# Developing Trust in Aggregated Government Data

*Provenance,*
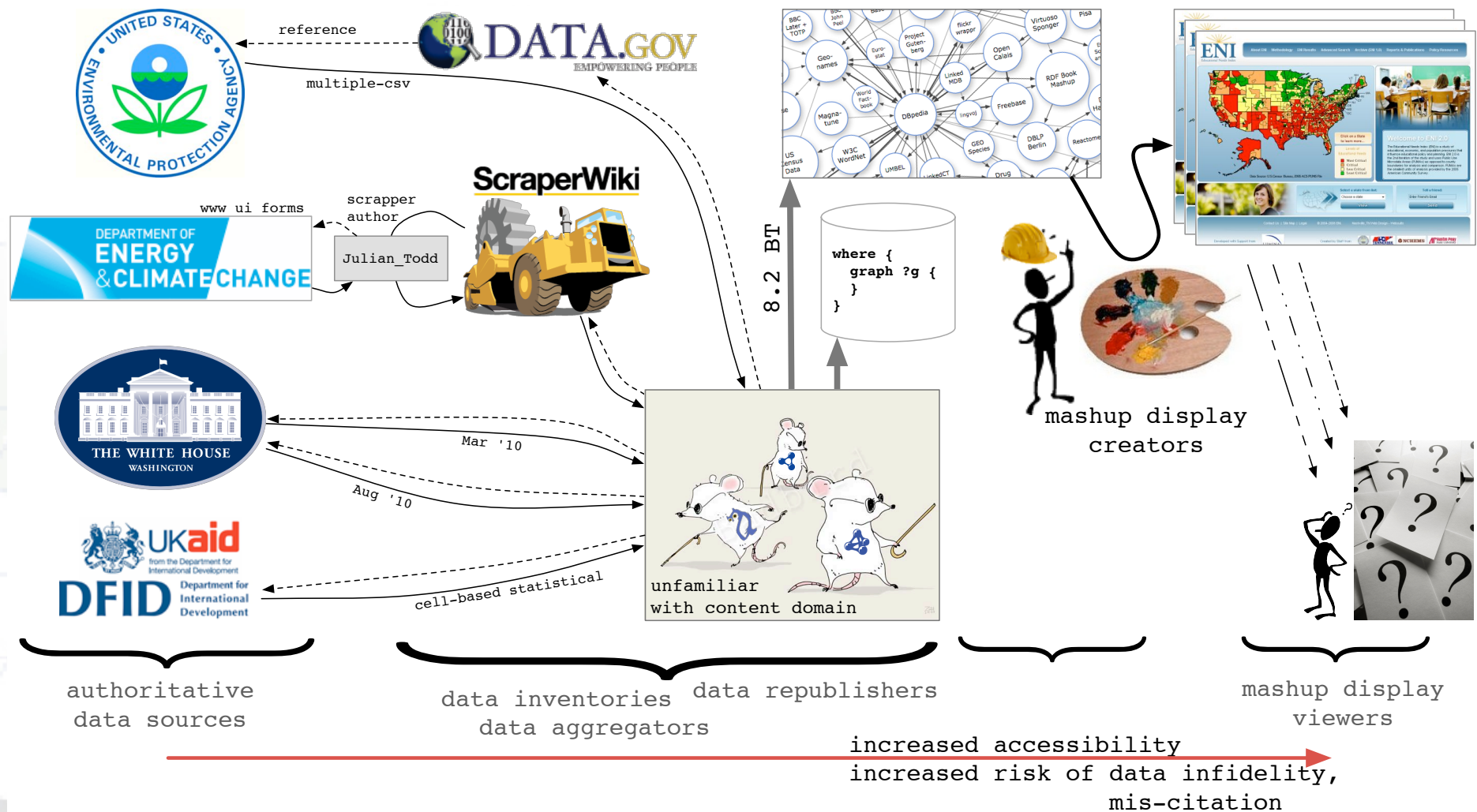
*Interpretation Knowledge*

*and URI Design*

*for*

*Incremental Enhancement of Tabular Data*

Timothy Lebo and Gregory Todd Williams
Tetherless World Constellation
Rensselaer Polytechnic Institute

# Challenges for Data Aggregators



reference

multiple-csv

www ui forms

scrapper author

Julian_Todd

8.2 BT

```
where {
    graph ?g {
    }
}
```

Mar '10

Aug '10

cell-based statistical

unfamiliar with content domain

mashup display creators

authoritative data sources

data inventories data aggregators

data republishers

mashup display viewers

increased accessibility
increased risk of data infidelity, mis-citation

# Challenges for Data Aggregators

Disconnected, Not on Semantic Web

On Semantic Web

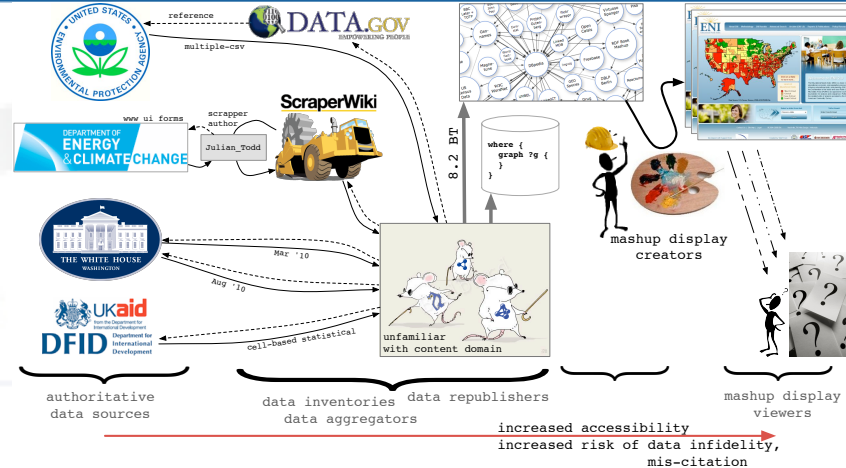http://www.whitehouse.gov/files/disclosures/
visitors/WhiteHouse-WAVES-Released-0910.csv

Trust
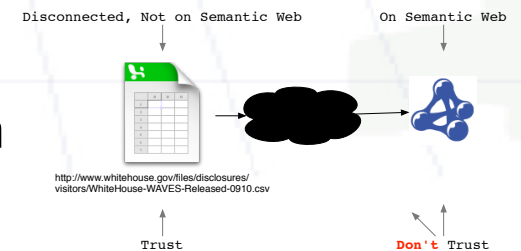
Don't Trust

3

# Assumptions



- Most data are from third-party sources

- Data are updated regularly and irregularly

- Complete interpretation is not immediately possible

- Subsequent interpretations should be backward-compatible

- Provenance is essential
  - Distinguishing among sources
  - Minimizing manual modifications during conversion
  - Tracing to source data
  - Attribution



4

# Outline

- *Challenges for Data Aggregators*
- Trust in aggregated product
  - Provenance capture
  - Parameterized interpretations
  - Intuitively-structured RDF
- Interpretation knowledge vocabulary
  - Naming a Dataset
  - Naïve CSV conversion
  - Specifying an enhancement
  - A few simple examples
- URIs designed for dataset maintenance
  - Predicate Layering
  - Subject Versioning
  - VoID subset Hierarchy
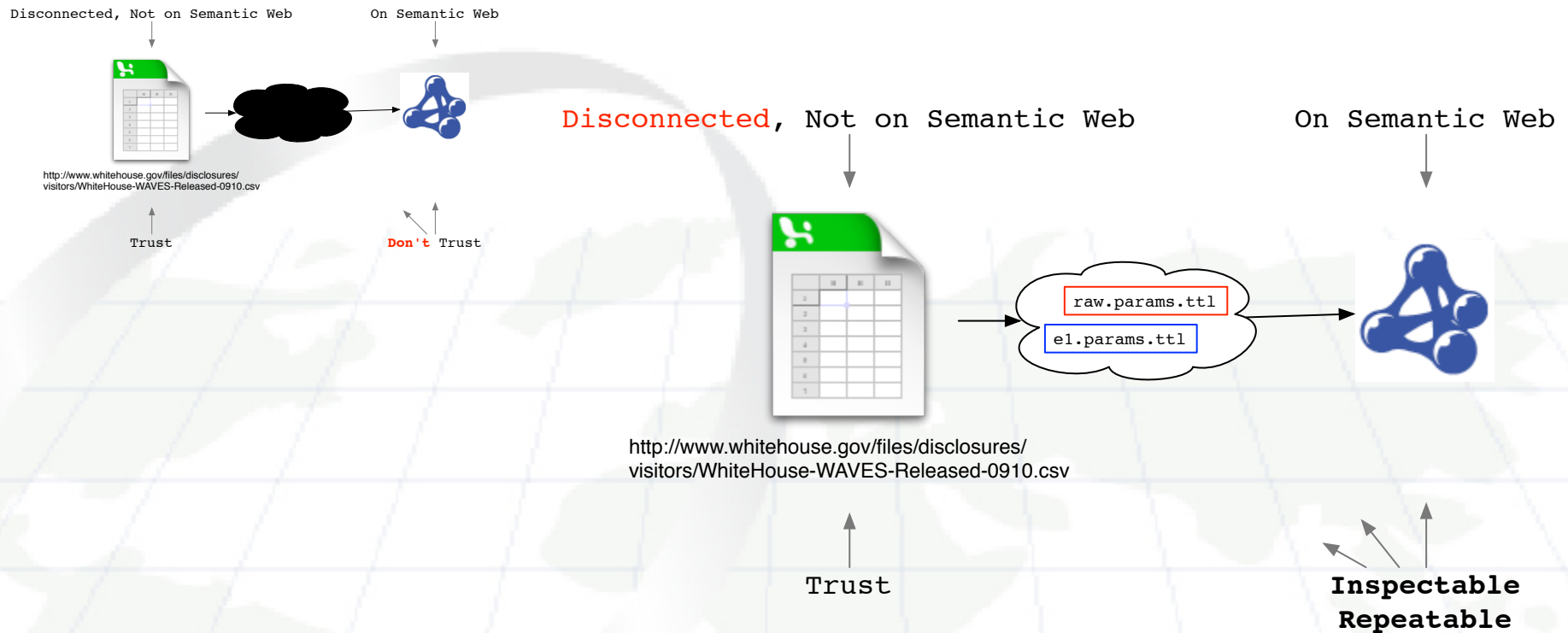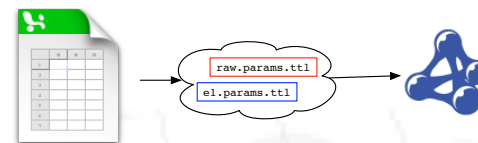  - Hierarchical Integration

# Capturing Conversion Provenance



data.gov/details/1008

1

epa.gov/state_single_pw.zip

2

web

file

/source/data-gov/1008/
version/2010-Aug-30/

source/

state_single_pw.zip
state_single_pw.txt

3

manual/

state_single_pw.txt.csv
state_single_pw.txt.csv.e1.params

4

automatic/

state_single_pw.txt.csv.rdf

5    6    7

publish/

data-gov-1008-2010-Aug-30.ttl.tgz

csv2rdf4lod
.jar

/csv2rdf4lod/bin/*.sh

8

logd.tw.rpi.edu/sparql

http://logd.tw.rpi.edu/source/data-gov/
provenance_file/1008/version/2010-Aug-30/
source/state_single_pw.zip

http://logd.tw.rpi.edu/source/data-gov/
provenance_file/1008/version/2010-Aug-30/
source/state_single_pw.zip.pml.ttl

http://logd.tw.rpi.edu/source/data-gov/
provenance_file/1008/version/2010-Aug-30/
manual/STATE_SINGLE_PW.CSV.e1.params.ttl

http://logd.tw.rpi.edu/source/data-gov/
file/1008/version/2010-Aug-30/conversion/
data-gov-1008-2010-Aug-30.ttl.tgz

1 – Following redirects

2 – Retrieving data file

3 – Unzipping

4 – Manual tweaks

5 – Converter invocation

6 – Predicate lineage

7 – Tracing triple to cell

8 – Populating endpoint

6

Disconnected, Not on Semantic Web          On Semantic Web

http://www.whitehouse.gov/files/disclosures/
visitors/WhiteHouse-WAVES-Released-0910.csv

Trust          Don't Trust

Disconnected, Not on Semantic Web          On Semantic Web

raw.params.ttl

e1.params.ttl

http://www.whitehouse.gov/files/disclosures/
visitors/WhiteHouse-WAVES-Released-0910.csv

Trust          Inspectable Repeatable

- Citing original data location

- Exposing intermediate steps

- Parameterized interpretations w/o modifying original dataset

- Consumer chooses the interpretation layer that they trust

- Producing intuitively-structured RDF

7

# Outline

- Challenges for Data Aggregators
- Trust in aggregated product
  - Provenance capture
  - Parameterized interpretations
  - Intuitively-structured RDF
- Interpretation knowledge vocabulary
  - Naming a Dataset
  - Naïve CSV conversion
  - Specifying an enhancement
  - A few simple examples
- URIs designed for dataset maintenance
  - Predicate Layering
  - Subject Versioning
  - VoID subset Hierarchy
  - Hierarchical Integration

# Naming a Dataset

`<http://base.edu` /source/ `SSS` /dataset/ `DDD` /version/ `VVV>`

| Name Component | Naming Convention | e.g. |
|---|---|---|
| source | organization's dns | "nci-nih-gov" |
| | | "census-gov" |
| | | "nber-org" |
| dataset | ID from organization | "353" |
| | | "tus-cps" |
| | | "stack-heights" |
| version | broad classification | "1st-anniversary" |
| | official release date | "2010-Jan-15" |
| | HTTP last_mod date | "2008-Jul-03" |

"Who?"

"What?"

"Which?"

**Dataset's URI:**

`http://logd.tw.rpi.edu` /source/ `census-gov` /dataset/ `tus-cps` /version/ `2008-Jul-03`

9

```
@prefix void: <http://rdfs.org/ns/void#> .

< http://base.edu /source/ SSS /dataset/ DDD /version/ VVV >
    rdf:type void:Dataset;

    conversion:source_identifier   "SSS" ;
    conversion:dataset_identifier  "DDD" ;
    conversion:version_identifier  "VVV" ;
.
```
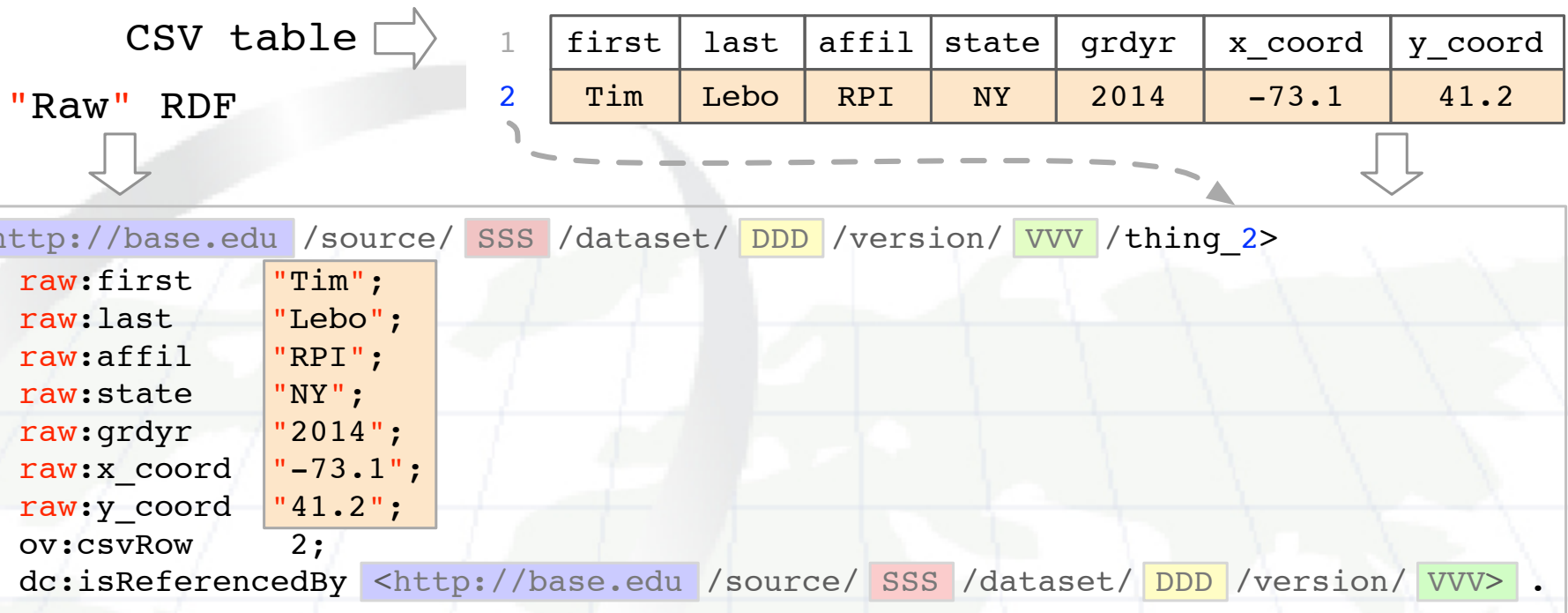
10

# Naïve CSV Conversion

CSV table ➡

"Raw" RDF

| | first | last | affil | state | grdyr | x_coord | y_coord |
|---|-------|------|-------|-------|-------|---------|---------|
| 1 | | | | | | | |
| 2 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

```
<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
    raw:first          "Tim";
    raw:last           "Lebo";
    raw:affil          "RPI";
    raw:state          "NY";
    raw:grdyr          "2014";
    raw:x_coord        "-73.1";
    raw:y_coord        "41.2";
    ov:csvRow          2;
    dc:isReferencedBy <http://base.edu /source/ SSS /dataset/ DDD /version/ VVV> .
```

- Dataset URI re-purposes as subject namespace
- All subjects point to dataset
- Position in original table preserved
- All values are untyped literals
- Rows are rarely conceptually normalized
- Columns may not express a binary relationship

11

# Problem:
# The Real World is Messy

## whitehouse-gov's visitor-records

| NAMELAST | NAMEFIRST | NAMEMID | UIN | BDGNBR | ACCESS_TYPE | TOA | POA | TOD | POD | APPT_MADE_DATE | APPT_START_DATE | APPT_END_DATE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CRUMBLY | ANGELIQUE | | U07467 | | VA | | | | | 5/18/10 9:15 | 5/18/10 9:10 | 5/18/10 23:59 |
| ABRAHAM | ABEBE | | U08781 | | VA | | | | | 5/21/10 10:05 | 5/21/10 14:15 | 5/21/10 23:59 |
| ABRAHAM | ABEBE | | U07781 | 79975 | VA | 5/21/10 14:00 | A0101 | 5/21/10 14:49 | A1 | 5/18/10 17:00 | 5/21/10 14:15 | 5/21/10 23:59 |
| ABRAHAM | AZENEGASH | | U08781 | | VA | | | | | 5/21/10 10:05 | 5/21/10 14:15 | 5/21/10 23:59 |
| ABRAHAM | AZENEGASH | | U07781 | 79125 | VA | 5/21/10 14:00 | A0101 | 5/21/10 14:49 | A1 | 5/18/10 17:00 | 5/21/10 14:15 | 5/21/10 23:59 |
| ABRAHAM | JIITU | | U08781 | | VA | | | | | 5/21/10 10:05 | 5/21/10 14:15 | 5/21/10 23:59 |
| ABRAHAM | JIITU | | U07781 | 79959 | VA | 5/21/10 14:01 | A0101 | | | 5/18/10 17:00 | 5/21/10 14:15 | 5/21/10 23:59 |
| BECK | MIRIAM | | U05979 | 0 | VA | 5/15/10 13:58 | B0402 | | | 5/12/10 12:41 | 5/15/10 14:00 | 5/15/10 23:59 |
| BECK | ROBERT | | U05979 | 0 | VA | 5/15/10 13:58 | B0402 | | | 5/12/10 12:41 | 5/15/10 14:00 | 5/15/10 23:59 |

## data-gov's 1554

| country_name | program_name | | FY1946 | FY1947 | FY1948 | FY1949 | FY1950 | FY1951 |
|---|---|---|---|---|---|---|---|---|
| Afghanistan | Child Survival and Health | | | | | | | |
| Afghanistan | Department of Defense Security Assistance | | | | | | | |
| Afghanistan | Development Assistance | | | | | | | |
| Afghanistan | Economic Support Fund/Security Support Assistance | | | | | | | |
| Afghanistan | Food For Education | | | | | | | |
| Afghanistan | Global Health and Child Survival | | | | | | | |
| Afghanistan | Inactive Programs | | | | | | 1000 | 100000 |

## dfid-gov-uk's statistics-on-international-

| | | | Education | Health | Social Services | Water Supply & Sanitation | Government & Civil Society | Economic | Environment Protection |
|---|---|---|---|---|---|---|---|---|---|
| **Africa: North of Sahara** | | | | | | | | | |
| Algeria | 2004/05 | | - | - | - | - | - | - | - |
| | 2005/06 | | - | - | - | - | - | - | - |
| | 2006/07 | | - | - | - | - | - | - | - |
| | 2007/08 | | - | - | - | - | - | - | - |
| | 2008/09 | | - | - | - | - | - | - | - |
| Egypt | 2004/05 | | 725 | - | - | - | - | - | 2 029 |
| | 2005/06 | | 6 | - | - | - | - | - | 72 |

# Solution:
# Interpretation Parameters

```
@prefix void: <http://rdfs.org/ns/void#> .

<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV>
    rdf:type void:Dataset;
.
```

```
@prefix conversion: <http://purl.org/twc/vocab/conversion/> .

:dataset a void:Dataset;
    conversion:base_uri            "http://base.edu"^^xsd:anyURI;
    conversion:source_identifier   "SSS";
    conversion:dataset_identifier  "DDD";
    conversion:dataset_version      "VVV";
    conversion:conversion_process [
        a conversion:ConversionProcess;
        conversion:enhancement_identifier "1";
        conversion:enhance [
            ...
        ];
        conversion:enhance [
        ];
        ...
    ];
.
```

14

| | first | last | affil | state | grdyr | x_coord | y_coord |
|---|---|---|---|---|---|---|---|
| 1 | | | | | | | |
| 2 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
                    raw:grdyr   "2014" .

```
conversion:enhance [
    ov:csvCol    6;
    ov:csvHeader      "grdyr";

    conversion:label "Completion Year";
];
```

<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
                    e1:completion_year   "2014" .

15

| | first | last | affil | state | grdyr | x_coord | y_coord |
|---|-------|------|-------|-------|-------|---------|---------|
| 1 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
        raw:grdyr   "2014" .

```
conversion:enhance [
    ov:csvCol      6;
    ov:csvHeader      "grdyr";
    conversion:label "Completion Year";
    conversion:range  xsd:gYear;
];
```

<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
        e1:completion_year   "2014"^^xsd:gYear .

16

# Typing with Patterns

```
raw:appt_made_date    "1/5/2010 14:49" .

conversion:enhance [
    ov:csvCol    6;
    ov:csvHeader      "appt_made_date";

    conversion:range              xsd:dateTime;
    conversion:datetime_pattern "M/d/yy HH:mm";
    conversion:datetime_timezone_offset  -300;
];

e1:appt_made_date  "2010-01-05T14:49:00-05:00"^^xsd:dateTime .
```

17

# Promoting a value to a Resource

| | first | last | affil | state | grdyr | x_coord | y_coord |
|---|---|---|---|---|---|---|---|
| 1 | | | | | | | |
| 2 | Tim | Lebo | RPI | NY | 2014 | −73.1 | 41.2 |

```
<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
            raw:affl   "RPI" .
```

```
conversion:enhance [
    ov:csvCol      3;
    ov:csvHeader       "affil";

    conversion:range   rdfs:Resource;
];
```

```
<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>
        e1:affil
            <http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /value-of/affil/RPI> .
```

18

| first | last | affil | state | grdyr | x_coord | y_coord |
|-------|------|-------|-------|-------|---------|---------|
| Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

`<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>`
```
        raw:affl  "RPI" .
```

```
conversion:enhance [
    ov:csvCol    3;
    ov:csvHeader     "affil";

    conversion:range  rdfs:Resource;
    conversion:range_template "[.] [#4]";
];
```

`<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>`
```
        e1:affil
        <http://base.edu /source/ SSS /dataset/ DDD /version/ VVV / value-of/affil/RPI_NY> .
```

19

# Template Variables

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | first | last | affil | state | grdyr | x_coord | y_coord |
| 2 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

| | | |
|---|---|---|
| current value | [.] | RPI |
| row number | [r] | 2 |
| column number | [c] | 3 |
| value of column | [#4] | NY |
| value of property | [@state] | NY |
| property name | [@] | affil |
| property domain | [D] | Graduation |
| property range | [R] | Institution |
| enhancement layer | [e] | 1 |
| versioned dataset namespace | [/sdv] | http://base.edu /source/ SSS /dataset/ DDD /version/ VVV / |
| dataset namespace | [/sd] | http://base.edu /source/ SSS /dataset/ DDD / |
| source namespace | [/s] | http://base.edu /source/ SSS / |
| base namespace | [/] | http://base.edu / |

| | first | last | affil | state | grdyr | x_coord | y_coord |
|---|---|---|---|---|---|---|---|
| 1 | | | | | | | |
| 2 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

`<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>`
         `raw:affl "RPI" .`

```
conversion:enhance [
    ov:csvCol     3;
    ov:csvHeader       "affil";
    conversion:range   rdfs:Resource;
    conversion:links_via   <http://www.rpi.edu/~lebot/lod-links/universities.ttl>;
    conversion:subject_of dcterms:identifier;
];
```
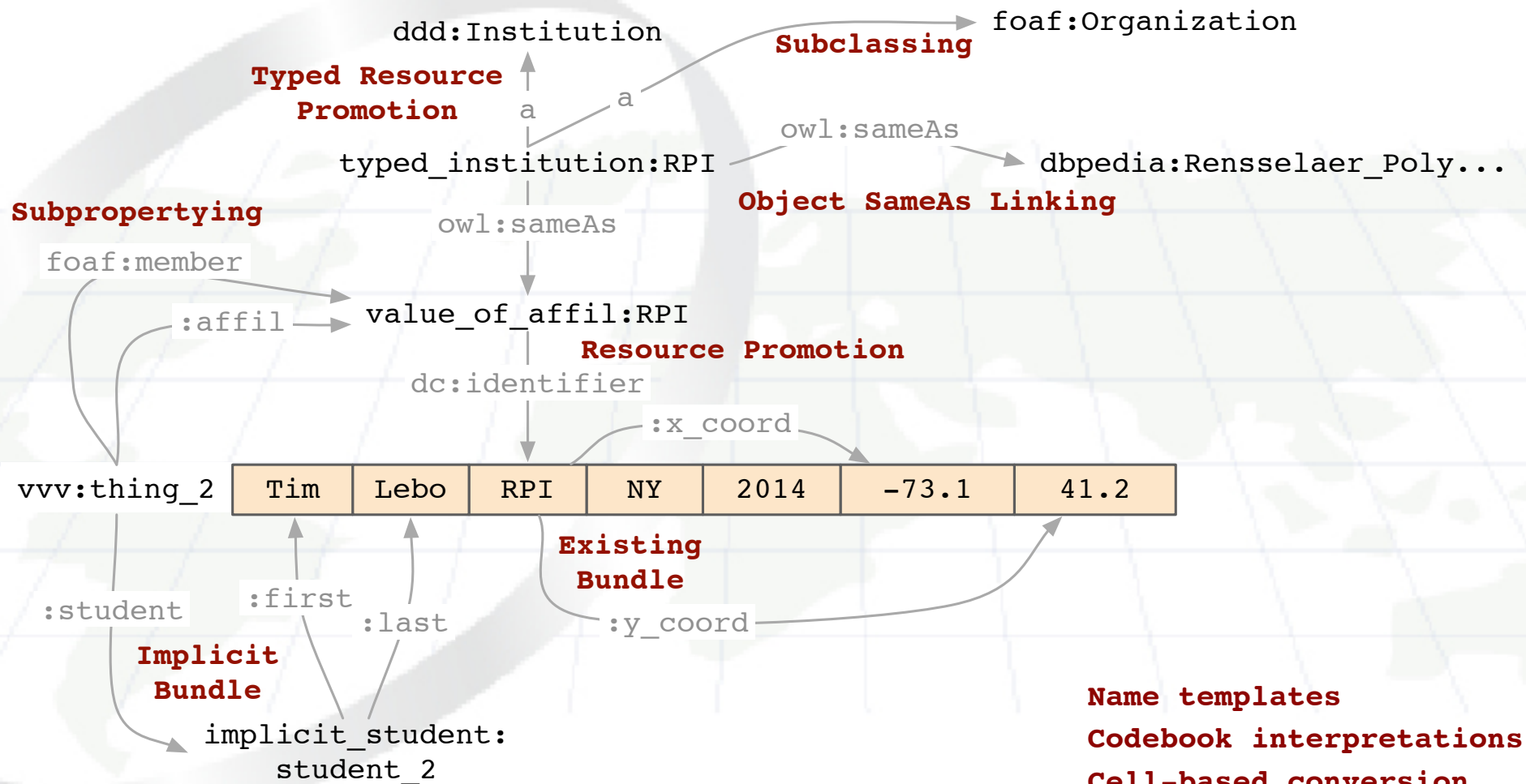
`<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /thing_2>`
     `e1:affil`
        `<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /value-of/affil/RPI> .`

`<http://base.edu /source/ SSS /dataset/ DDD /version/ VVV /value-of/affil/RPI>`
      `owl:sameAs <http://www.dbpedia.org/resource/Rensselaer_Polytechnic_Institute> .`

21

# A Few More Enhancements...

| 1 | first | last | affil | state | grdyr | x_coord | y_coord |
|---|-------|------|-------|-------|-------|---------|---------|

ddd:Institution → **Subclassing** → foaf:Organization

**Typed Resource Promotion**

a ↑   a

typed_institution:RPI — owl:sameAs → dbpedia:Rensselaer_Poly...

**Object SameAs Linking**

owl:sameAs

**Subpropertying**

foaf:member

:affil → value_of_affil:RPI

**Resource Promotion**

dc:identifier

:x_coord

| vvv:thing_2 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |
|-------------|-----|------|-----|----|----|-------|------|

**Existing Bundle**

:student    :first    :last    :y_coord

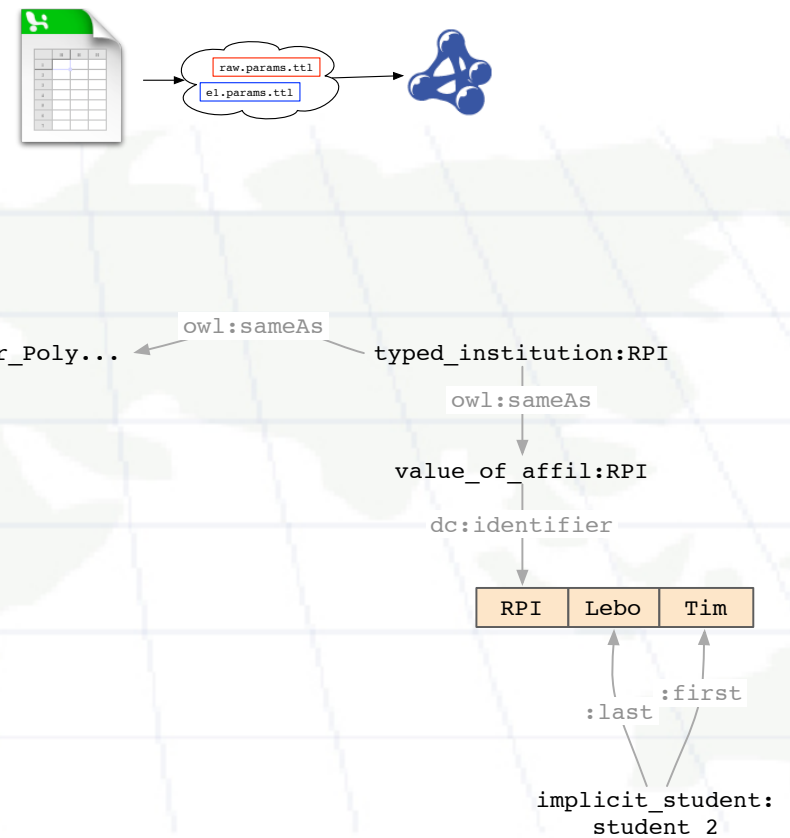**Implicit Bundle**

implicit_student: student_2

**Name templates**
**Codebook interpretations**
**Cell-based conversion**
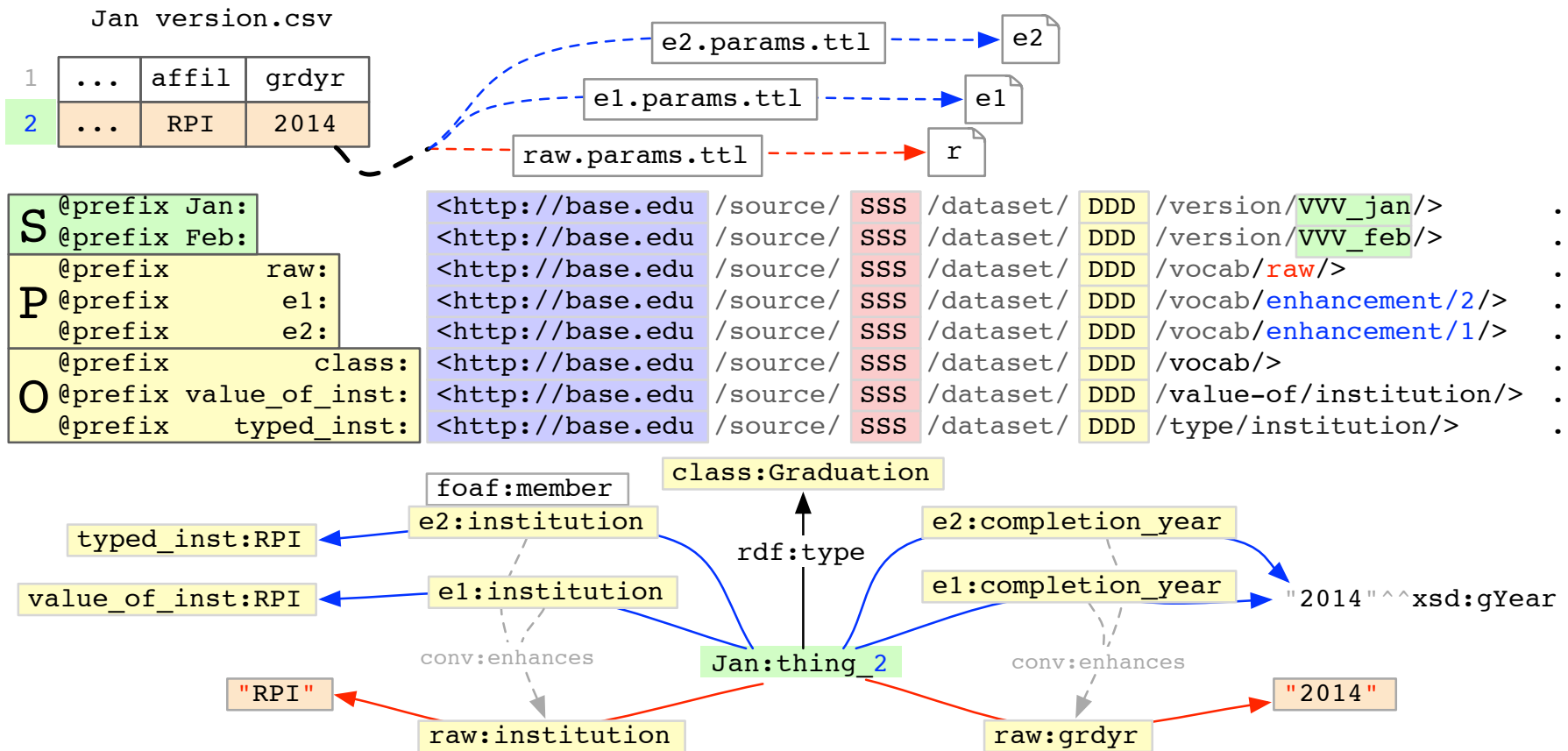**Structure assistance**

# Outline

- Assumptions and Design Objectives
- Trust in aggregated product
  - Provenance capture
  - Parameterized interpretations
  - Intuitively-structured RDF
- Interpretation knowledge vocabulary
  - Naming a Dataset
  - Naïve CSV conversion
  - Specifying an enhancement
  - A few simple examples
- URIs designed for dataset maintenance
  - Predicate Layering
  - Subject Versioning
  - VoID subset Hierarchy
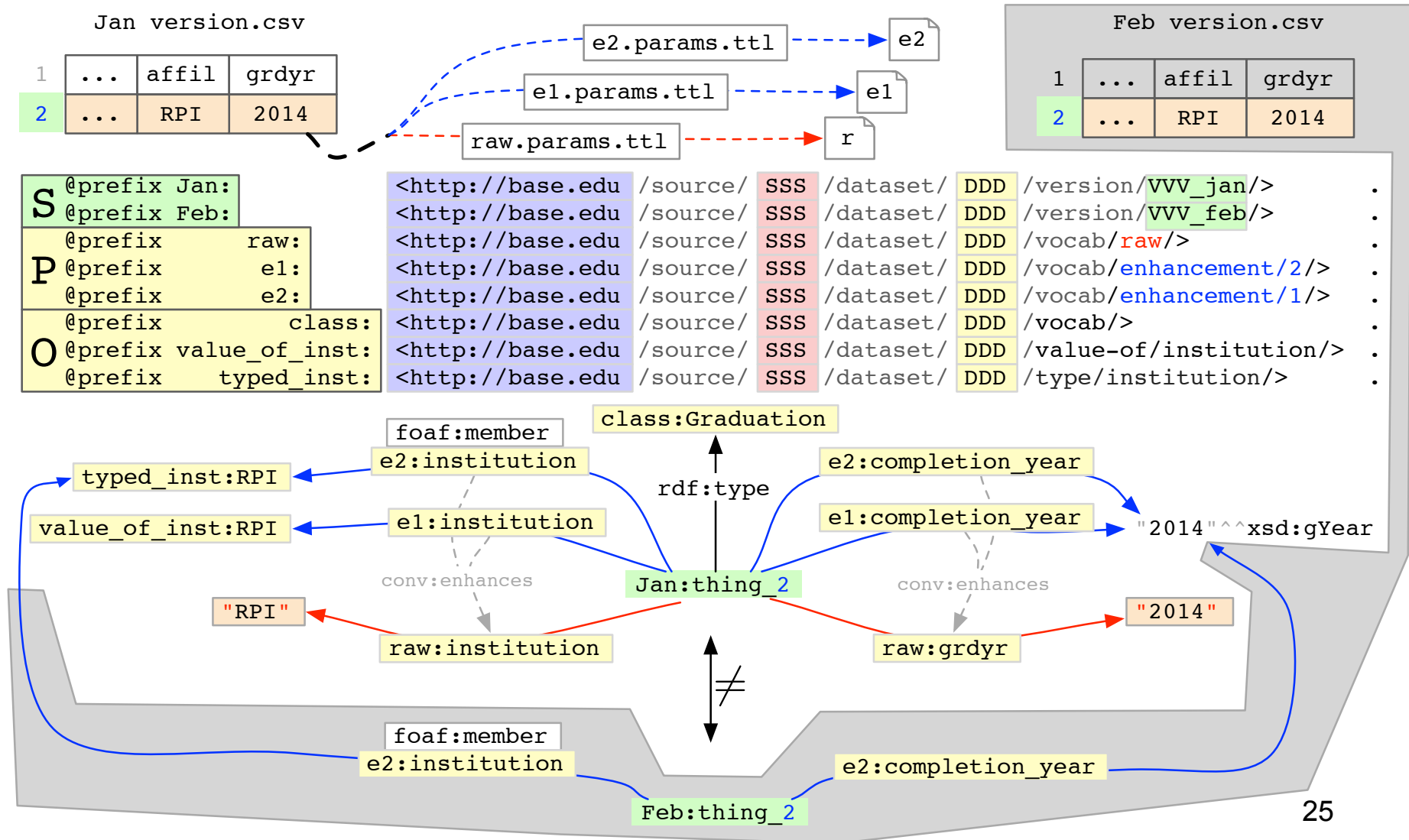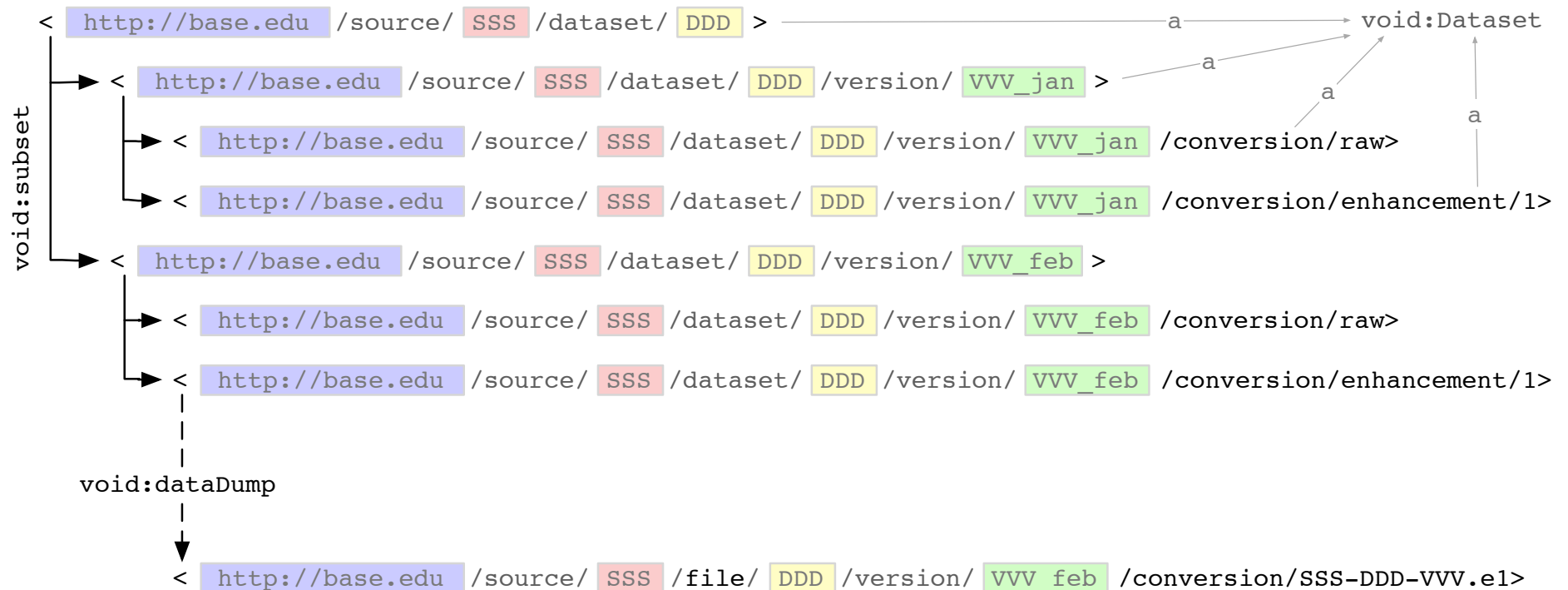  - Hierarchical Integration



23

```
<  http://base.edu  /source/ SSS /dataset/ DDD >  ─────────────────────────────────── a ──────────→ void:Dataset

        <  http://base.edu  /source/ SSS /dataset/ DDD /version/ VVV_jan >  ────────── a ───↗
                                                                                              a
            ▶ <  http://base.edu  /source/ SSS /dataset/ DDD /version/ VVV_jan /conversion/raw>
                                                                                              a
            ▶ <  http://base.edu  /source/ SSS /dataset/ DDD /version/ VVV_jan /conversion/enhancement/1>

        <  http://base.edu  /source/ SSS /dataset/ DDD /version/ VVV_feb >

            ▶ <  http://base.edu  /source/ SSS /dataset/ DDD /version/ VVV_feb /conversion/raw>

            ▶ <  http://base.edu  /source/ SSS /dataset/ DDD /version/ VVV_feb /conversion/enhancement/1>
```

void:subset

void:dataDump

```
              <  http://base.edu  /source/ SSS /file/ DDD /version/ VVV_feb /conversion/SSS-DDD-VVV.e1>
```

26

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | first | last | affil | state | grdyr | x_coord | y_coord |
| 2 | Tim | Lebo | RPI | NY | 2014 | -73.1 | 41.2 |

`http://base.edu` /source/ `SSS` /dataset/ `DDD` /version/ `VVV` /thing_2> raw:affil "RPI" .

```
conversion:enhance [                          conversion:enhance [
   ov:csvCol                3;                   conversion:class_name "Graduation";
   ov:csvHeader            "affil";              conversion:subclass_of
   conversion:domain_name "Graduation";             "[/sd]vocab/RiteOfPassage";
   conversion:range        rdfs:Resource;          "[/s]vocab/RiteOfPassage";
];                                                 "[/]vocab/RiteOfPassage";
                                               <http://www.dbpedia.org/resource/Rite_of_passage>;
                                            ];
```

<http://www.dbpedia.org/resource/Rite_of_passage> ◄──────── rdf:type

< `http://base.edu` /vocab/RiteOfPassage> ◄──────── rdf:type

< `http://base.edu` /source/ `SSS` /vocab/RiteOfPassage> ◄──────── rdf:type

< `http://base.edu` /source/ `SSS` /dataset/ `DDD` /vocab/RiteOfPassage> ◄── rdf:type

< `http://base.edu` /source/ `SSS` /dataset/ `DDD` /vocab/Graduation> ◄── rdf:type

< `http://base.edu` /source/ `SSS` /dataset/ `DDD` /version/ `VVV` /thing_2>

27

# Summary

- Challenges for Data Aggregators
- Trust in aggregated product
    - Provenance capture
    - Parameterized interpretations
    - Intuitively-structured RDF
- Interpretation knowledge vocabulary
    - Naming a Dataset
    - Naïve CSV conversion
    - Specifying an enhancement
    - A few simple examples
- URIs designed for dataset maintenance
    - Predicate Layering
    - Subject Versioning
    - VoID subset Hierarchy
    - Hierarchical Integration