# DataTags, Data-Handling Policy Spaces, and the Tags Language

Michael Bar-Sinai
Computer Science Dept.
Ben-Gurion University of the Negev
Be'er-Sheva, Israel

Latanya Sweeney
Data Privacy Lab
Harv
Ca

Mercè Crosas
Institute for Quantitative Social Science

p
Mich
mbarsina
@n

| Tag Type | Description | Security Features | Access Credentials |
|---|---|---|---|
| Blue | Public | Clear storage, Clear transmit | Open |
| Green | Controlled public | Clear storage, Clear transmit | Email- or OAuth Verified Registration |
| Yellow | Accountable | Clear storage, Encrypted transmit | Password, Registered, Approval, Click-through DUA |
| Orange | | ...ed storage, ...ed transmit | Password, Registered, Approval, Signed DUA |
| Red | Fully accountable | ...ed storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |

CRCS Center for Research on Computation and Society
at Harvard School of Engineering and Applied Sciences

*Ben–Gurion University of the Negev*

NSF

DATA PRIVACY LAB

We present a framework for **formally describing**, **reasoning about**, and **arriving at** data-handling policies

We present a framework for **formally describing**, **reasoning about**, and **arriving at** **data-handling policies**

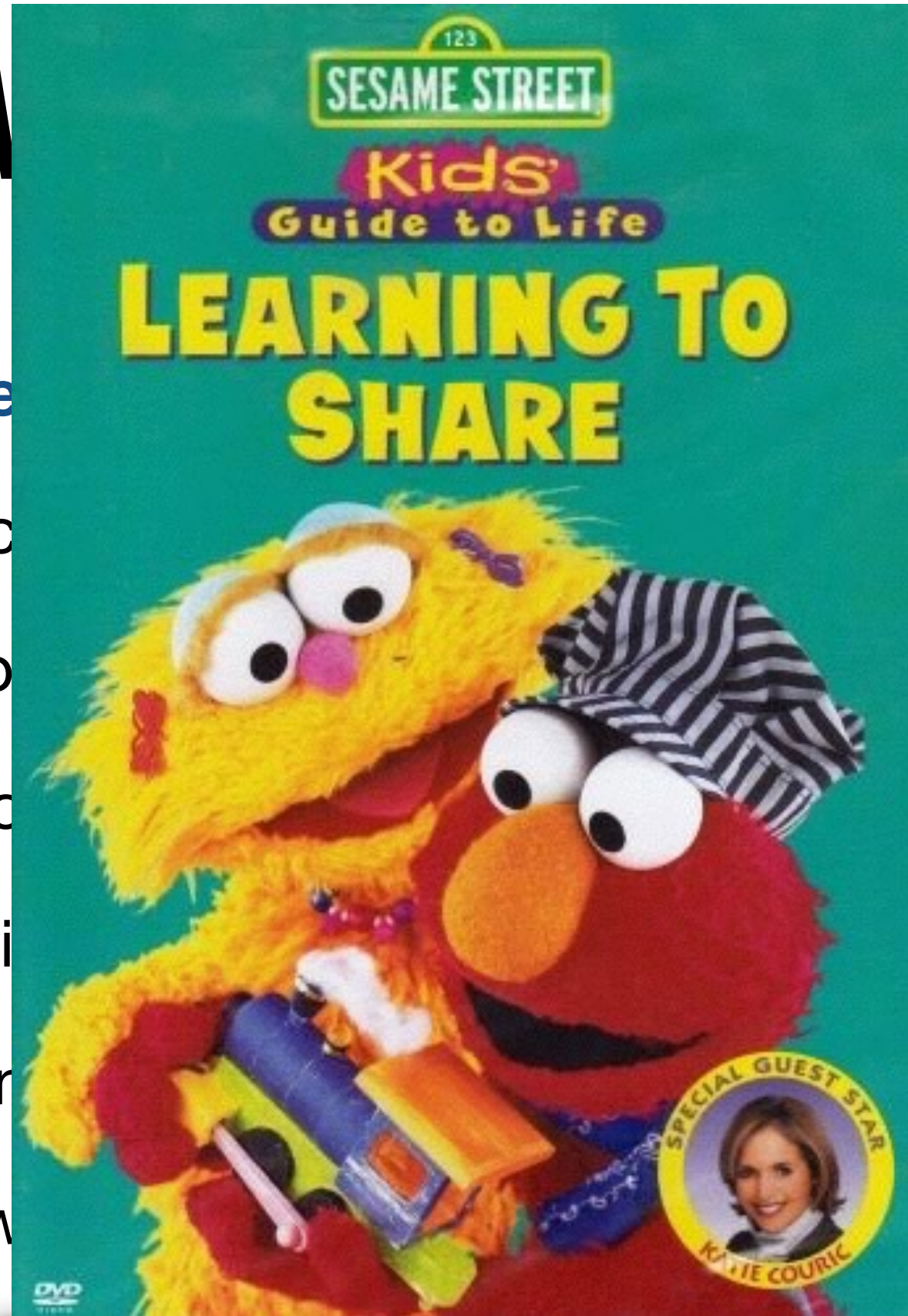*Making it Easier to store and share scientific datasets*

# Why Share Data?

◉ **Good Science**

  ◉ Transparency

  ◉ Collaboration

  ◉ Research acceleration

  ◉ Reproducibility

  ◉ Data citation

◉ **Compliance** with requirements from sponsors and publishers

# W ?

- **Good Science**

  - Transparenc

  - Collaboratio

  - Research ac

  - Reproducibi

  - Data citation

- **Compliance** w s and publishers

# Sharing Data is Nontrivial

◉ *Sharing may harm the data subjects*

◉ **Law** is complex

  ◉ 2187 privacy laws in the US alone, at federal, state and local level, usually context-specific [Sweeney, 2013]

◉ **Technology** is complex

  ◉ E.g. encryption standards change constantly, as new vulnerabilities are found

◉ **Specific dataset provenance** (may be) complex

**Dataset handling policies play the critical role of balancing <span style="color:red">privacy risks</span> and <span style="color:green">scientific value</span> of sharing datasets.**

# Here are some
# Data Handling Policies

-4-

7. The restrictions, if any, co
of this agreement shall last for:

[   ] 5 years;   [   ] 10 years:

[   ] other _____
(subject to approval by Cent

8. Over the years, literary ric
problems for a data-archive. A scholar wh
utilize a questionnaire, a code book or ot
after its contributor has died may not be
may be other difficulties, including illne
which will make it impossible for a social
to use the material. Therefore, the Cent
transferred to the Henry A. Murray Researc

A contributor can only transfer copyi
data-set (e. g. questionnaire, codebook, s
etc.) which were personally created by ti
created for the contributor(s) as a work
was transferred to the contributor. If se
is material in which persons other than th
the Center must be informed so as not to

Under the copyright law which took ef
written transfer of copyright is needed i
the physical property. Copyright lasts fc
years.

The contributor makes the following

(a) Contributor warrants that the m
his/her own, except for those contributed
copyright notice in the name of a person
that they do not infringe upon the rights
are those as set forth herein:

_____
_____
_____

(b) Contributor agrees that the mat
shall become the property of Radcliffe Co
including copyright, of the contributor a
understood that each contributor shall ha
data in any future research or publicatio
rights by the contributor are those set f

_____
_____

---

RADCLIFFE COLLEGE   Ten Garden Street, Cambridge, Massachusetts 02138   (617) 495-8140

The Henry A. Murray Research Center: A Center for the Study of Lives

MEMORANDUM OF AGREEMENT BETWEEN

THE HENRY A. MURRAY RESEARCH CENTER OF RADCLIFFE COLLEGE

AND DATA CONTRIBUTORS

This agreement is made between _____ (contributor)
and Radcliffe College regarding the data set entitled _____

_____
The Henry A. Murray Research Center is a division of Radcliffe College.
_____ deposit(s) in the Henry A. Murray
Research Center (Center) the following materials:

Completed Questionnaire Surveys from it xxx

1. The Center will pay for all costs involved in acquiring the
materials specified above including the costs of removing the names and
such other identifying information as determined by the Center or the
contributor. The total cost shall not exceed $ _____ .

2. The contributor
A. [x] believes there is reason to maintain the anonymity
of each individual respondent.

B. [ ] believes there is no reason to maintain the anonymity
of each individual respondent.

3. If Box A in 2 is checked, the following information shall be de-
leted from the materials. Check all that apply.

---

-3-

ata may be used only at the Center. (If this
copies of machine-readable data may be trans-
other locales. However, in no case will non-
ata be released for use elsewhere.

if any, on use of the material: _____

_____

g provisions relate to the follow-up of the
)

contributor will allow the sample to be followed-
y researchers affiliated with the Center subject
he following conditions:

[   ] A follow-up study may only be performed
with the collaboration of the contributor.

[   ] The contributor will provide the Center
with the names and addresses of the sub-
jects, with:

[   ] no further restriction. These
identifiers may be made available to
affiliated researchers who may be per-
mitted to make contacts with subjects
at the discretion of the Center.

[   ] the restriction that these identifiers
may be used only by the Center staff.

[   ] the restriction that any contacts of
the subjects must be made through the
contributor unless he/she gives written
permission to the researcher to make
such contacts.

[   ] The contributor will not provide the Center
with the names and addresses of the subjects.
Any contacts of the subjects must be made
through him/her.

contributor will not allow the sample to be
owed-up by researchers affiliated with the Center.

# Here are some new
# Data Handling Policies

# Formal$_{cs}$ DHPs

**W3C's Privacy Preference Project (P3P)**

Focuses on web data collection

**Open Digital Rights Language (ODRL)**

Models DRM, supports privacy and rule-based assertions

**PrimeLife Policy Language (PPL)**

Focuses on downstream usage, using rules

**Data-Purpose Algebra** [Hanson, Berners-Lee, Kagal, Sussman, Weitzner]

Models restriction transformation along data processing path

# DataTags

| Tag Type | Description | Security Features | Access Credentials |
|---|---|---|---|
| Blue | Public | Clear storage, Clear transmit | Open |
| Green | Controlled public | Clear storage, Clear transmit | Email- or OAuth Verified Registration |
| Yellow | Accountable | Clear storage, Encrypted transmit | Password, Registered, Approval, Click-through DUA |
| Orange | More accountable | Encrypted storage, Encrypted transmit | Password, Registered, Approval, Signed DUA |
| Red | Fully accountable | Encrypted storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |
| Crimson | Maximally restricted | Multi-encrypted storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |

*DataTags and their respective policies*

Data-handling policies consist of independent *aspects*.
*Encryption at rest, transfer type, access credentials, etc.*

Data-handling policies consist of independent *aspects*.
  *Encryption at rest, transfer type, access credentials, etc.*

Each aspect has multiple *possible requirements*, and can be defined such that these requirements are ordered.

# DHPs: From Text to Space

Data-handling policies consist of independent *aspects*.
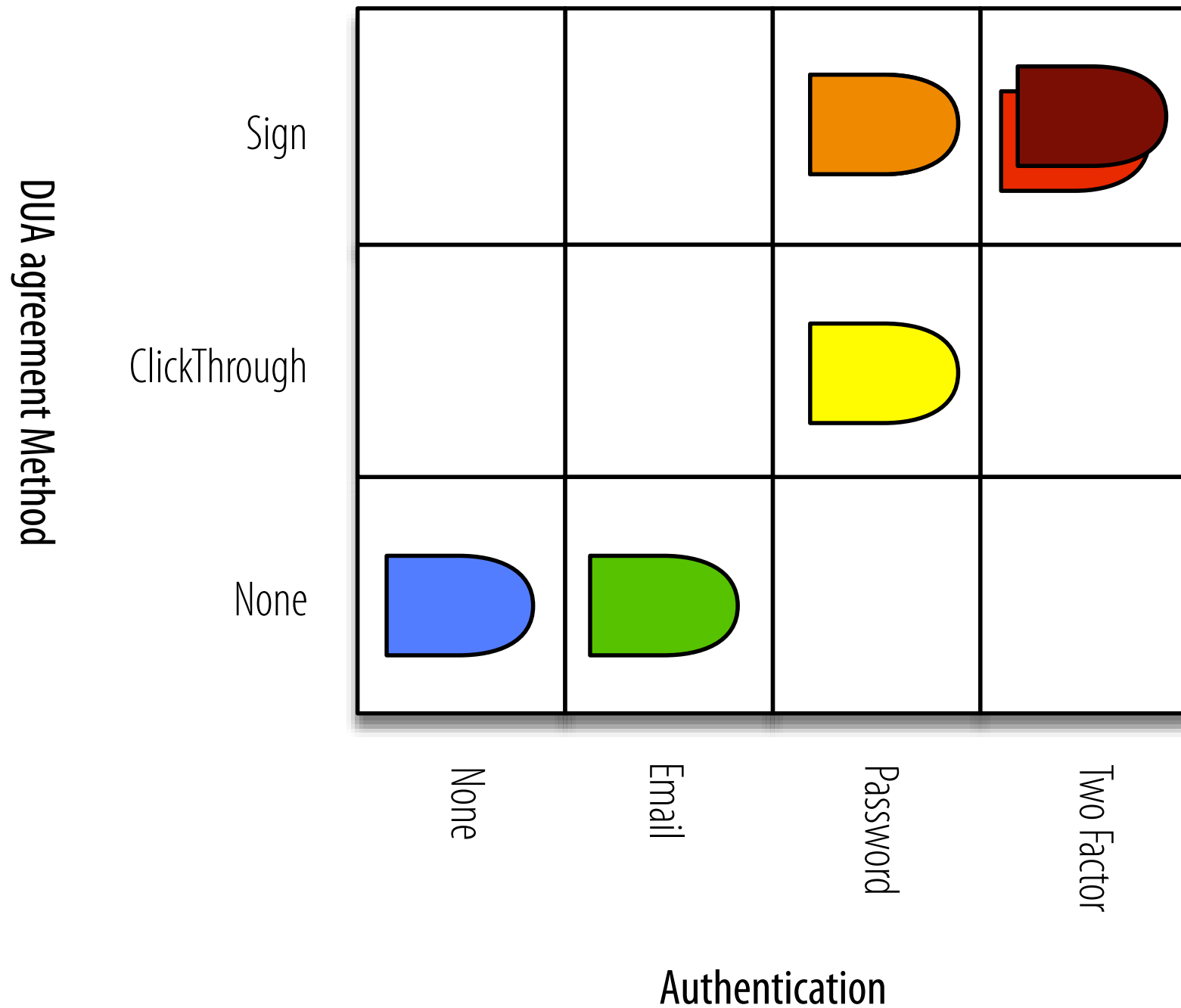*Encryption at rest, transfer type, access credentials, etc.*

Each aspect has multiple *possible requirements*, and can be defined such that these requirements are ordered.

*We can construct a data-handling policy space by viewing aspects as axes, where each aspect's possible requirements serves as its coordinates.*

# Going from this…

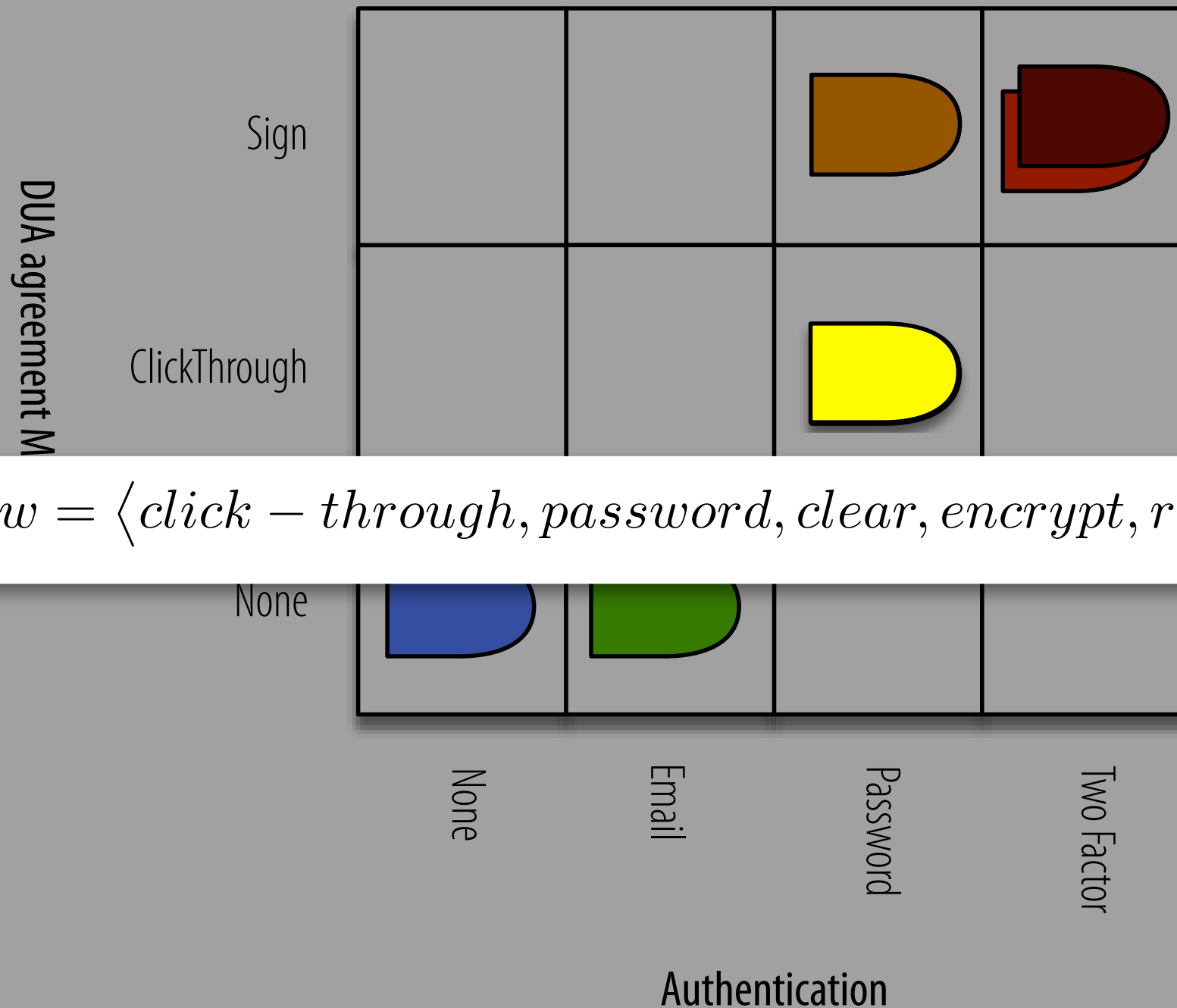| Tag Type | Description | Security Features | Access Credentials |
|---|---|---|---|
| Blue | Public | Clear storage, Clear transmit | Open |
| Green | Controlled public | Clear storage, Clear transmit | Email- or OAuth Verified Registration |
| Yellow | Accountable | Clear storage, Encrypted transmit | Password, Registered, Approval, Click-through DUA |
| Orange | More accountable | Encrypted storage, Encrypted transmit | Password, Registered, Approval, Signed DUA |
| Red | Fully accountable | Encrypted storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |
| Crimson | Maximally restricted | Multi-encrypted storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |

# ...to This*

DUA agreement Method

Sign

ClickThrough

None

Authentication

None  Email  Password  Two Factor

* Shown here is a 2-D projection over the DUA Agreement Method and Authentication axes.

# …to This*



$$Yellow = \langle click - through, password, clear, encrypt, registered, approval \rangle$$

DUA agreement M

Sign

ClickThrough

None

None    Email    Password    Two Factor

Authentication

*Shown here is a 2-D projection over the DUA Agreement Method and Authentication axes.*

# Strictness

# Strictness



compliance(P)

DUA agreement Method

- Sign
- Click-Through
- Implied

Authentication

- None
- Email/OAuth
- Password
- Two Factor

# Strictness



compliance(P)

All policies that do not breach P

P

DUA agreement Method

Sign

Click-Through

Implied

Authentication

None

Email/OAuth

Password

Two Factor

# Lenience



Lenience

DUA agreement Method

Sign

Click-Through

Implied

support(P)

P

*All policies that P does not breach*

None

Email/OAuth

Password

Two Factor

Authentication

# A Dataset and a Repository walk into a DHP space…

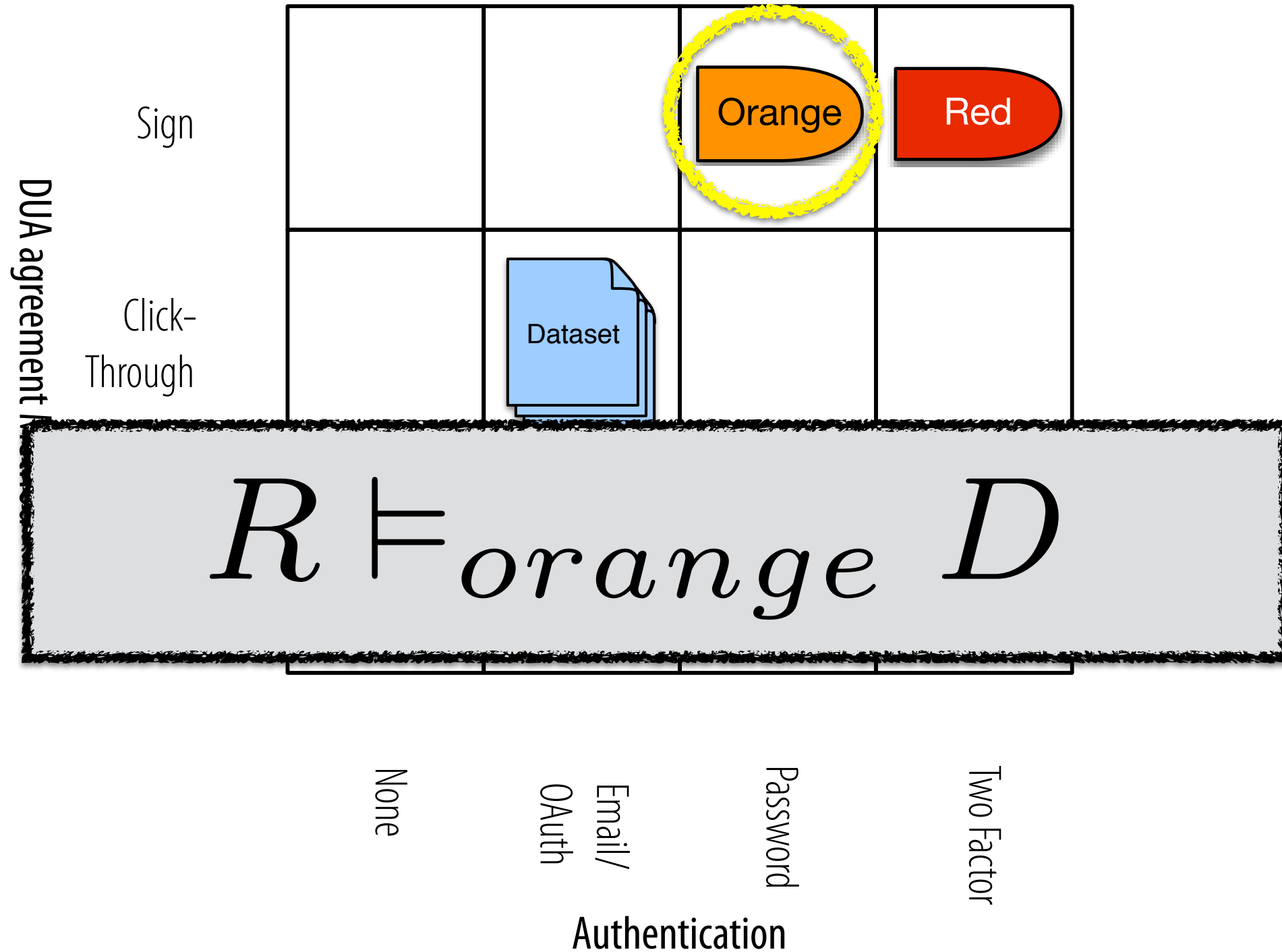# A Dataset and a Repository walk into a DHP space…

# A Dataset and a Repository walk into a DHP space…

compliance(Dataset)

DUA agreement Method

Sign — Orange | Red

Click-Through — Dataset

Implied — Blue | Green

Authentication: None | Email/OAuth | Password | Two Factor

A Dataset and a Repository walk into a DHP space…

# A Dataset and a Repository walk into a DHP space…



Sign

DUA agreement

Click-Through

Orange

Red

Dataset

$$R \vDash_{orange} D$$

None

Email/OAuth

Password

Two Factor

Authentication

# DataTagsTools

```
t1/0.8/definitions.ts ../WORK/dt1/0.8/questionnaire.dg
Reading definitions: ../WORK/dt1/0.8/definitions.ts
 (full:  /Users/michael/Documents/PhD/IQSS/Data-Tags/DataTaggingLibrary/DataTagsLi
b/dist/../WORK/dt1/0.8/definitions.ts)
Reading decision graph: ../WORK/dt1/0.8/questionnaire.dg
 (full:  /Users/michael/Documents/PhD/IQSS/Data-Tags/DataTaggingLibrary/DataTagsLi
b/dist/../WORK/dt1/0.8/questionnaire.dg)

     +--------------
    +|               \
   +||               o)
  +|||               /
  |||+--------------
  ||+--------------
  |+--------------
  +-------------    _  _   _   _
             |  _ | |_ _ |_ _ |_ _   _
             | | || | (_| |_ | (_| (_|  (_  \
             |_|_|\/ \_,_|\_|\_|\_,_|\_,  |_|
                     datatags.org|__/

# Run Started
Do the data concern living persons?
Possible Answers:
 - yes
 - no
answer (? for help): 
```

# Tag Space

```
    BlueToCrimson.ts              •

1   <*¬
2   This is the tag space for the DataTags set proposed at:¬
3   ·Latanya Sweeney, Mercè Crosas, and Michael Bar-Sinai. Sharing sensitive data with conf
·   │ Science, 2015.¬
4   *>¬
5   ¬
6   DataTags: consists of Security, AccessCredentials. <-- This is the top-level slot¬
7   ¬
8   Security: consists of Storage, Transmit.¬
9   ¬
10  AccessCredentials: consists of Authentication, Registration, Approval, DUAAcceptance.¬
11  ¬
12  Storage[How are data stored on disk]: one of¬
13  ·clear [No encryption used],¬
14  ·encrypt [Data are stored encrypted on disk],¬
15  ·multiEncrypt [Data are encrypted on disk, in a way that the server cannot unencrypt th
16  │¬
17  Transmit[How are data travelling through networks]:¬
18  │··one of clear, encrypt.¬
19  ¬
20  Authentication: one of none, password, twoFactor.¬
21  ¬
```

# Tag Space

```
   BlueToCrimson.ts                    •

1   <*¬
2   This is the tag space for the DataTags set proposed at:¬
3   ┊Latanya Sweeney, Mercè Crosas, and Michael Bar-Sinai. Sharing sensitive data with conf
•   ┊Science, 2015.¬
4   *>¬
5
6   DataTags: consists of Security, AccessCredentials      <-- This is the top-level slot¬
7   ┊                                                       
8   Security: consists of Storage, Transmit
9   ¬
10  AccessCredentials: consists of Authentication, Registration, Approval, DUAAcceptance.¬
11  ¬
12  Storage[How are data stored on disk]: one of¬
13  ┊clear [No encryption used],
14  ┊encrypt [Data are stored encrypted on disk],¬
15  ┊multiEncrypt [Data are encrypted on disk, in a way that the server cannot unencrypt th
16  ┊
17  Transmit[How are data travelling through networks]:¬
18  ┊one of clear, encrypt.¬
19  ¬
20  Authentication: one of none, password, twoFactor.¬
21  ¬
```

Block Comment
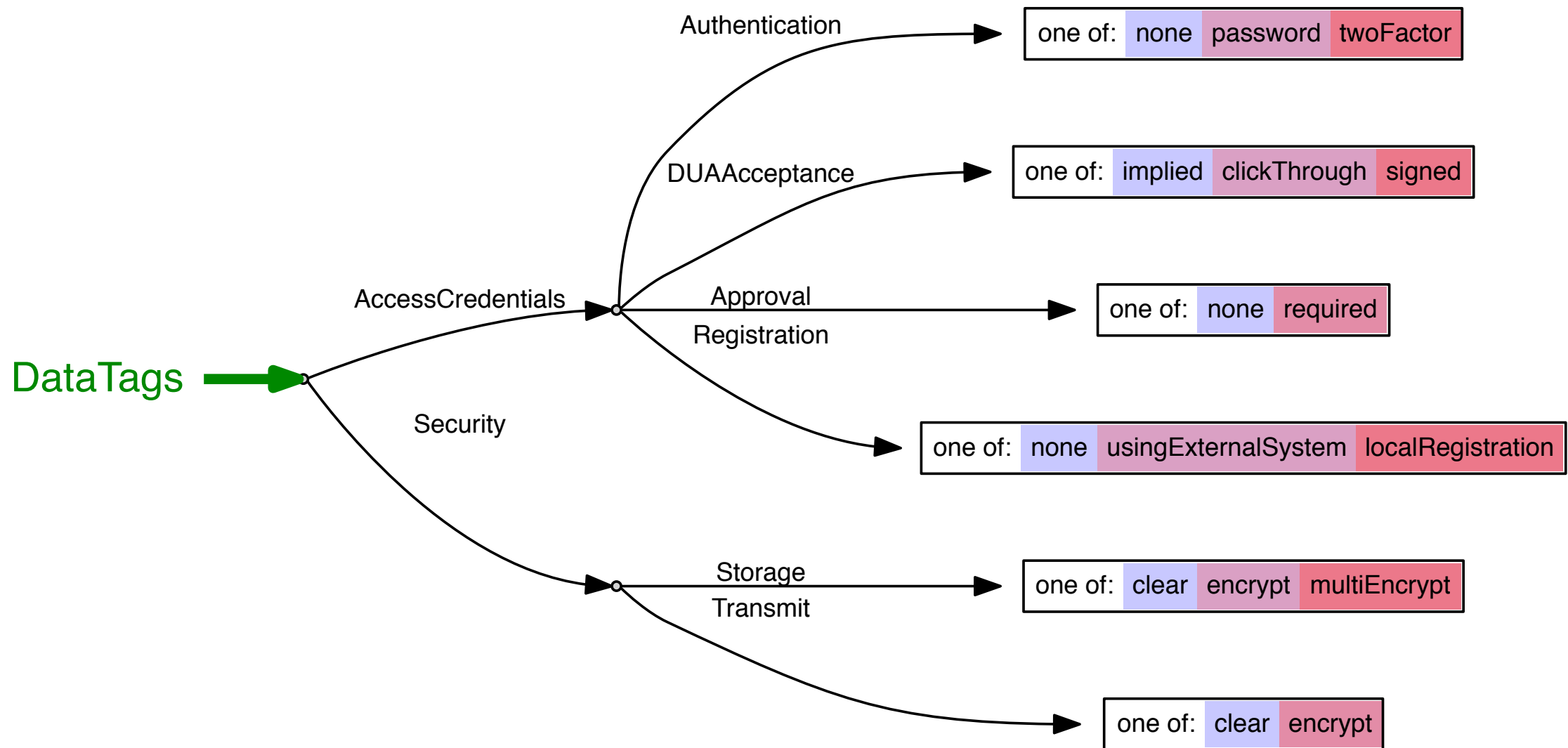
Compound Slot

Line Comment

Description

Atomic Slot

# Tag-Space Visualized



*Visualization using CliRunner (on a later slide) and Graphviz ([www.graphviz.org](www.graphviz.org)).*

# Arriving at a DHP

# Arriving at a DHP

# Tags Questionnaire

❖ "Interview with an expert" metaphor

❖ Consists of a tag space and a *decision graph*

```
                simple.dg

 1   [ask:
 2     {text: Do the data concern humans?}
 3     {answers:
 4       {yes: [set: Assertions+=humanData]
 5             [ask:
 6               {text: Does the data contain educational records?}
 7               {answers:
 8                 {yes: [call: eduCompliance]}
 9                 {no:  [set: Transit=encrypt]}
10               }]}
11       }
12     ]
13   [set: Storage=clear; Transit=clear] <-- defaults
14   [todo: Handle IP issues here]
15   [end]
16   <* Educational Compliance Section *>
17   [>eduCompliance< ask:
```

# Decision Graph - Visualized

# Decision Graph - Visualized



*HIPAA Compliance - Decision Graph*

# Decision Graph - Visualized



HIPAA Compliance - Decision Graph

Nodes and labels visible:

- consent
- agreement
- HIPAASafeHarbor
- HIPAAStatistician
- HIPAALimitedDataset
- HIPAACoveredEntity
- HIPAABusinessAssociate

some of

DataTags → TODO

Labels on edges:
- Basis
- Harm — one of: noRisk, minimal, shame, civil, criminal, maxControl
- Effort — one of: notApplicable, identified, identifiable, deIdentified, anonymous
- DataType
- Approval — one of: none, email, signed
- Use — one of: noRestriction, research, IRB, noProduct
- Publication — one of: noRestriction, notify, preApprove, prohibited
- Reidentify — one of: contact, reidentify, noProhibition, noPeople, noEntities, noMatching
- Acceptance
- TimeLimit — one of: implied, click, signed, signWithID
- DUA
- IP
- Sharing — one of: none, _5years, _2years, _1year
- one of: anyone, notOnline, organization, group, none
- Handling
- Storage — one of: clear, encrypt, multiEncrypt
- Transit — one of: clear, encrypt
- Authentication — one of: none, contactable, password, twoFactor

# Arriving at a DHP



*HIPAA, C.F.R Part 2, FERPA, PPRA,*
*Education Science Report Act (2002), Privacy Act (1974),*
*CIPSEA, Title 13, DPPA*

# CliRunner

❖ Questionnaire Development Console

❖ Run, debug, visualize

❖ Query:
e.g. *what answer sequences result in encryption=clear, harm=severe?*

[set: Thank+=you]
[end]

I ♥ Data

**http://datatags.org**

*http://datascience.iq.harvard.edu/about-datatags*