



THE CENTER FOR ADVANCED STUDIES  
IN SCIENCE AND TECHNOLOGY POLICY

# FROM DATA MINING TO COMPUTATIONAL SOCIAL SCIENCE FOR COUNTERTERRORISM



**K. A. TAIPALE**

EXECUTIVE DIRECTOR, CENTER FOR ADVANCED STUDIES  
SENIOR FELLOW, WORLD POLICY INSTITUTE  
ADJUNCT PROFESSOR OF LAW, NYLS

PRESENTED TO:

THE COMMITTEE ON TECHNICAL AND PRIVACY DIMENSIONS OF  
INFORMATION FOR TERRORISM PREVENTION AND OTHER NATIONAL GOALS  
THE NATIONAL ACADEMIES, APRIL 27, 2006

# Presentation: Policy issues in applying computational social science to CT

- Definitions
  - What is DM in this context?
- Methodologies
  - How is DM used in CT?
- Policy concerns
  - What is the problem?
    - Preemption
    - Privacy, 4th A, DP, &c.
    - Legal significance
- Observations
  - What is (where is?) the solution?
- Conclusion(s)
  - Where do we go from here?
- *References*

## Threshold problem: Definitions

- What is data mining? (policy vs. technical definition)
- Aggregation (access) and automated analysis (sorting) = “data mining” for policy purposes
  - Data matching
  - Link analysis
  - Pattern generation
  - Pattern matching
- TAPAC and DM reporting amendment (Feingold-Leahy)
  - Any query to third party data or secondary use of gov data
- What is “data mining” for CT? Computational social science used to diagnose social pathologies, i.e. hostile, covert networks
- Legal/policy issue: use of DM as predicate for further action

# Computational social science in CT

- Mathematics, statistics, economics, political science, cultural anthropology, sociology, psychology, psychiatry, neuroscience, and computer science (modeling, visualization, and simulation)
- Counterterrorism as complex adaptive system - emergence, self-organization, small worlds
- Development and use of nonlinear, nondeterministic theories/models of complex human phenomena (at all scales)
- Move from modeling the individual (AI) to the group (network)
- How do groups organize and act covertly in concert? (~LoneW)

# Automated analysis methodology

- Subject based (find out *more about* entity) (“psychology”)
  - Confirm identity
  - Disambiguate identity (~ semantic reconciliation)
  - Evaluate “reputation” to establish trust or suspicion
- Pattern based (find *other related* or *like* entities) (“sociology”)
  - Observed (MO) (**too few examples?**)
  - Hypothesized (red teamed) (**social science**) (**force CM sigs**)
  - “Data mined” (KDD) (**identify avoidance sweet spot**)
  - Forward and inverse problem (model/spot) (“loose-thread”)
- Consequence of “match” (use/inference) (CI for decision-making)
  - Investigative (reasonableness of predicate)
  - Decisional/adjudicatory consequences (fairness of process)
  - Watch lists ... (consequence of “soft triggers”)

## Propositional v. relational DM-ing

- Propositional DM (associative/correlations) (homogenous databases/unrelated subjects)
  - “Terrorists order pizza with credit cards” (VV)
  - “My TiVo Thinks I’m Gay” (WSJ)
- Relational DM (connections/interactions) (heterogeneous databases/related subjects)
  - Visualization
  - Social network theory
  - Communications
  - Organization
  - Funding
- Bottom Line: not all shared frequent flyer accounts belong to terrorists, but sharing FF# with a known terrorist is a pretty good predicate for investigation (cf. CASE STUDY)

# Outcomes for DM in CT

- Allocate resources (preventative measures)
  - Risk pools (risk management based)
  - Systems and layered defense
  - Recognize and plan for failures
    - Over- and under-inclusions
    - System failures
    - Unintended consequences of success
- System self-awareness (situational awareness)
  - How does/should an information mediated society maintain SA
  - Monitor/control for abuse/misuse (~ can be turned on itself)
  - Discovery and learning (vs. normalization and central limit theorem)
- Predicate for further action - thus, depends on the confidence interval for particular use

## Caveat: the engineering example

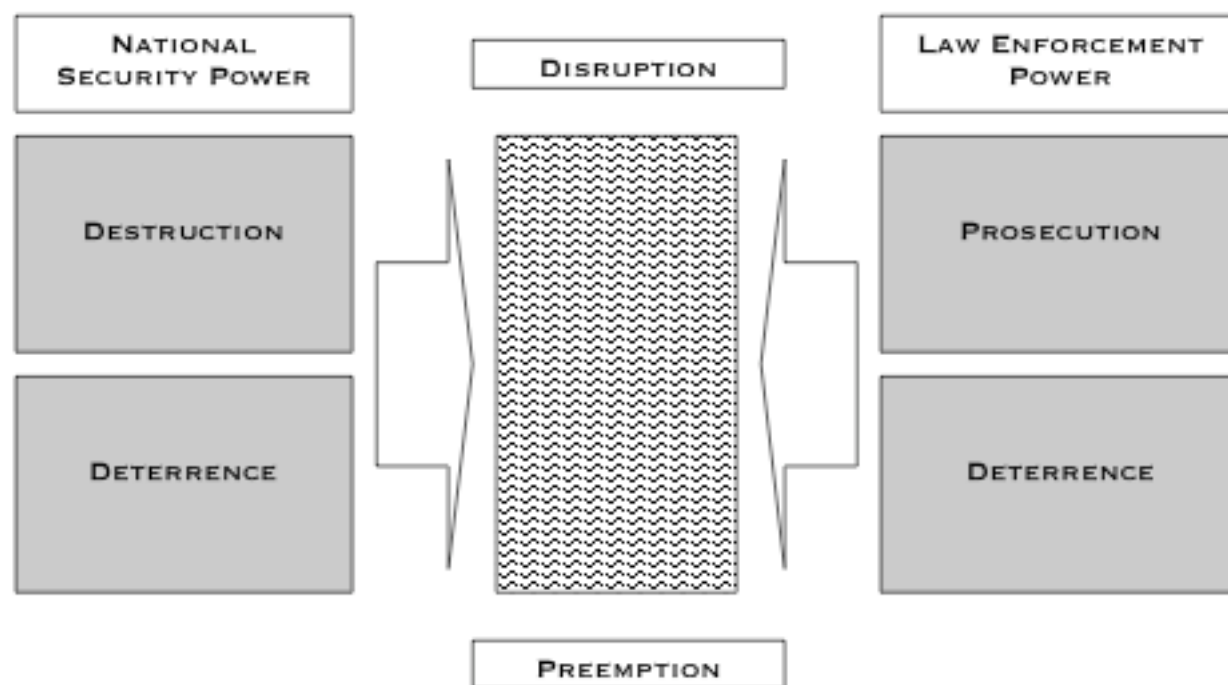
- Setting ratio of false positives to false negatives
  - Default system bias / decision heuristic
  - Variable to threat environment and particular circumstance
  - Cost to security (“real” vs. “theater”) (cost of secondary screening)
  - Cost to functionality (friction, limited degrees of freedom)
  - Probability-based security will have failures
- Engineering models and catastrophic collapse
  - Don’t trust models, you over-engineer (costs)
  - Trust your models, you “push the envelope” (collapse)
- DM is only one tool
  - Need “defense in depth”
  - Need to build in error correction and elegant failure modes



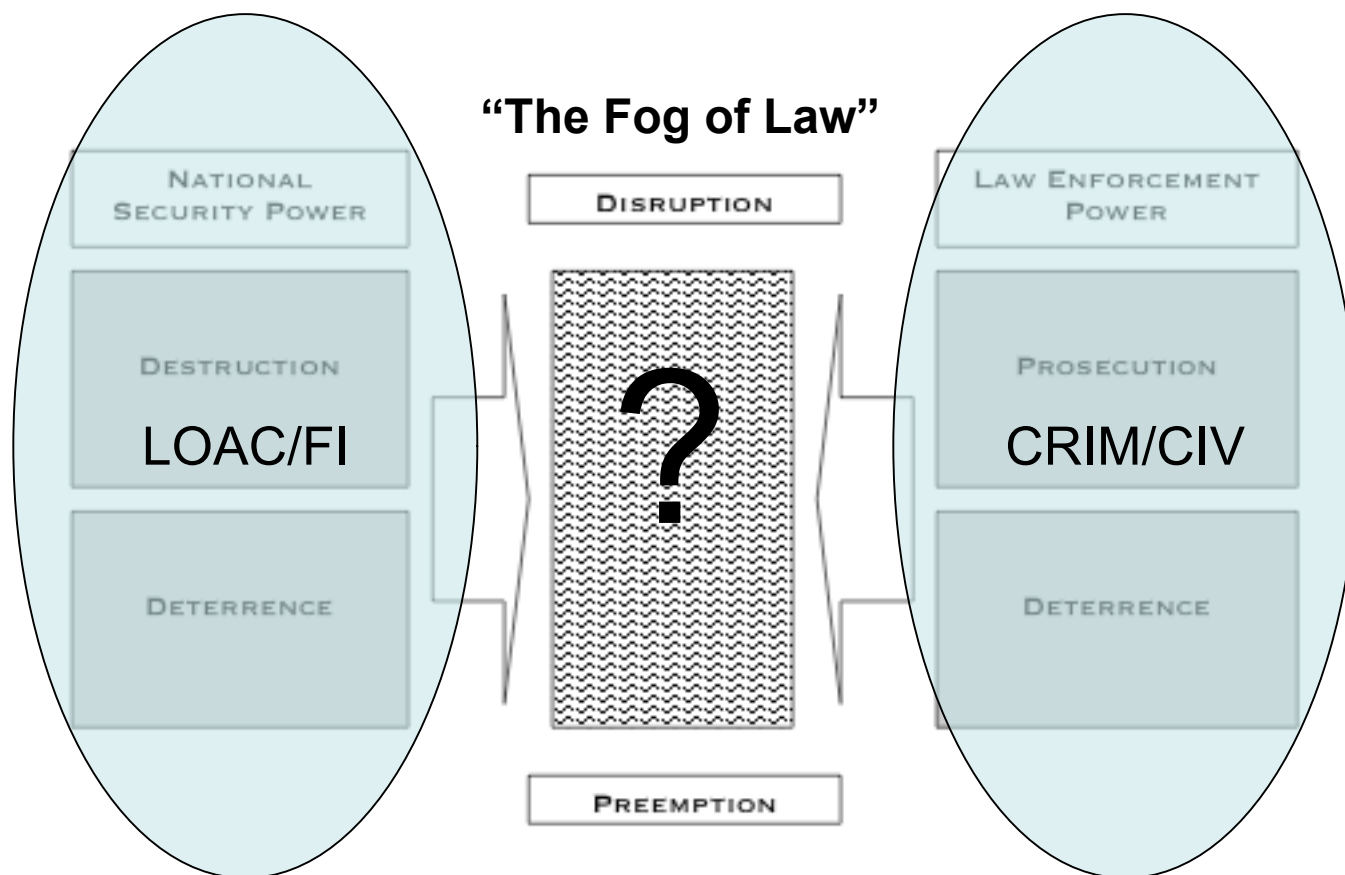
## Policy concerns - preemption

- Reactive law enforcement vs. preemptive national security strategies (counter nation-threatening potentials)
  - Disparate doctrinal regimes for foreign v domestic threats
  - “Truth” finding vs. disruption (action and countermeasures)
- *Beccarian model* - punishment and deterrence of deviant individuals by sanctioning commission of criminal acts
- *Foucauldian model* - general social compliance through ubiquitous preventative surveillance and through system constraints
- Changing role of security services - from policing (arrest and prosecution) to *risk management* through surveillance, exchange of information, auditing, communication, and classification.

# Converging missions to prevent potential catastrophic (nation-threatening) outcomes



# Need for policy reconciliation: doctrinal conflict and vacuum



## Policy concerns - “privacy”

- Privacy *interests* (The Private “I”)
  - Secrecy (place or things) 4th Amendment
  - Autonomy (behavior) 5th Amendment
  - Anonymity (identity and speech) 1st Amendment
- Privacy protection
  - Communication -- *Katz v US* (1967) (people not places)
  - Information - *Whalen v. Roe* (1977) (secrecy and autonomy)
  - Third party exception “shared” -- *US v Miller* (1976) (cancelled checks) *Smith v Maryland* (1979) (telecom “attributes”)
  - Identity -- Cf. *Talley, McIntire, NAACP w/ Gilmore, Hiibel*
  - Compare particular request v. general search / general warrant
- “Privacy” is generally a data access not data mining issue (?)

## 4th Amendment

- “The right of the people to be secure in their persons, houses, papers, and effects, against *unreasonable searches* and seizures, shall not be violated, and no *Warrants* shall issue, but upon *probable cause*, supported by Oath or affirmation, and *particularly describing* the place to be searched, and the persons or things to be seized.” 4th Amendment, U.S. Constitution.
- Two clauses
  - No “unreasonable searches” (protects citizens)
  - Warrants “upon probably cause” (protects officials)

## 4th Amendment - [un]reasonable search

- The “prohibition on unreasonable searches is not accorded more weight than the permission to conduct reasonable searches”
- “Reasonableness” (“totality of circumstances”)
  - Predicate for action (probative value not probabilistic nature is key)
  - Alternatives (balance of interests)
  - Consequences (intrusion vs. state interest)
  - [Opportunity for error correction?]
- No requirement for notice or hearing
- If reasonable at the time then error cost fall on false positives and only consequence of unreasonableness is exclusionary rule

## 4th Amendment - [probable cause]+[warrant]

- No universal requirement for warrant
  - Plain view / incident to arrest / stop and frisk
  - Administrative search (~EPA, OSHA)
  - Fixed checkpoint / Airport search / Courtroom search
  - “Special needs” - primary purpose is not LE (drug testing)
- Pre-FISA (1978) - National security and foreign intelligence exception for both surveillance (*Keith*) and search (i.e., 4th Amendment did not apply)
- In upholding FISA (*Megahey* 1982 EDNY) the court held that the 4th A probable cause requirement is not “fixed or static” but depends on the circumstances (thus, FISA was constitutional even if/when you apply 4th A requirements)

## Profiling or non-particularized suspicion

- Test is probative value - confidence interval for particular use
- Drug courier profiles - upheld
  - Six Sup. Ct. cases upholding - assessed whether specific facts amounted to reasonable suspicion of drug trafficking
  - Requires “explaining the algorithm” (independent v. dependent variables) (ensemble classifiers) (Bayesian) Etc.
  - statistical expert witness to explain model? (~ police hunch)
- Behavioral profiles (hijacker profile - US v. Lopez) - upheld
  - “in effect ... system itself ... acts as informer ...”
  - 6% weapon discovery in targets upheld based on attempt to single out hijackers (94% false positive OK)
  - But cf. Edmond v Indianapolis (drug sweep traffic stops struck down with 9% arrests) (primary purpose LE) (but cf. drunk driving and other special needs)
  - social science expert witness to explain model?
- BL: Need to explain the model (but cf. *validity* with *accuracy*)



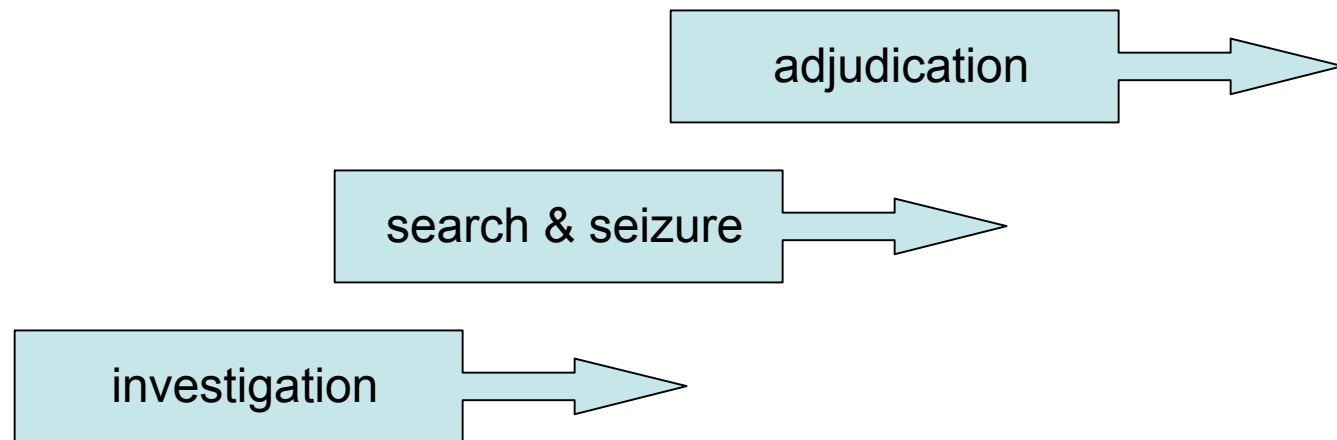
## 5th Amendment - Due Process

- Deprivation of life or liberty (consequences)
- Mathews v. Eldridge (1976) (“three factor balancing test”)
  - Private interest affected by government action
  - Government interest (“including the function involved” and burdens of additional process)
  - Analysis/affect of greater or lesser process (on accuracy, consequence, or error correction/compensation)
- Notice and fair hearing
- Compensation for error

## Additional policy concerns

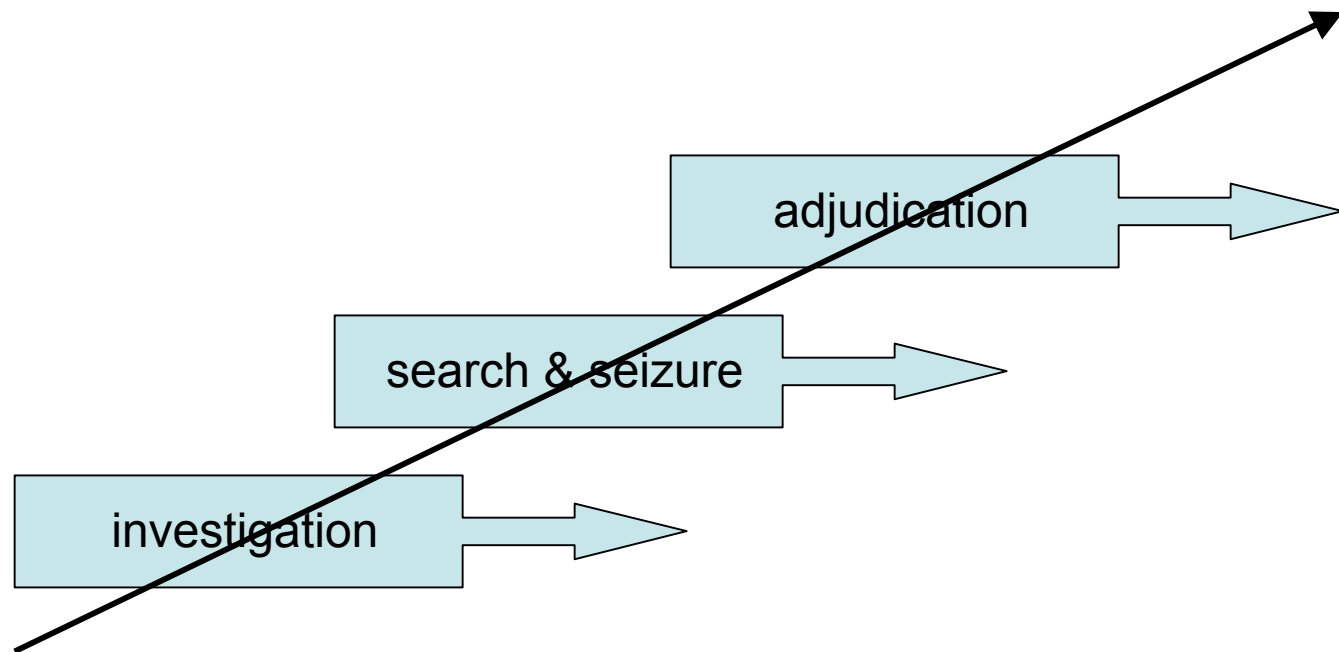
- 1st Amendment “chilling effect” (*but* there is no 1st Amendment right to plot terrorist acts in secret)
  - Material support statutes
  - Need to redefine “clear and present danger” for network effect?
  - See references in “Info War” presentation at <http://seekingsymmetry.info>
- Slippery slope
  - Mission creep
  - Data appetite
  - Bureaucratic expansion
- Selective enforcement (“equal protection” - 14th Amendment)
  - Precursor or low level activities targeted in suspect population
  - Disproportionate impact
- Trusted system problem (see IEEE article in materials)
  - Authorization (totalitarian) (DEFAULT=DENY)
  - Accountability (freedom) (DEFAULT=PERMIT)

## Categories vs. rankings/probabilities



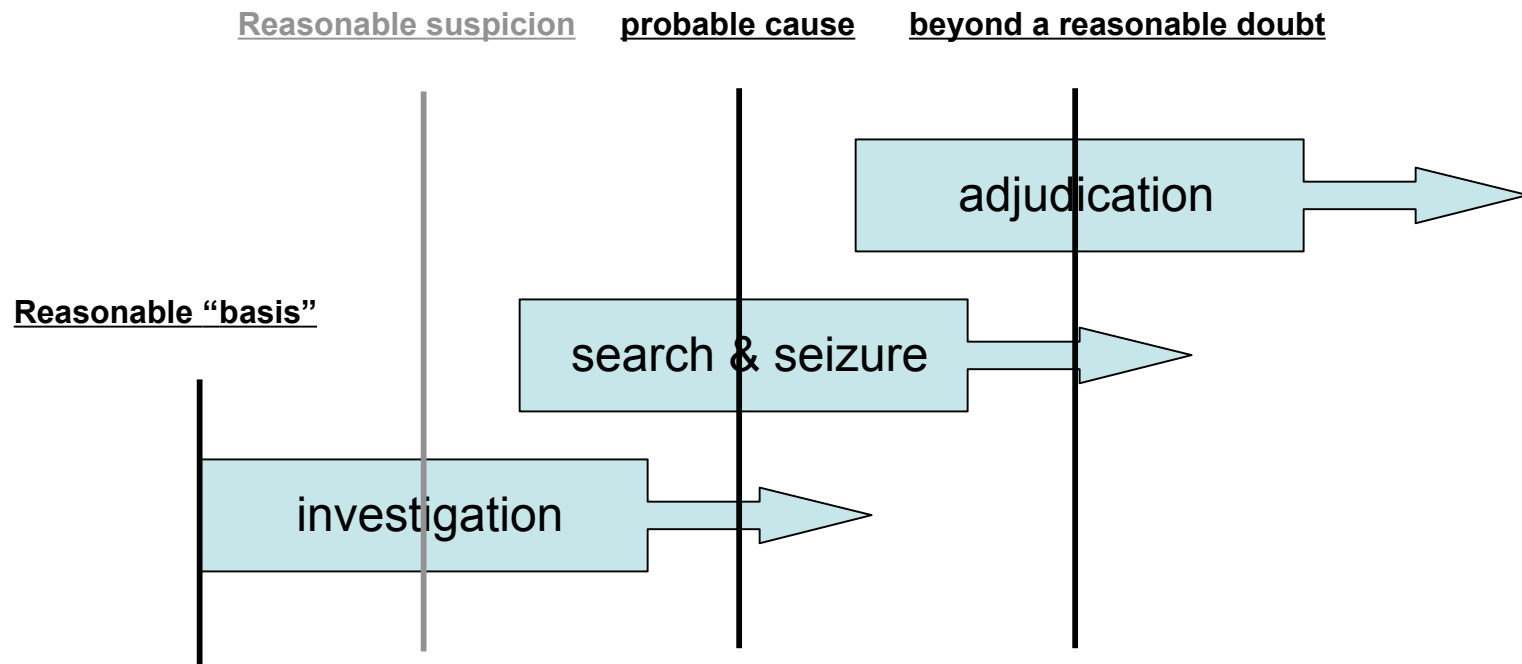
The law (and policy) is categorical (discrete) but reality is continuous ...

## “Real-world” suspicion (~ confidence)



The law (and policy) is categorical (discrete) but reality is continuous ...

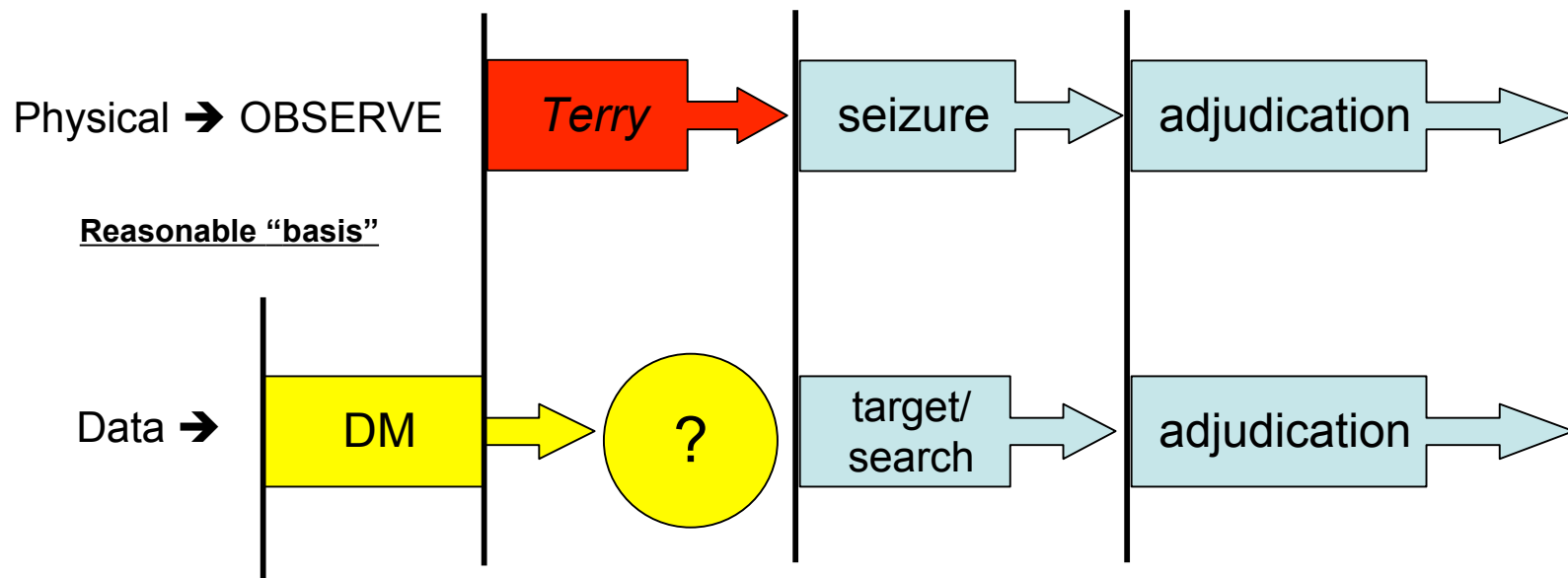
# “Legal” suspicion (~ significance)



The law (and policy) is categorical (discrete) but reality is continuous ...

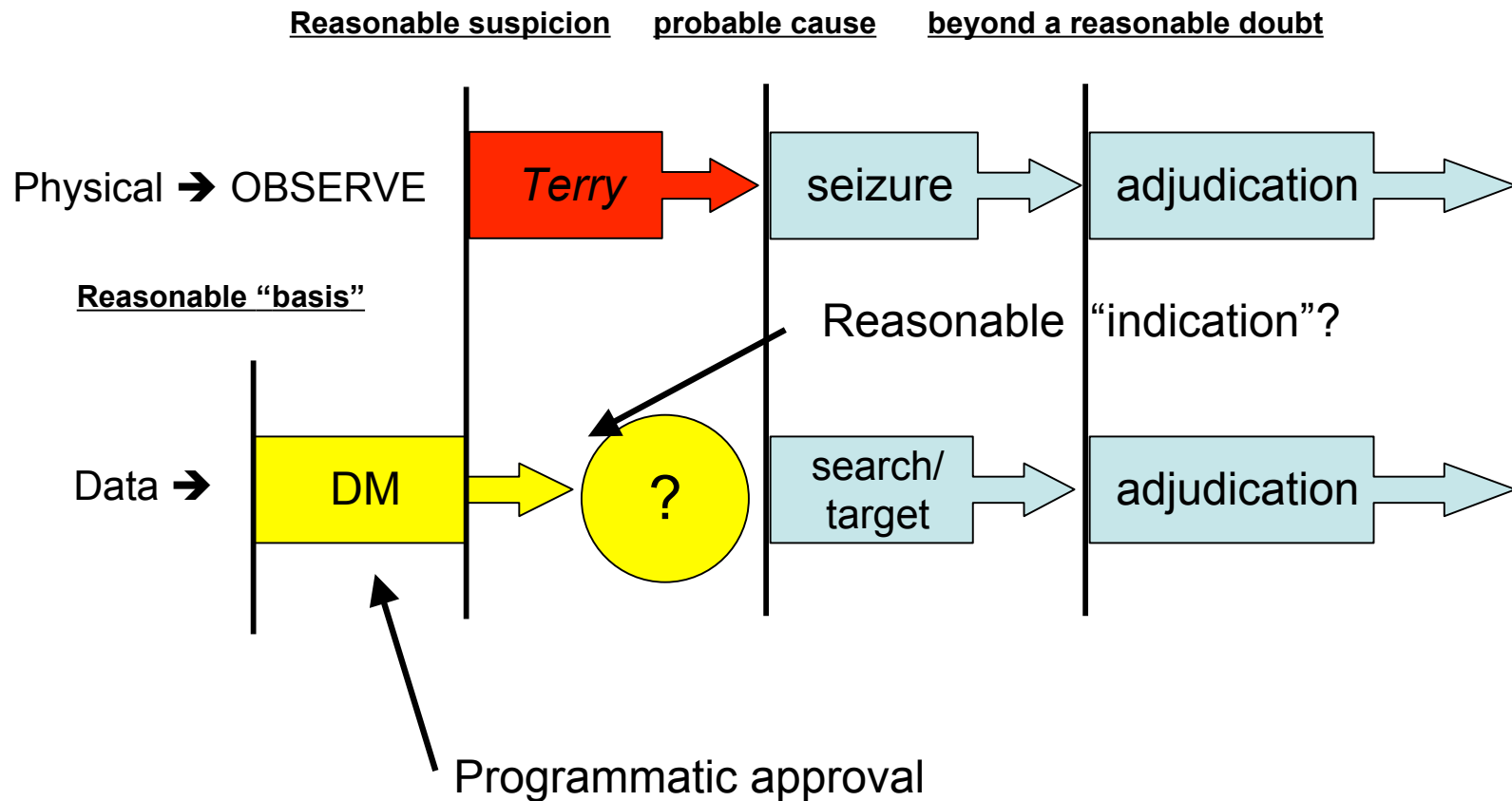
# Legal significance testing in the physical world vs. the dataverse

Reasonable suspicion    probable cause    beyond a reasonable doubt



The law (and policy) is categorical (discrete) but reality is continuous ...

# Required: the *e-Terry* stop



## The e-*Terry* Stop

- Need the equivalent of a “*Terry* stop” for dataveillance and electronic surveillance (*Terry v. Ohio* 1968)
- Reasonable basis to initiate automated screening (programmatic approval related to threat targeted) (scalarity)
- Reasonable suspicion (i.e., articulable or statistical facts) to permit limited follow up investigation to eliminate or establish probable cause
  - Disambiguate identity (cf. CAPPs II) (not you)
  - Invalidate relationship (you, but not terrorist)
    - Additional independent variables in common
    - Additional independent variables not in common
  - VALIDATE > probable cause > (target you)
- Need for explicit authority and procedures for use of computational social science as predicate for investigation



# Calculus of reasonableness

- due process - fairness (Dworkin) (justice as fairness - Rawls)
  - Predicate for action
  - Alternative
  - Consequences (of false negative and false positive)
  - Error correction
- Threat environment (heuristic bias)
  - Need for dynamic assessment
  - Need to dial up (or down) in response to circumstance
  - Authorized purpose driven (demand) not data (supply)
- Dynamic/runtime - reasonableness judged by totality of circumstances at the time of decision-making
- See discussion in materials (“Why Can’t We All Get Along”)

# Requirements

- Develop clear standards - legal and administrative for both R&D and implementation - IC community driven (anticipatory & preemptive)
- Technical features to support due process
  - Mechanisms for intervention and oversight
    - Rule-based processing, selective revelation, “circuit breakers” (build in human or policy intervention points), authentication and audit
  - Design for failure (at systems, applications, and data level)
- Transparency of methodology
  - Statistics/CI's
  - Algorithms/mining/KDD
  - Social science modeling
  - Data reliability, etc.
- Research and “prove” the value of a relationship between the social sciences and computational science in spotting (or predicting) threats

*CAVEAT: Can't burden R&D with success  
or technology/social science with perfection*

- R&D
  - Unrealistic expectations to prove outcomes prior to R&D
  - Opposition to research on the basis that it "*might not work*" is an example of what has been called the "*zero defect*" culture of punishing failure, a policy that stifles bold and creative ideas - "*downright un-American.*" David Ignatius, ***Back in the Safe Zone***, Wash. Post, Aug. 1, 2003, at A:19 (discussing the knee-jerk opposition to a "terrorist futures market")
- Development/deployment
  - Perfection shouldn't be the enemy of the good (or the better)
  - "Does it work" should be judged against existing and/or alternative methods not against perfection

## Conclusion(s)

- Not a balancing act, no fulcrum, need to *maximize both security and liberty* within the constraints imposed by the other
- Technology development is not deterministic, but it is *inevitable*
- *Law alone* offers little security and brittle privacy protection
- Technology *cannot provide either security or privacy*, but:
  - Properly employed - can better allocate security resources
  - Properly designed - enable existing mechanisms (or analogues)
- The use of computational social science for counterterrorism follows the natural historical development of *social immunology*

## DM&CT Background References

ISAT 2002 Study, ***Security with Privacy*** (Dec. 2002)

Paul Rosenzweig, ***Proposals for Implementing the Terrorism Information Awareness System***, Legal Memorandum No. 8, Heritage Foundation (Aug. 2003) 2 Geo. J. L. & Pub. Pol'y 169 (2004).

K. A. Taipale, ***Data Mining and Domestic Security: Connecting the Dots to Make Sense of Data***, 5 Columbia Sci. & Tech. L. Rev. 2 (December 2003).

Department of Defense (DOD) Technology and Privacy Advisory Committee (TAPAC), ***Safeguarding Privacy in the Fight Against Terror***, Final Report (Mar. 2004).

Mary DeRosa, ***Data Mining and Data Analysis for Counterterrorism***, CSIS Press (Mar. 2004).

K. A. Taipale, ***Technology, Security and Privacy: The Fear of Frankenstein, the Myth of Privacy and the Lessons of King Ludd***, 7 Yale J. L. & Tech. 123; 9 Intl. J. Comm. L. & Pol'y 8 (March 2004).

## DM&CT Background References II

Charles Weiss, *The Coming Technology of Knowledge Discovery: A Final Blow to Privacy Protection?* 2004 U. Ill. J. L. Tech. & Pol'y 253 (2004)

ABA SCL&NS, *Cantigny Principles on Technology, Terrorism, and Privacy*, Nat. Sec. Law Report (Feb. 2005).

Daniel J. Steinbock, *Data Matching, Data Mining, and Due Process*, 40 Georgia L. Rev. (2005).

K. A. Taipale, *The Trusted System Problem: Security Envelopes, Statistical Threat Analysis, and the Presumption of Innocence*, IEEE Intelligent Systems, Trends and Controversies (Sept./Oct. 2005).

Fred H. Cate, *Legal Standards for Data Mining*, in Emergent Information Technologies and Enabling Policies for Counter-Terrorism (Robert Popp and John Yen, eds., IEEE Press 2006).

K. A. Taipale, *Whispering Wires and Warrantless Wiretaps: Data Mining and Foreign Intelligence Surveillance*, NYU Rev. L. & Security (Spring 2006).

## About the Center for Advanced Studies

- The Center is an independent, non-partisan research firm focused on information, technology, and national security policy.
- The Center seeks to inform and influence national and international policy- and decision-makers in both the public and private sectors by providing sound, objective analysis, insight, and advice; in particular by identifying and articulating issues that lie at the intersection of technologically-enabled change (both opportunities and challenges) and existing practice in public policy, law, and technology development.
- The Center has ongoing research projects in: *Law Enforcement and National Security in the Information Age; Telecommunications and Cybersecurity Policy; Information Operations and Information Warfare;* among others.
- *More info and contact at [www.advancedstudies.org](http://www.advancedstudies.org).*

*More: <http://taipale.info/>*  
*</end>*