



Supporting C2 Research and Evaluation: An Infrastructure and its Potential Impact

James Law, Ph.D. and Marion Ceruti, Ph.D.

Space and Naval Warfare Systems Center Pacific (SSC Pacific)

16th ICCRTS, Quebec City, Canada
21-22 June, 2011

Topic 6: Experimentation, Metrics, and Analysis

Why?

- ▼ The incredible, time-consuming, difficulty of getting test data.
 - Searching for data...
 - Organizing data...
 - Storing data...
 - Reusing data...
 - Maintaining data...
- ▼ Bottleneck for:
 - Development.
 - Testing.
 - Integration.
 - Fielding.
 - Maintenance.

C4ISR Empirical Studies & Experiments

- ▼ Four (4) empirical studies.
- ▼ One (1) controlled experiment.
- ▼ **Our Thesis:**
 - A common body of test data that can be used for varying experimentation will provide researchers and practitioners with greater power to evaluate
 - algorithms,
 - architectures,
 - implementations.
- **Our Proposal:**
 - A common repository of unrestricted, open source data sets and supporting infrastructure.

Challenges for Research and Evaluation

- ▼ Replicability of findings.
- ▼ Supporting aggregation of findings.
- ▼ Reducing costs.
- ▼ Obtaining representative operational profiles.
- ▼ Isolating the effects of individual factors.

Survey of Data Sources

- [1] P. Baldi, P. Frasconi, and P. Smyth, *Data sets for the Book Modeling the Internet and the Web*, <http://ibook.ics.uci.edu/datasets.html>
- [2] Carnegie Mellon University Language Technologies Institute, *The ClueWeb09 Dataset*, <http://boston.lti.cs.cmu.edu/Data/clueweb09/>
- [3] CDISS, *Database of Terrorist Incidents, 1940 – 1999*, http://www.cdiss.org/pages/Programmes/Revolutionary_Warfare_Counter_Insurgency/Publications.asp
- ...
- [13] H. Chen and C. Yang, eds. *Terrorism Informatics: Knowledge Management and Data Mining for Homeland Security*, Springer, New York, NY. 2008. <http://ai.arizona.edu/research/terror/>
- [14] DARPA, *Message Understanding: Evaluation and Conference: Proceedings of the 3rd-6th DARPA Workshops*, Morgan Kaufman Publishers, 1996.
- ...
- [16] J. O. Engene, *Terrorism in Western Europe: Events Data (TWEED)*, <http://folk.uib.no/sspje/tweed.htm>
- [17] Fondation pour Recherche Stratégique, *Base De Données Sur Les Actes Terroristes*, Paris, France, <https://bdt.frstrategie.org/>
- [18] A. Frank and A. Asuncion, *UCI Machine Learning Repository*, University of California, Irvine, School of Information and Computer Science, 2010. <http://archive.ics.uci.edu/ml/>
- [19] A. Frank and A. Asuncion, *UCI Machine Learning Repository*, University of California, Irvine, School of Information and Computer Science, 2010.
- ...
- [23] L. Jasny, C.S. Marcum, and C.T. Butts, “improv: A Dataset of Disaster Response Microevents,” 2009, R package version 0.05. Contact Carter Butts (buttsc@uci.edu) for availability.
- ...
- [25] J. Leskovec, *The Stanford Large Network Dataset Collection*, <http://snap.stanford.edu/data/>
- [26] J. Leskovec, L. Backstrom and J. Kleinberg, *MemeTracker*, <http://www.memetracker.org/data.html>
- ...
- [30] B. Mendelsohn, *Global Terrorism Resource Database* <http://people.haverford.edu/bmendels/> Mendelsohn also maintains a list of terrorism data sources: http://people.haverford.edu/bmendels/terror_attacks
- ...
- [32] National Counterterrorism Center (NCTC), *Worldwide Incidents Tracking System (WITS)*, <https://wits.nctc.gov/>
- [33] National Institute of Justice (NIJ), *Terrorism Databases for Analysis*, <http://www.ojp.usdoj.gov/nij/topics/crime/terrorism/databases.htm>
- [34] National Institute of Standards and Technology (NIST), *Information Access Division (IAD), Introduction to Information Extraction, MUC Archive Site*, http://www-nlpir.nist.gov/related_projects/muc/index.html
- [35] M.E. Nissen, “Enterprise Command, Control, and Design: Bridging C2 Practice and CT Research,” *The International C2 Journal*, Vol. 1, No. 1, 2007. Defense Technical Information Center, <http://handle.dtic.mil/100.2/ADA486841>
- ...
- [37] RAND Corporation, *Database of Worldwide Terrorism Incidents*, <http://www.rand.org/nsrd/projects/terrorism-incidents/index.html>
- ...
- [40] SEMVAST Project, *Scientific Evaluation Methods for Visual Analytics Science and Technology*, <http://www.cs.umd.edu/hcil/semvast/>
- [41] SEMVAST Project, *Visual Analytics Benchmarks Repository*, <http://hcil.cs.umd.edu/localphp/hcil/vast/archive/index.php>
- [42] University of Chicago, *Joint Threat Anticipation Center (JTAC)*, <http://jtac.uchicago.edu/resourceTerror.shtml>
- [43] University of Michigan, *Inter-University Consortium for Political and Social Research (ICPSR)*, <http://www.icpsr.umich.edu/icpsrweb/ICPSR/>
- [44] University of Michigan, *Terrorism and Preparedness Data Resource Center (TPDRC)*, <http://www.icpsr.umich.edu/icpsrweb/TPDRC/>
- [45] University of Pennsylvania, *The Linguistic Data Consortium*, <http://www ldc.upenn.edu/>
- [46] U.S. Department of Homeland Security, *National Consortium for the Study of Terrorism and Responses to Terrorism (START)*, <http://www.start.umd.edu/start/>
- ...

Survey of Data Sources

Terrorism Data Sets	Sources	Size	Formats
<i>Dark Web Terrorism Research Project</i> [13]	Web crawls	Large	Text, video
<i>Haverford Database of Terrorist Acts</i> [30]	News, propaganda	Medium	Text, pdf
<i>(CDISS), Database of Terrorist Incidents, 1940 – 1999</i> [3]	CDISS	Large	Text
<i>French Database of Terrorist Acts</i> [17]	French agencies	Medium	Text
<i>Message Understanding Conference (MUC) Archives</i> [34]	Synthetic	Small	Text
<i>RAND Database of Worldwide Terrorism Incidents</i> [11]	RAND	Medium	Text
<i>Scientific Evaluation Methods for Visual Analytics Science and Technology (SEMVAST)</i> [40, 41]	Synthetic	Medium	Formatted text
<i>National Consortium for the Study of Terrorism and Responses to Terrorism (START)</i> [46]	START	Large	Databases, text
<i>Terrorism and Preparedness Data Resource Center</i> [44]	START	Medium	Text
<i>NIJ Terrorism Databases for Analysis</i> [33]	NIJ	Small	Databases
<i>Terrorism in Western Europe: Events Data (TWEED)</i> [16]	TWEED	Medium	Database
<i>Joint Threat Anticipation Center (JTAC)</i> [42]	JTAC	Medium	Database
<i>Worldwide Incident Tracking System (WITS)</i> [32]	NCTC	Medium	Database
<i>Ali Baba Data Set</i> [22]	Synthetic	Small	Text
<i>Counter-Terror Social Network Analysis and Intent Recognition (CT-SNAIR)</i> [48]	Govt agencies	Medium	Text
<i>World-Trade Center Event Sequence Data</i> [23]	UC Irvine [23]	Small	R-archive

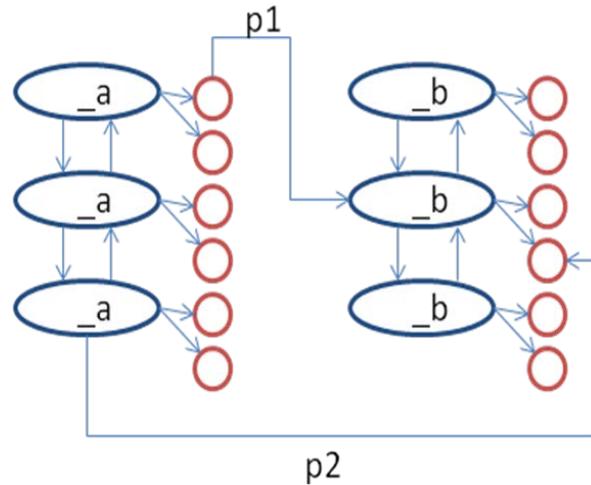
Survey of Data Sources

Non-Terrorism Data Sets	Sources	Size	Formats
<i>Inter-University Consortium for Political and Social Research, ICPSR [43]</i>	Research studies	Many small sets	Various
<i>Modeling the Internet and the Web: Probabilistic Methods and Algorithms [1]</i>	Research studies	Small	Formatted text
<i>ClueWeb09 Dataset [2]</i>	Web crawls	Large	Text
<i>Linguistic Data Consortium (LDC) [45]</i>	LDC	Small	Text
<i>Stanford large network dataset collection</i>	Various	Large	Formatted text
<i>MemeTracker data</i>	News, blogs	Large	Formatted text
<i>UC Irvine Machine Learning Repository</i>	Research studies	Small	Formatted text

Required Infrastructure

Our proposed repository contains:

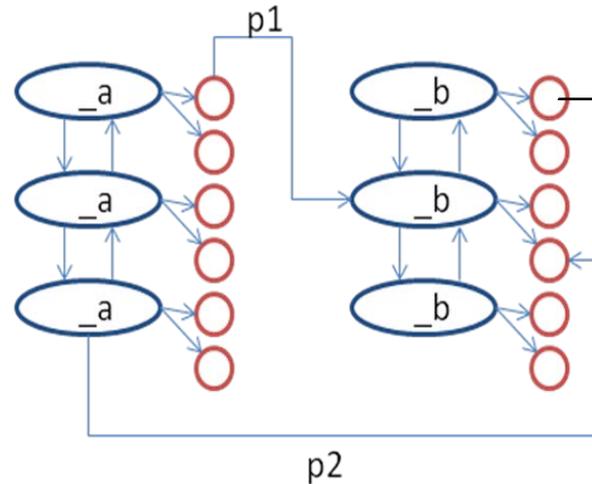
- Data Sets (RDF)



- Documentation and Supporting Tools

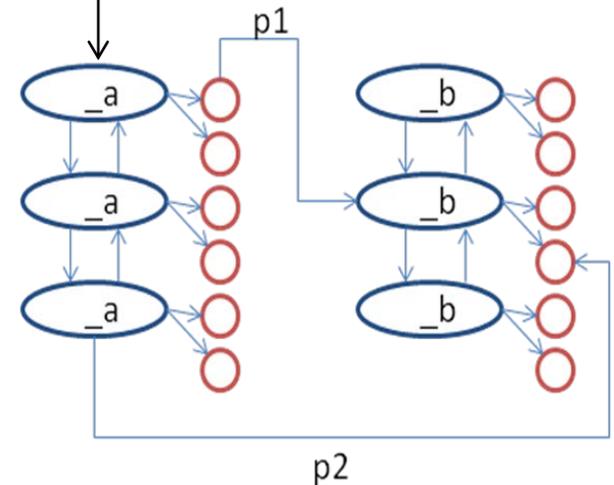
Data Sets

- Synthetic
- Mashups
- Real



Benefits of mashups:

- Can be created relatively quickly with scripts.
- Small data sets can be selectively enlarged.
- Prospective algorithms and systems can be tested for sensitivity to node set size and edge density.
- Data sets can be readily shared with collaborators.



Documentation and Supporting Tools

- A complete specification of formats (including standard formats)
- Date-time group (e.g. day, hour, minute, sec.) of data-element or data-set
- Latitude and longitude or other location specification of data collection
- Data pedigree, e.g. the origin of the data (e.g. sensor ID)
- The means (e.g. sensors, observations, etc.) used to collect the data sets
- The provenance or processes involved in producing or delivering the data
- The algorithms that were used to integrate and fuse the data
- Standard metrics such as data-set size
- Transformation scripts

Concerns

- **Threats to validity**

- **Internal** - how well the experiment controls the dependent variables and, therefore, how strong the causal relationships in the experiment can be trusted.
- **External** - extent to which the results of the experiment can be generalized to other experimental subjects.

- **Extension and integration**

- How will existing data sets be maintained?
- How will new data sets be integrated?
- How will the infrastructure be adapted as systems and practice change?

Repository State

Current:

- Several example mashups.
 - Expect to be available for unrestricted release soon.
- Example Python scripts for transforming to RDF.

Future:

- XML and RDF validation support.
 - Visualization support.
- Feedback from internal users.

jim.law@navy.mil

SSC *PACIFIC*
on Point
and at the Center of C4ISR