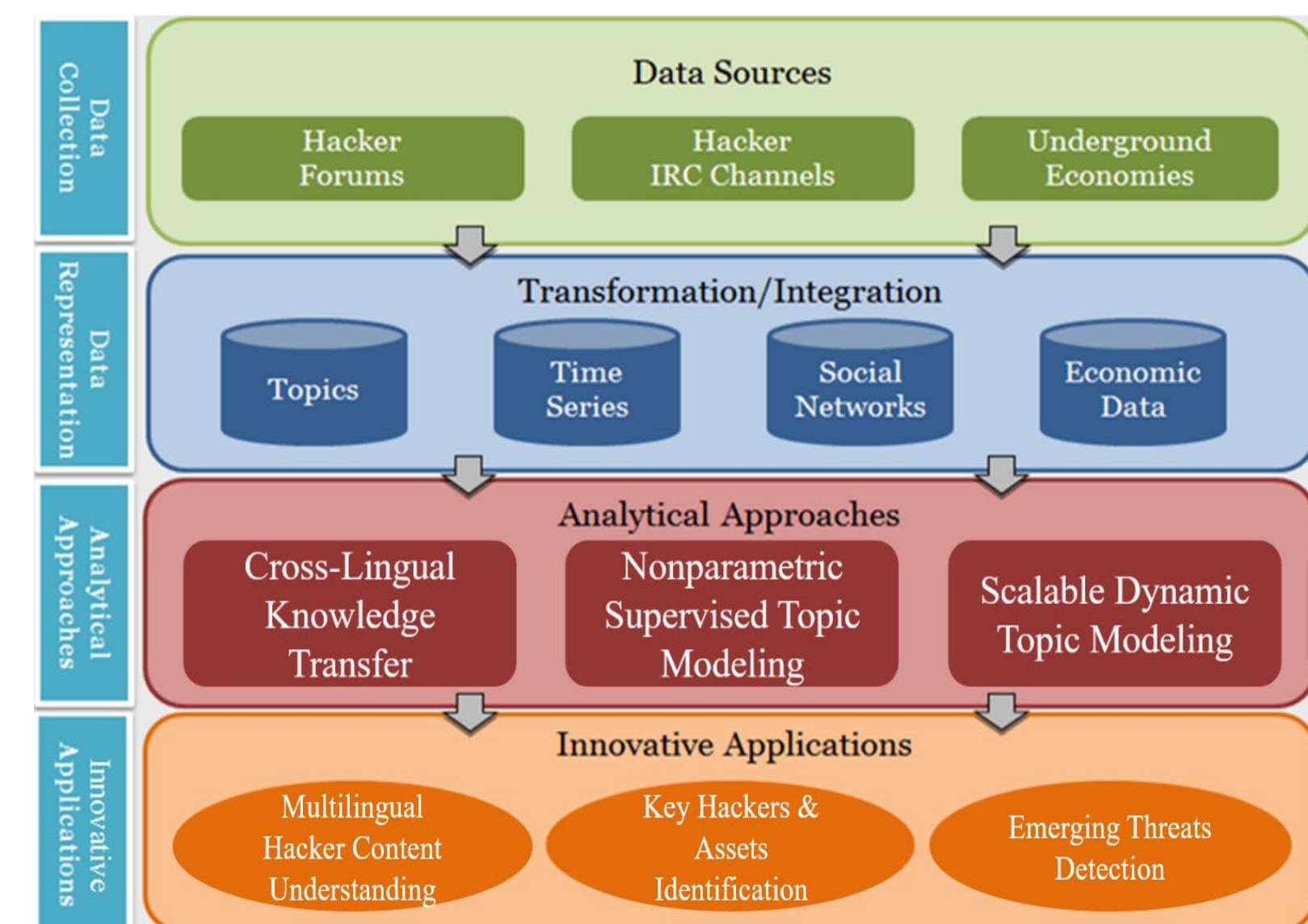


## Research Background

- Conventional Cyber Threat Intelligence (CTI) that analyzes historical attack data (e.g., network logs, honeypots, system logins, and IDS/IPS events) suffers two major limitations:
  (1) Focusing on *reactive* intelligence from previous attacks rather than *proactive* intelligence about threats that have the potential to cause damage (e.g., zero-day attacks) but have not yet been used for cyberattacks,
  (2) Ignoring the threat actors responsible for the attacks, missing the full picture of the hacker ecosystem (e.g., communities, specialties, tools, etc.)
- Hacker communities (e.g., forums, Internet-Relay-Chat (IRC) channels, and underground economies) are of particular interest as they allow hackers to share malicious assets such as hacking tools, malware source code, and hacking tutorials with one another.
- However, several technical difficulties in hacker community data collection and analytics exist:
  (1) The massive volume of data collection
  (2) The heterogeneity and covert nature of the data elements and their often not so obvious linkages
  (3) The ability to comprehend common hacker terms and concepts embedded in communities across multiple regions

## Key Challenges to be Addressed

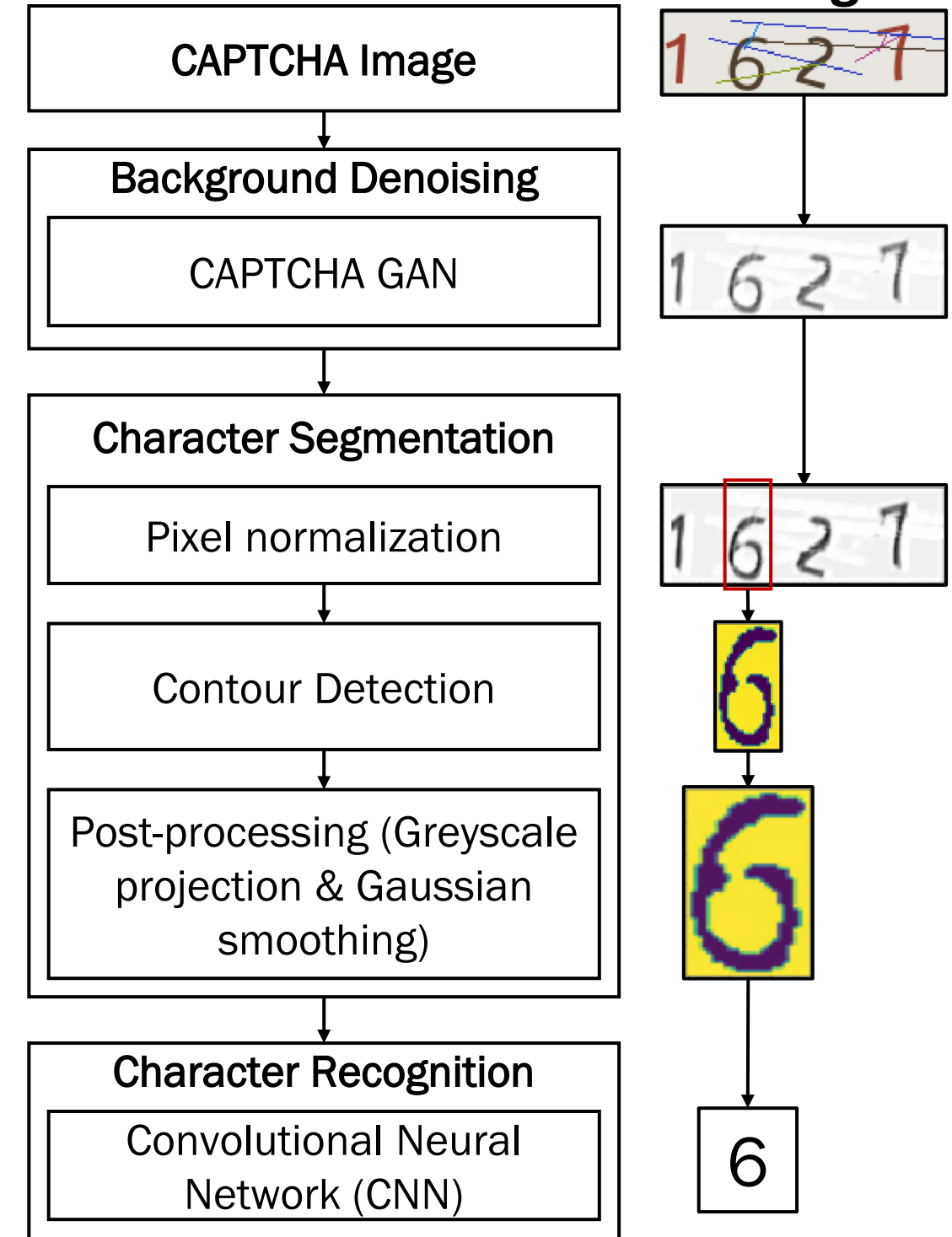This project aims to address two research goals:
1) to advance current capabilities for scalable identification, collection, and analysis of international hacker community contents
2) to make contributions to the cybersecurity community by developing new big data techniques that could enable security researchers to conduct analyses on hacker content and within other related domains.

To address these research goals, this project features four major research projects:
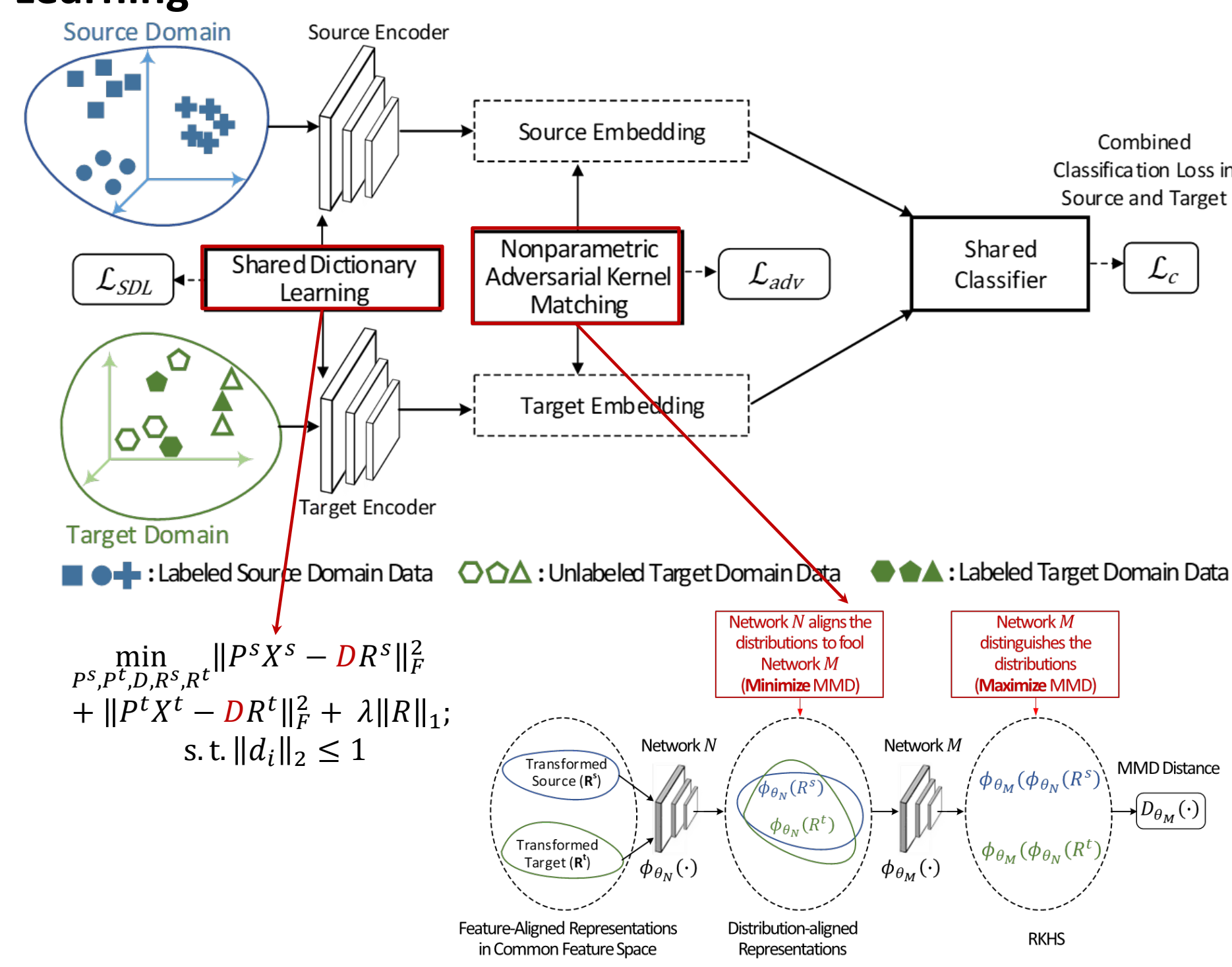1. to develop a large, comprehensive, and longitudinal testbed of international online hacker community contents, including forums, IRCs, underground economies, and other emerging hacker assets
2. to develop an automated method for understanding multilingual hacker community contents to facilitate multilingual CTI analytics
3. to develop a scalable framework for the identification of key hacker assets from the massive underground economy
4. to develop novel dynamic topic modeling capabilities to detect emerging topics from hacker communities, and
5. to address the emerging threat of leaked PII records circulating the hacker community through developing deep learning techniques for privacy risk assessment and a web portal for the general public to get educated and take preventive measures.

## Scientific Impact

- Methodological contributions in hacker content analysis to assist cybersecurity practitioners in studying hacker communication and developing proactive CTI from critical hacker communities
  - Exploration of hacker languages, topics, assets, and threats across different international hacker communities
- A rich, large-scale, longitudinal, open-source collection of diverse hacker community field data to support timely and data-driven cybersecurity and criminology exploration and hypothesis testing
  - Increasing the number of individuals working on CTI research, and subsequently advancing important and complex hacker community CTI inquiries
- Advancing computational knowledge and capabilities in deep transfer learning, deep generative modeling, supervised topic modeling, dynamic topic modeling, neural variational inference, and numerous other important domains.
- Applicability to other big data problems, especially in the social media context:
  - Community identification and collection procedures: applicable to other studies seeking to collect deep web data, especially from sources that employ anti-crawling mechanisms (common among DarkNet markets.)
  - Developed methodologies: applicable to analyzing multilingual community data in other contexts, especially for other legitimate electronic networks-of-practice where tutorials, source code, and other tools may be shared.
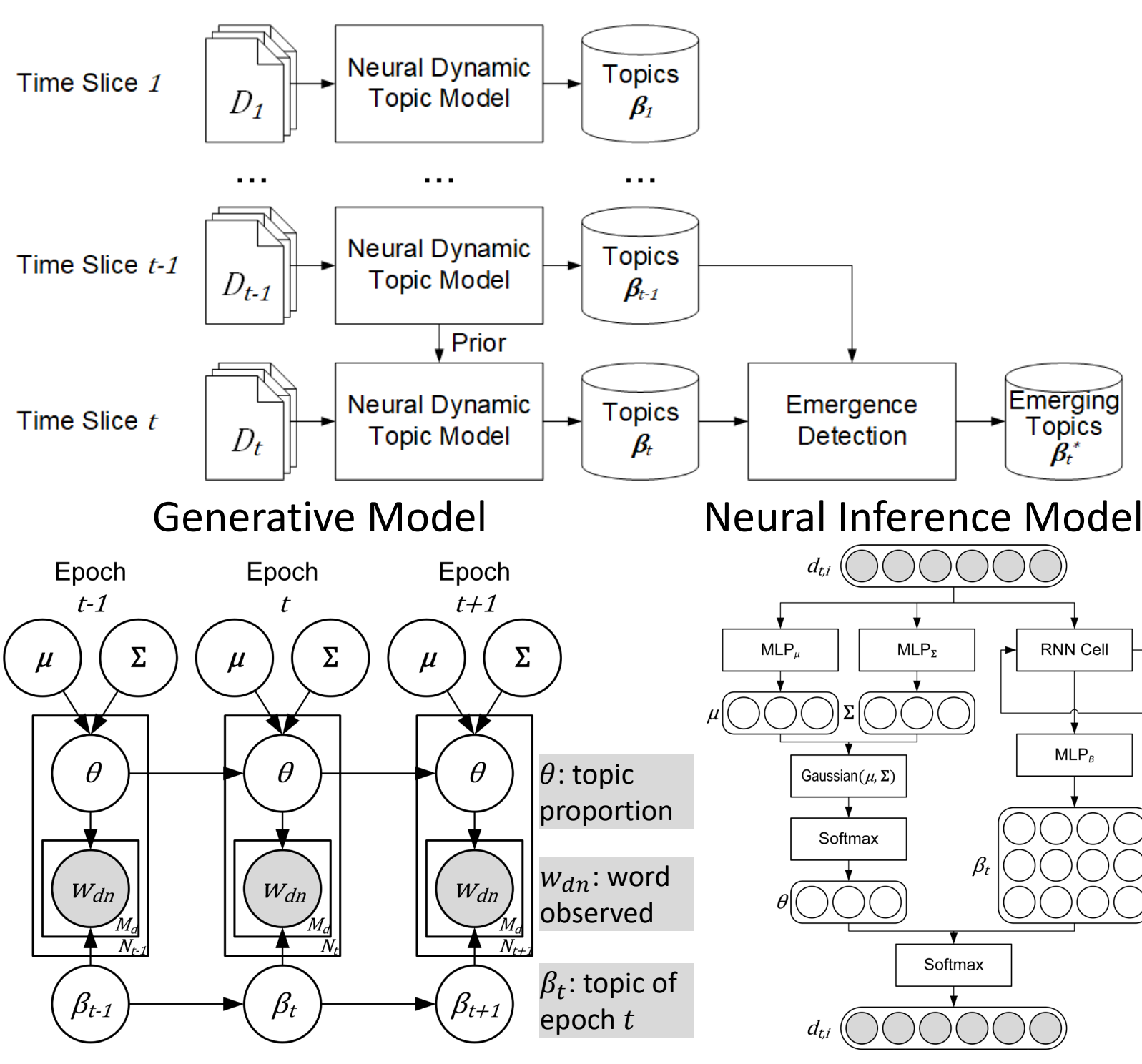
## Selected Technical Approaches:

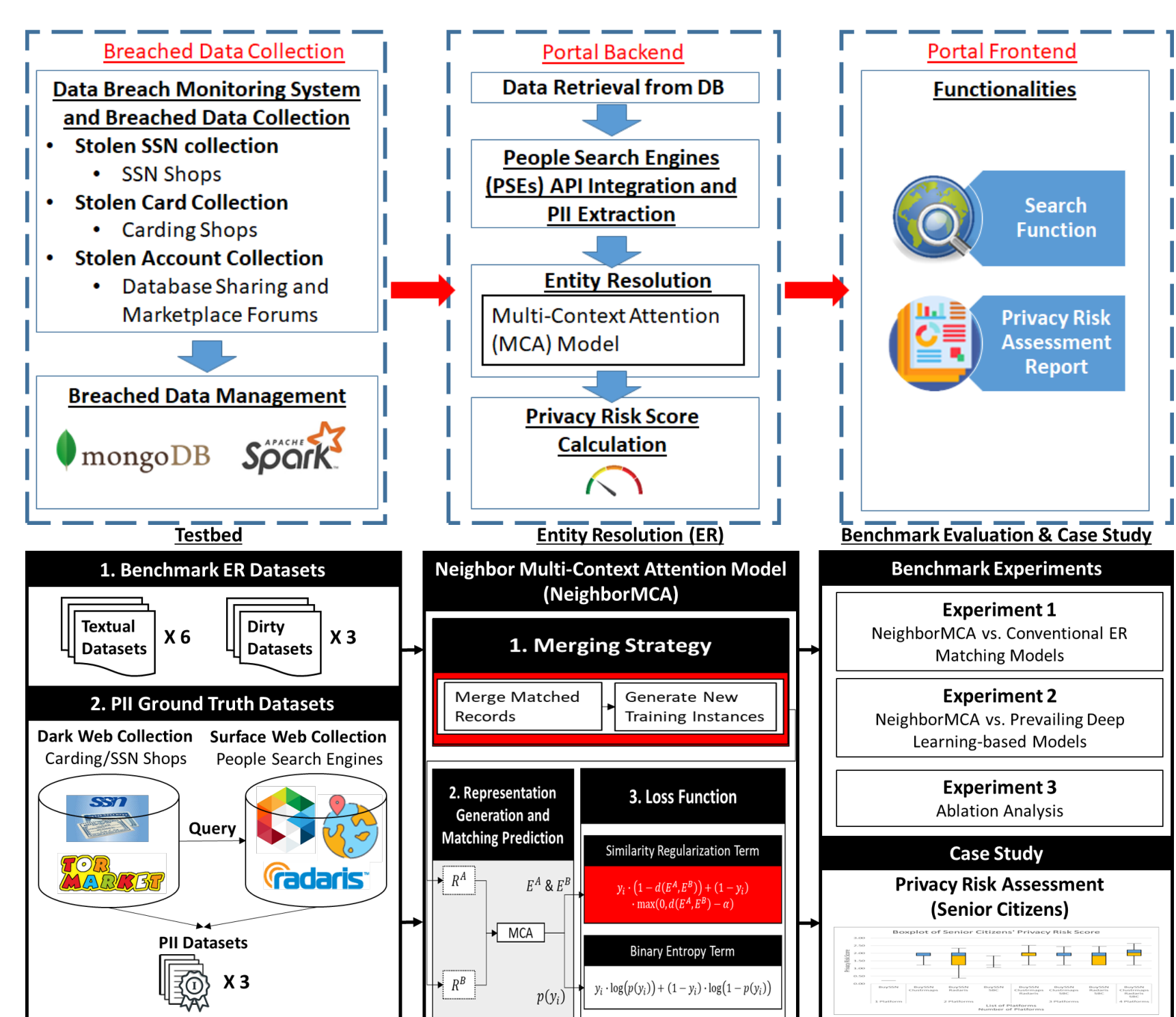### 1. Hacker Community Collection: Counteracting Dark Web CAPTCHA with Generative Adversarial Learning



### 2. Cross-lingual Analytics: Heterogeneous Domain Adaptation with Adversarial Neural Representation Learning



$$\min_{P^S, P^t, D, R^s, R^t} \|P^s X^s - DR^s\|_F^2 + \|P^t X^t - DR^t\|_2^2 + \lambda \|R^t\|_1;$$
$$s.t. \|d_i\|_2 \leq 1$$

### 3. Emerging Topic/Threat Detection: Dynamic Topic Modeling with Neural Variational Inference



Generative Model    Neural Inference Model

### 4. Hacker Community Breached Data Analytics: AZSecure Privacy Analytics



## Broader Impact on Society

### Practitioner communities
- Government agencies (e.g., US Postal Inspection Service and Canada Revenue Agency)
- Law enforcement (e.g., US Secret Service)
- Criminologists (e.g., The Society for the Policing of Cyberspace (POLCYB))
- Intelligence analysts (e.g., National Cyber-Forensics & Training Alliance (NCFTA))

### Dissemination of Collected Data
- https://www.azsecure-data.org/, our sustainable platform for retaining and sharing our selected data collections and analytical approaches (facilitated by our NSF-funded DIBBs project)

### Dissemination of Research Findings
- NCFTA and POLCYB: our primary domain advisory and end-user feedback groups
- Periodical knowledge sharing with practitioner community facilitated by NCFTA and POLCYB

## Broader Impact on Education and Outreach

### Education at University of Arizona
- Top-ranked MS in Cybersecurity program: integrating selected results into the Cyber Warfare class
- AZSecure Scholarship-for-Service (SFS) program: preparing master's and Ph.D. students with cybersecurity research experience as per the independent studies requirement of the SFS curriculum

### Education at University of Georgia
- Area of Emphasis on Information Security: integrating research findings from selected projects into the Cyber Threat Intelligence class (~300 students over three years)

### Outreach
- IEEE ICDM Deep Learning for Cyber Threat Intelligence (DL-CTI) Workshop
- ACM KDD Workshop on AI-enabled Cybersecurity Analytics
- IEEE TDSC Special Issue on Explainable AI for Cyber Threat Intelligence (XAI-CTI)
- IEEE ICDM Machine Learning for Cybersecurity (MLC) Workshop

## Broader Impact and Broader Participation

### Student Training
- University of Arizona: six graduate students and three undergraduate students trained in hacker community collection and analysis; contributed to the development of tools and tutorials
- University of Georgia: two undergraduate students trained in data collection

### Professional Development:
- Data collection (e.g., DBMS, modern programming languages)
- Infrastructure design (e.g., virtual machine management) and development
- Project management and team membership
- Assisting faculty with the development of practical, hands-on cybersecurity-related activities
- Contributing their knowledge about cutting-edge research to enhance the course curriculum
- Opportunity to interact with faculty and others in the field

Award ID#: CNS-1936370