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

PI: Dr. Hsinchun Chen, Regents' Professor, ACM/IEEE Fellow, U. of Arizona (UA), AI Lab Director

Co-PI: Dr. Weifeng Li, U. of Georgia (UGA)

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

There is a growing perspective among both cybersecurity researchers and professionals that understanding the human behaviors behind malicious online acts such as cybercrime and cyberterrorism is an important objective in gaining a better understanding of the cyber threat landscape. Hacker communities, which commonly exist as 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. While such communities make excellent candidates for research, the technical difficulties in data collection and analytics: the massive volume of data collection, the heterogeneity and covert nature of the data elements and their often not so obvious linkages, and the ability to comprehend common hacker terms and concepts embedded in communities across multiple regions present non-trivial challenges to researchers.

This project aims to address two research goals: (1) to advance current capabilities for large-scale identification, collection, and analysis of hacker communities, and (2) to make contributions to cybersecurity research by developing new big data and text analytics techniques that could enable cybersecurity researchers to conduct analyses on hacker community content and within other related domains.
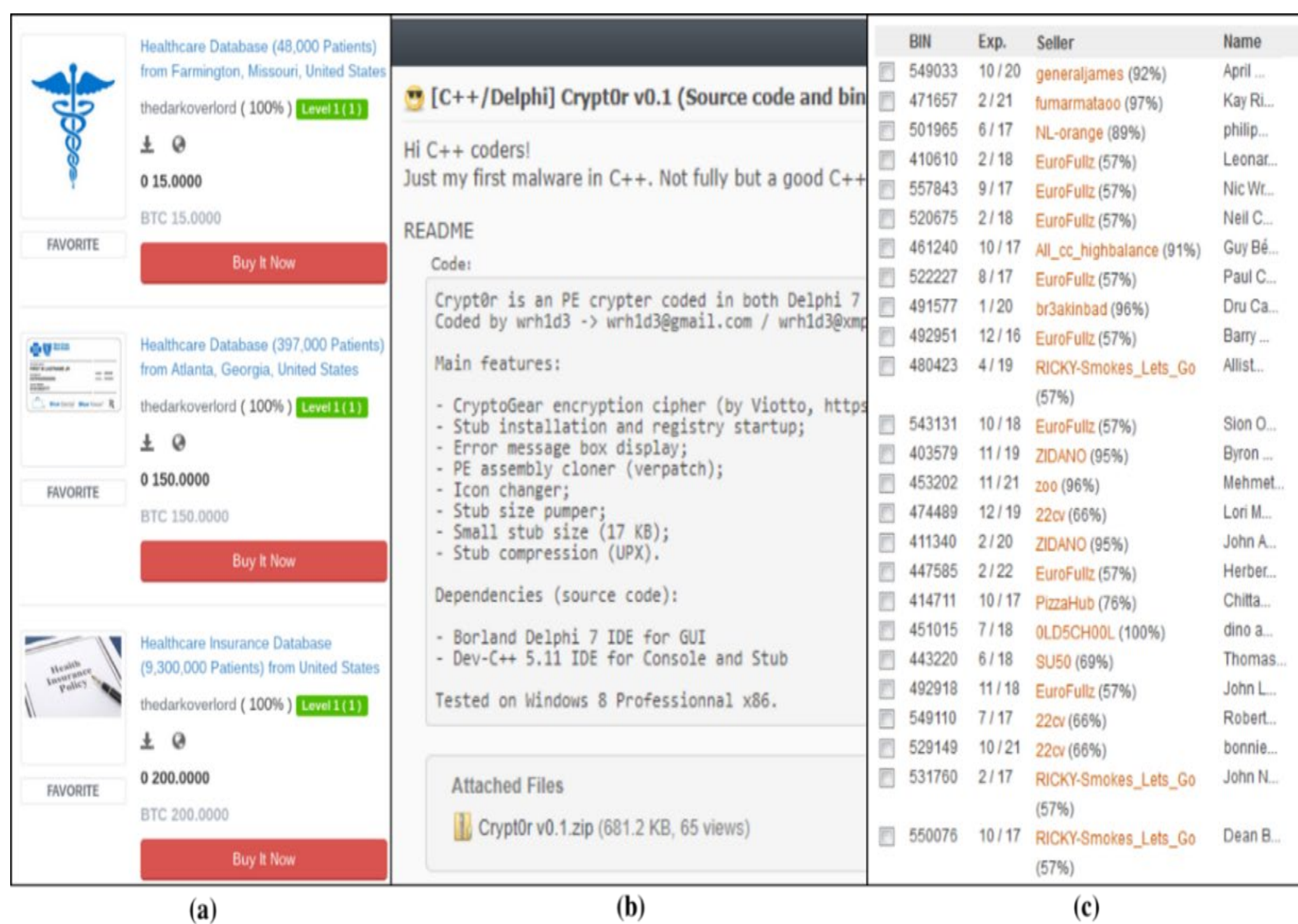


Figure 1. Potential Threats in Selected Hacker Community Content: (a) stolen health data on a DNM; (b) crypter, a technology for ransomware on forums; (c) credit/debit cards and SSNs for sale in a carding shop.
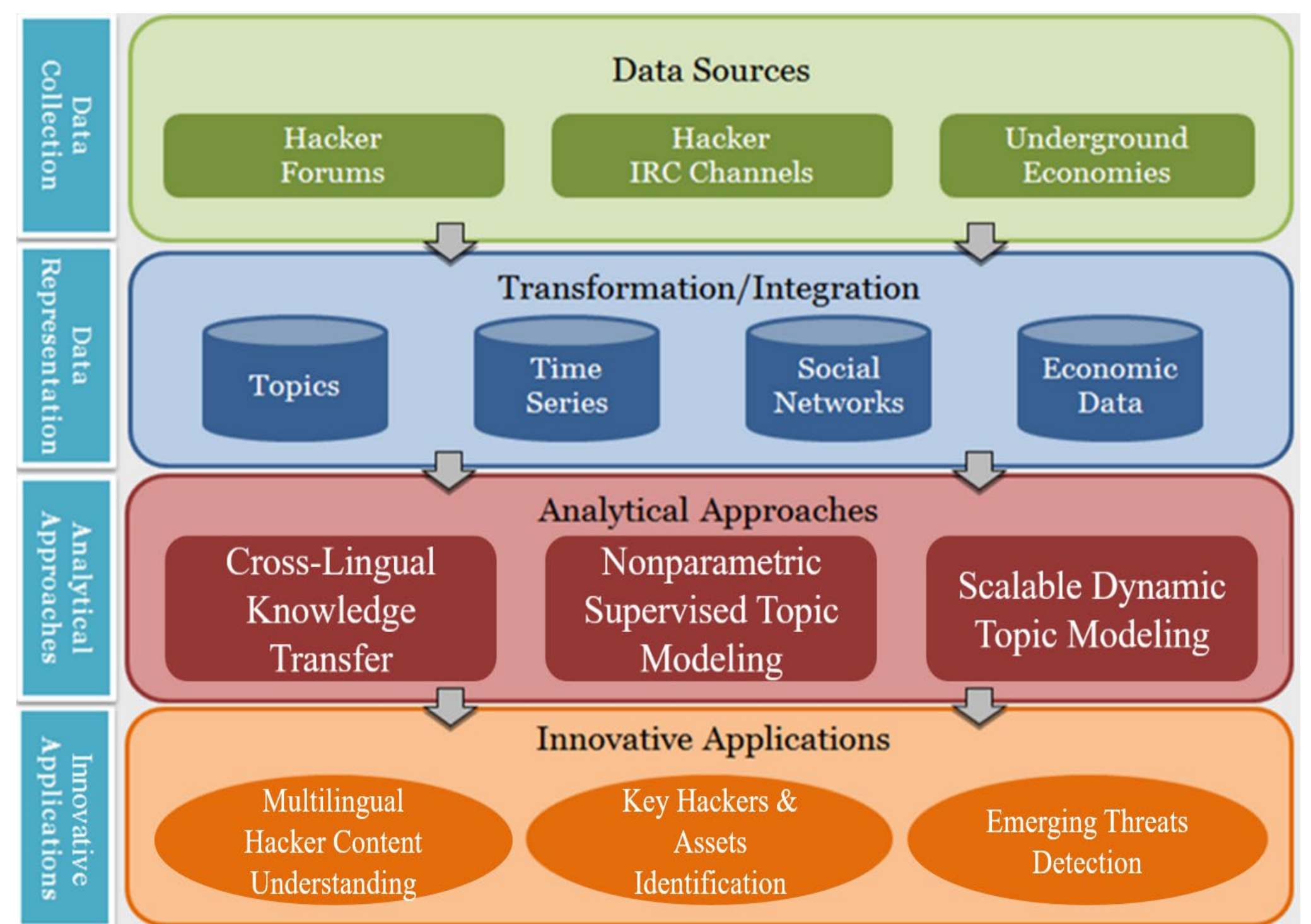


Figure 2. Proposed Framework of Scientific Impacts

## Intellectual Merit:

Developing a large, comprehensive and longitudinal testbed of significant international online hacker communities:

- Web crawlers equipped with counter anti-crawling mechanisms and IRC listeners capable of connecting to hidden services

Understanding multilingual hacker community content to facilitate multilingual analytics:

- An adversarial deep representation approach to learning the language-invariant representations from content in English and non-English hacker communities

Identifying the key hacker assets such as malware and stolen data dumps:

- A nonparametric supervised topic modeling method for examining customer reviews of hacker assets to infer their quality

Modeling temporal attributes of hacker community content and finding evidence of potential emerging threats:

- A scalable dynamic topic modeling technique designed for incorporating expert knowledge of hacker communities and converging efficiently

## Societal Impact:

- National Cyber-Forensics & Training Alliance (NCFTA) and The Society for the Policing of Cyberspace (POLCYB), serving as primary domain advisory and end user feedback groups, as well as liaisons to disseminate knowledge to others in the cybersecurity community
- Planning for extensive collaboration with the greater ISI communities and NSF SFS communities to evaluate and provide feedback on our proposed analytical approaches

## Integration into Education:

- Significant research roles for master's and Ph.D. students in the UA's cybersecurity-based AZSecure Scholarship-for-Service (SFS) program
- NSA designated Center of Academic Education in Cyber Defense (CAE-CD) courses at UA (150+ students)
- UA's online MS in Cybersecurity program
- UGA's undergraduate area of emphasis in Information Security (~40 students)

## Dissemination:

- Cybersecurity conferences and academic outlets such as IEEE ISI, Women in Cybersecurity (WiCyS) (800+ participants), and NSF SaTC PI meeting (400 participants)
- NCFTA's and POLCYB's Annual Conferences each draw 1,000+ attendees from law enforcement, government, and industry
- NSF-funded AZSecure DIBBs-ISI repository designed to share dozens of multi-million record security datasets and tools