# Data Infrastructure Building Blocks (DIBBs) for Intelligence and Security Informatics Research and Community

- Identifying and Building the ISI Community
- They Want Our Data! Developing the Data Science Testbed
- Promoting Collaboration, Education, and Community
- Demonstration, Plans, and Future Growth
- Join Us!

PI: **Dr. Hsinchun Chen**, University of Arizona. Co-PIs: **Dr. Mark Patton** and **Cathy Larson**, University of Arizona. Project Partners: **Dr. Ahmed Abbasi**, University of Virginia; **Dr. Paul Hu**, University of Utah; **Dr. Bhavani Thurasingham**, University of Texas at Dallas;    **Dr. Chris Yang**, Drexel University.

# Building the ISI Community

- ISI = "**Intelligence and Security Informatics**"

- ISI = R&D for **advanced information technologies** and systems for **national and international security-related applications**
  - Through an integrated technological, organizational, and policy-based approach

- ISI community first brought together at the inaugural *Intelligence and Security Informatics International Conference*
  - Initially founded by Dr. Hsinchun Chen (UA), sponsored by the **National Science Foundation** and held in Tucson in 2003
  - Now under the auspices of IEEE with multiple international events and a community of 4,000+ scholars in US, Europe, and Asia
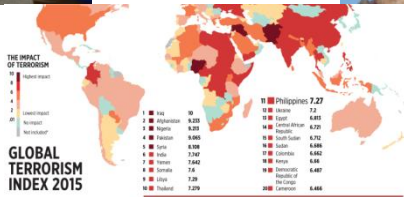
# Why ISI?



Since the year 2000, there have been over 73,000 terrorist attacks, killing more than 170,000 people.

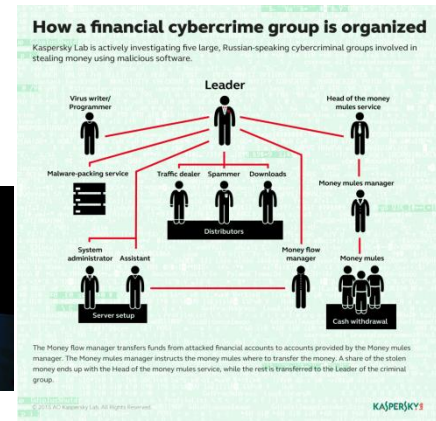The global economic impact of terrorism reached US$89.6 billion in 2015.

**-- Global Terrorism Index, 2015**

Cyber crime costs are estimated to be $400 billion a year, with estimates going as high as $2 trillion by 2019.

Large banks, retailers, and federal agencies make the headlllines, but all businesses are at risk.
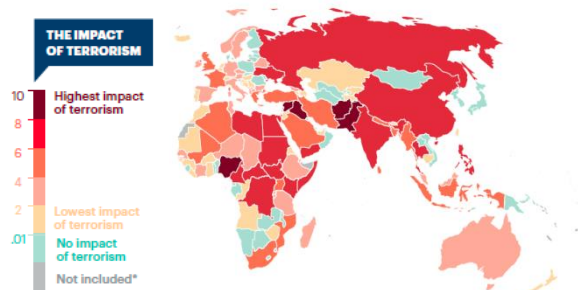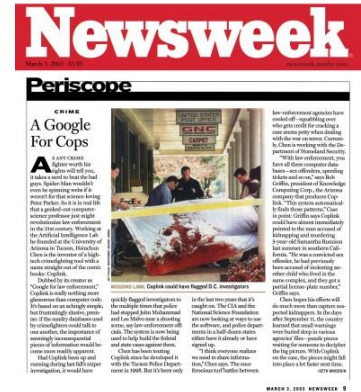
**-- Forbes, US edition, 1/17/2016**

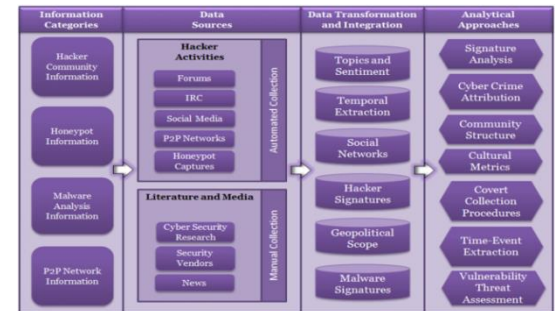# A Short Course in AI Lab **ISI Data History**:



- **"A Google for Cops"**: With funding from NIH, NSF, etc., the Lab develops **COPLINK** to support **information sharing, analysis, and visualization of law enforcement data**; COPLINK was commercialized and purchased by IBM and is now used by 3,500+ agencies. New approaches in data/text/web mining in COPLINK lead to Dark Web:



*From: Global Terrorism Index, 2015*

- The **Dark Web:** terrorists' and cyber criminals' use of the Internet. Generates **massive amounts of in-demand data: forums, websites, videos, blogs, etc.** From Dark Web to Hacker Web (cybersecurity) via **new approaches to spidering, social network analysis, authorship analysis, video analysis, sentiment analysis**, and more.

- The **HackerWeb** project uses **security big data analytics** to answer important questions about **hacker behaviors, markets, community structure, communication, and cultural differences**. (Research framework at right.)



**BIG DATA + SECURITY NEEDS + RESEARCHERS + NSF = DATA INFRASTRUCTURE BUILDING BLOCKS FOR INTELLIGENCE AND SECURITY INFORMATICS**

## Emphasize difficult to collect, novel data sources:

- **Dark Web Data** and **Geo Web Data**: the most **in-demand** and **easiest/fastest to make available**
  - Originally requested by over 800 researchers, students, analysts, others, in computer science, international relations, intelligence
  - **39M+ postings** from general, extremist, and terrorist forums in their original languages

- **PhishMonger**: A **new tool** developed as part of DIBBs to collect **live phishing websites** to support the study and analysis of these ephemeral exploits[1]
  - Currently the most requested data: **227 GB** downloaded in the recent period (September-November)

[1] For system details, see Dobolyi & Abbasi, 2016. "PhishMonger: A Free and Open Source Public Archive of Real-World Phishing Websites," *in* Proceedings of the IEEE ISI 2016 Conference.

# Build, promote, and support the ISI community

- **Data Challenges**: Organized around a question or series of questions which can only be answered by the use of the AZSecure-data set provided.
  - Encourages **new analytical techniques** and **promotes the use of the portal**
  - Successful data challenges @ **IEEE's ICDM** (International Conference on Data Mining), **ASONAM** (Advances in Social Network Analysis and Mining), **KDD** (ACM's Knowledge Discovery and Data Mining conference), reaching a potential audience of more than 3,000 attendees

- **New Workshops:** Two **Security Big Data and Analytics Sharing Workshops** will be held in 2017 and 2018 to get user input to build on the DIBBs-supported platform

- **Quarterly Updates:** Sent to list of **2,300+ potential ISI users**, compiled from ISI proceedings and directories of institutes, centers, labs, and other security related organizations

6

# Demonstration, Plans, and Future Growth

## AZSecure-data.org
### Intelligence and Security Informatics Data Sets

Data Infrastructure Building Blocks for ISI. A Project of the University of Arizona (NSF #ACI-1443019), Drexel University, University of Virginia, University of Texas at Dallas, and University of Utah

Home | About | Get Data | ISI Events | Citing Data | Papers | Policies

### Get Data

DARK WEB FORUMS | GEOWEB FORUMS | PHISHING SITES | TWITTER DATA

OTHER FORUMS | OTHER DATA

The AZSecure-data repository currently provides access to **Web forums, Internet phishing websites, Twitter data, and other data.** Most files are available to download from the "Get Data" buttons above; other files can be requested through the project manager. To request access to restricted data, send email to ailab@eller.arizona.edu and state which data set you would like to use and the purpose for which it will be used; also provide complete contact information including name, affiliation and mailing address, email, and telephone number.

The **Dark Web forums** w[...]
project on the study of i[...]
and Russian. **GeoWeb** ge[...]
assessing country risk. G[...]
Afghanistan, Algeria, Egy[...]
and Yemen. The collectio[...]
members. **Other forums**[...]
page. Additional informa[...]
accompanying each data[...]
which may then be opene[...]
click on the name of the[...]
PROVIDED "AS IS."

The **Internet phishing w**[...]
as downloadable zip or r[...]

**Websites**

- **Patriot, Militia, Hate and Linked Websites** - Collected by the Artificial Intelligence Lab, Management Information Systems Department, University of Arizona, the Patriot, Militia, Hate and Linked Websites collection presented here contains 74 websites belonging to groups identified by the Southern Poverty Law Center in 2009 as belonging to groups promoting extreme social perspectives. The collection also contains 123 additional websites linked to by the initial set of websites. The full list of websites in this collection is in the ReadMe.txt file. Due to the size of this collection, it has been divided into 20 portions to make downloading easier. Each bundle of websites contains the ReadMe.txt and About.pdf files.

ReadMe.txt
About.pdf
PatriotMilitiaHate.zip (394MB)
PatriotMilitiaHate2.zip (1GB)
PatriotMilitiaHate3.zip (2.3GB)
PatriotMilitiaHate4.zip (22MB)
PatriotMilitiaHate5.zip (9.3GB)
PatriotMilitiaHate6.zip (585MB)
PatriotMilitiaHate7.zip (227MB)
PatriotMilitiaHate8.zip (4GB)
PatriotMilitiaHate9.zip (562MB)
PatriotMilitiaHate10.zip (1.1GB)
PatriotMilitiaHate11.zip (865MB)
PatriotMilitiaHate12.zip (957MB)
PatriotMilitiaHate13.zip (584MB)
PatriotMilitiaHate14.zip (577MB)
PatriotMilitiaHate15.zip (434MB)
PatriotMilitiaHate16.zip (491MB)
PatriotMilitiaHate17.zip (781MB)
PatriotMilitiaHate18.zip (425MB)
PatriotMilitiaHate19.zip (833MB)
PatriotMilitiaHate20.zip (500MB)

**Social sharing links**

**Data types currently available**

- Dark Web and Geo Web forums
- Phishing sites
- Twitter data
- Other forums
- Other data:
  - Malware
  - Network traffic
  - News
  - Web chat
  - Websites

**ReadMe** file explains provenance, organization, and use

**About** file provides additional info about the published paper

**Data** available in easily downloaded zip files

**EXAMPLE: Patriot, Militia, Hate, & Linked Websites**

7

## *Coming Soon!*

1) **Hacker Community Data** - a selection of more than **185,000,000 records** from **79 platforms** to **study the international hacker community and behavior**, including:

- Darknet Marketplaces
- Hacker Forums
- IRC Channels
- Bitcoin transactions
- Carding shops



Figure 1. Hacker forum posting where a member shares malicious credential stealing tool.



Figure 2. Dumps of compromised credentials in a hacker forum posting.

2) Dunning **ISI researchers** to **share their data**!

- Cold-calling and otherwise contacting our community of 4,000+ researchers to prepare and share their data via the AZSecure-data portal.

3) New **repository platform** to provide a stable environment with increased searching, browsing, and analytical functionality

# *Join Us!*



- 270 members in LinkedIn ISI group
- Building up Facebook and Twitter activity
- Sending newsletter quarterly

*FIND OUT MORE:*

http://www.azsecure-data.org/

http://ai.arizona.edu/research/dibbs