# Neuroinformatic techniques for provenance & data sharing

Camille Maumet

GlaxoSmithKline - Neurophysics Workshop on
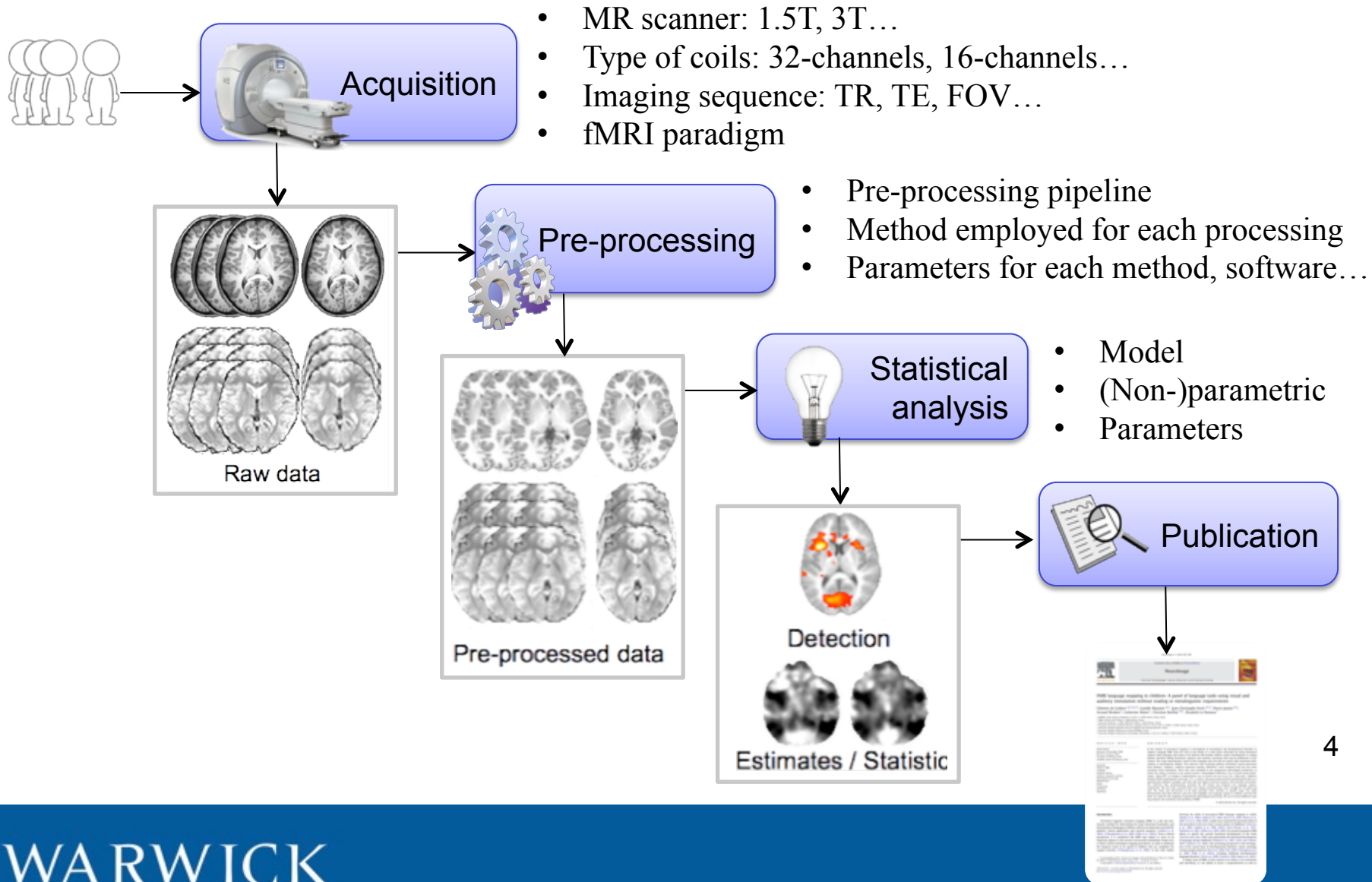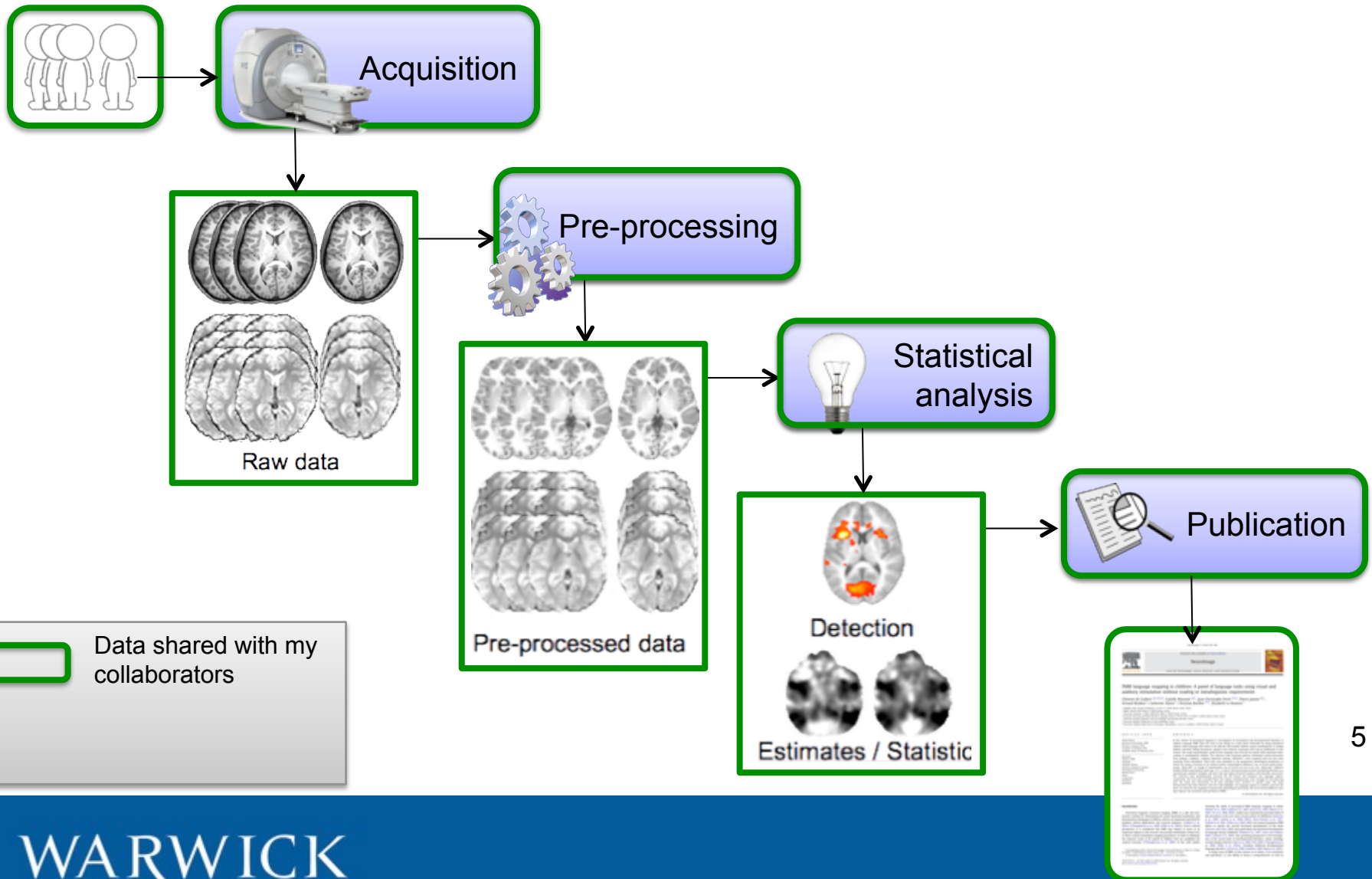Skeptical Neuroimaging

January 14th, 2014

# Outline

1. Data sharing: current practice in neuroimaging
2. How to become less skeptical?
3. Neuroinformatics techniques for provenance and data sharing

WARWICK

# Outline

1. Data sharing: current practice in neuroimaging
2. How to become less skeptical?
3. Neuroinformatics techniques for provenance and data sharing

3

WARWICK

# Overview of a neuroimaging study



**Acquisition**

- MR scanner: 1.5T, 3T…
- Type of coils: 32-channels, 16-channels…
- Imaging sequence: TR, TE, FOV…
- fMRI paradigm

Raw data

**Pre-processing**

- Pre-processing pipeline
- Method employed for each processing
- Parameters for each method, software…

Pre-processed data

**Statistical analysis**

- Model
- (Non-)parametric
- Parameters

Detection

Estimates / Statistic

**Publication**

WARWICK

4

# Neuroimaging and data sharing



Acquisition

Raw data

Pre-processing

Pre-processed data

Statistical analysis

Detection

Estimates / Statistic

Publication

Data shared with my collaborators

WARWICK

5

# Sharing data with my collaborators

WARWICK

# Neuroimaging and data sharing



Acquisition

Raw data

Pre-processing

Pre-processed data

Statistical analysis

Detection

Estimates / Statistic

Publication

Data shared with my collaborators

Data shared with the whole community

WARWICK

# A neuroimaging publication

- *Methods* section: metadata in free-form text.



Table

- *Results* section:



2D plot(s) of the detections



Description of the detections



Table of local maxima

8

WARWICK

# Outline

1. Data sharing: current practice in neuroimaging
2. How to become less skeptical?
3. Neuroinformatics techniques for provenance and data sharing

WARWICK

# Reproducibility



WARWICK

# *Full* provenance

WARWICK

# Meta-analysis: analyzing the analyses

- Coordinate-Based Meta-Analysis (CBMA)



Paper 1    Paper 2    ...    Paper n    New results!

- Image-Based Meta-Analysis (IBMA).



Study 1    Study 2    ...    Study n    New results!

# How to become less skeptical?

- Reproducibility
  - Confirm results by re-running an analysis

- Provenance
  - Needed for reproducibility
  - Avoid selection bias.

- Meta-analysis
  - Strengthen results by combining studies.

- What do we need?
  - Sharing data, meta-data and provenance.

WARWICK

# Data sharing: obstacles

- Psychological
  - "My" data
- Ethical constraints
- Technical: difficulties to share data with enough metadata to be really useful
  - *Available* data versus *usable* data.

"Less than a few percents of acquired neuroimaging data is available in public repositories" [Poline 2012]

WARWICK

# Outline

1. Data sharing: current practice in neuroimaging
2. How to become less skeptical?
3. Neuroinformatics techniques for provenance and data sharing

WARWICK

# Data sharing tools

WARWICK

# A standard format for meta-data

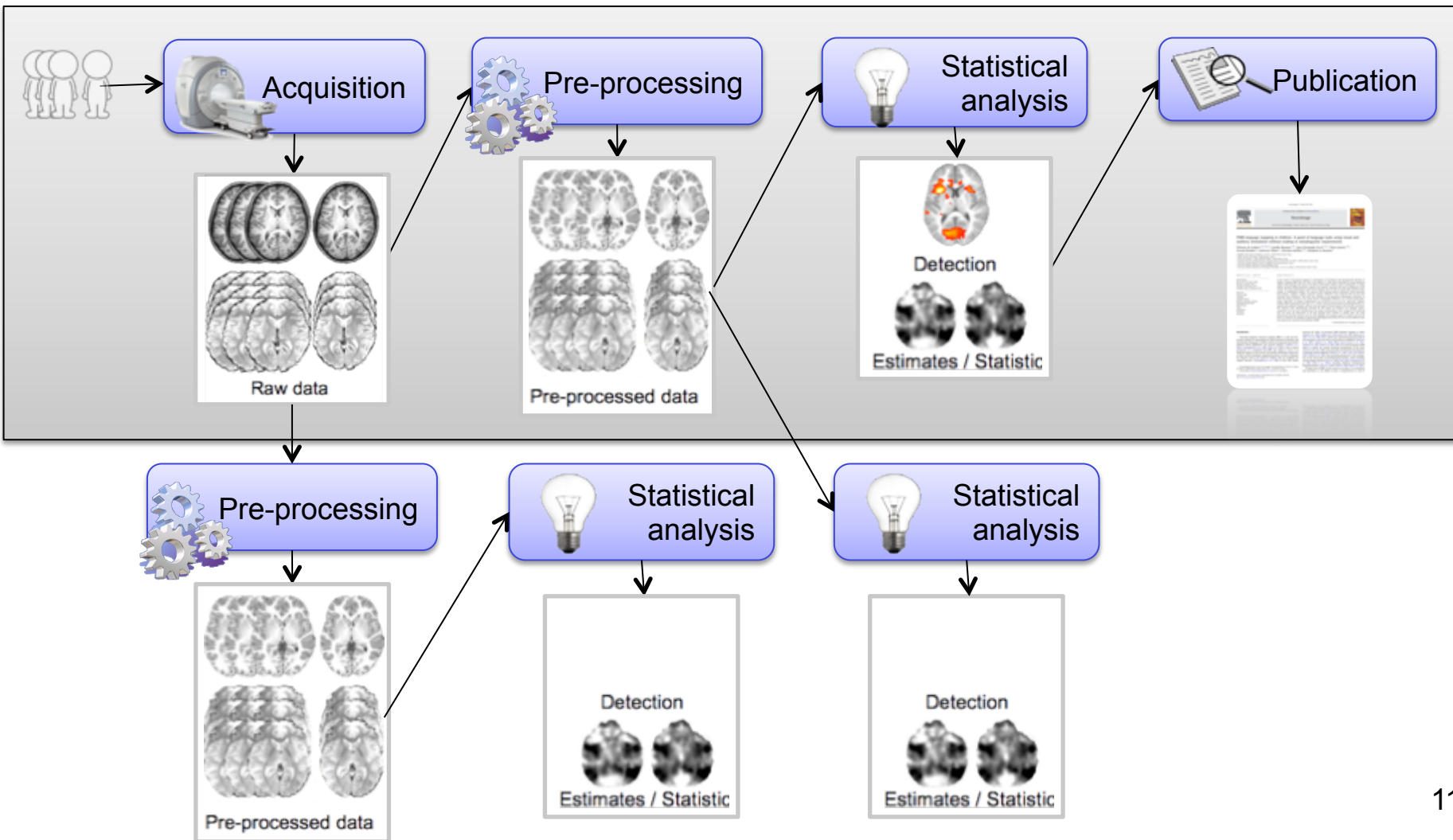- Sharing data across the data sharing tools…
- First attempt of an agnostic format: **X**ML-Based **C**linical **E**xperiment **D**ata **E**xchange Schema (XCEDE): www.xcede.org
  - Describes subject, study, activation
  - Limited provenance encoding
  - Initiative of the BIRN
- **N**euro**I**maging **D**ata **M**odel NI-DM: www.nidm.nidash.org
  - Based on web-semantic tools.
  - Initiative of the BIRN and INCF

WARWICK

# Three major players

- Bottom-up approach.
- Lean on **existing analysis software (SPM, FSL, AFNI)** to disseminate the standard.



Software Package

Automatically created with Neurotrends based on over 16 000 journal articles

18

# Work in progress

- Define a format to represent the results of a neuroimaging study with a focus on meta-analysis.



Vocabulary



Data model

WARWICK

# Neuroimaging terms

- Define a vocabulary to support the format.

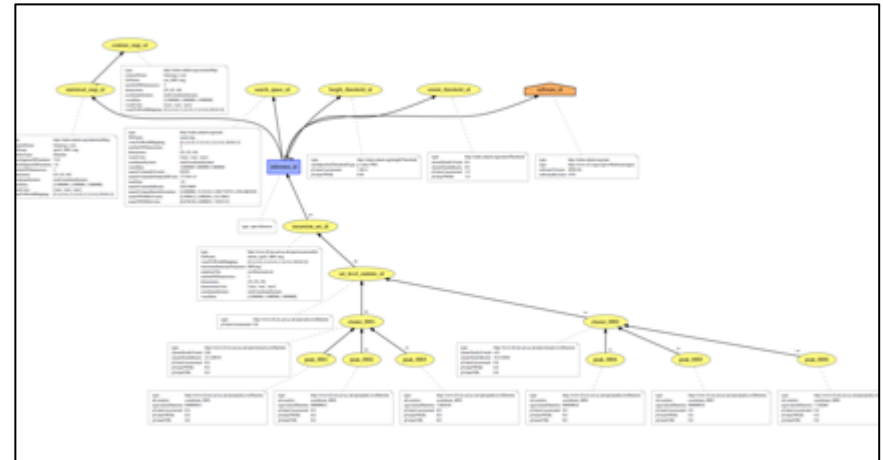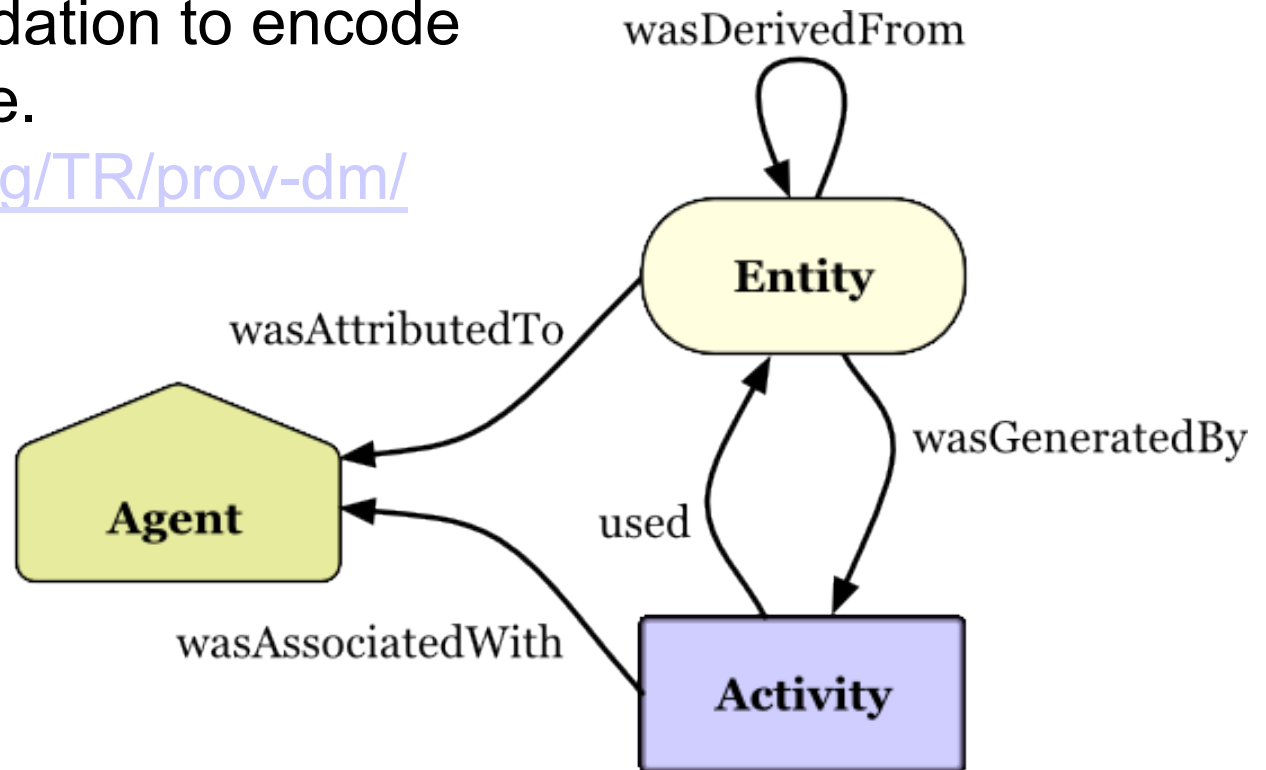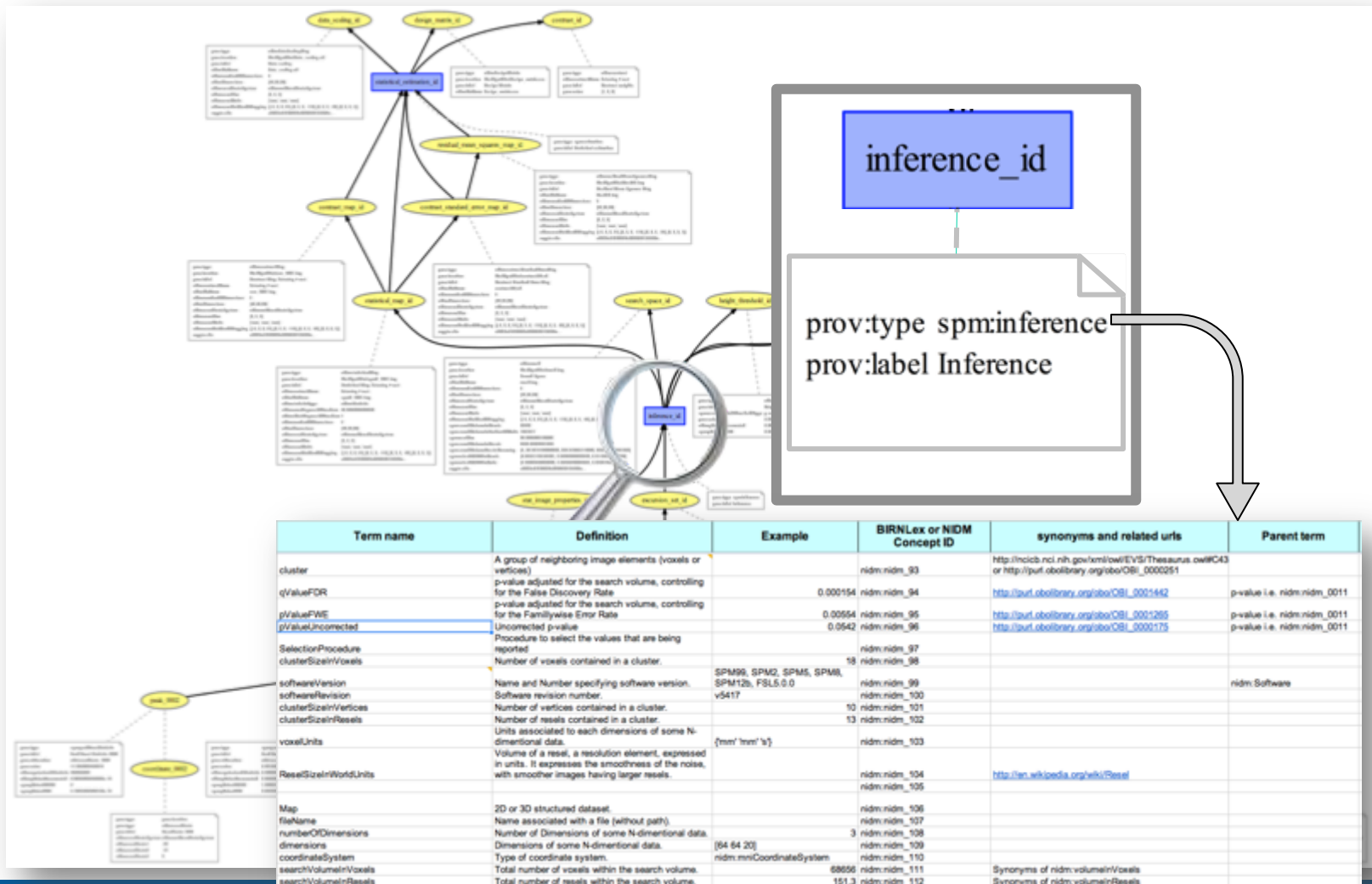| Term name | Definition | Example | BIRNLex or NIDM Concept ID | synonyms and related urls | Parent term |
|---|---|---|---|---|---|
| cluster | A group of neighboring image elements (voxels or vertices) | | nidm:nidm_93 | http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C43 or http://purl.obolibrary.org/obo/OBI_0000251 | |
| qValueFDR | p-value adjusted for the search volume, controlling for the False Discovery Rate | 0.000154 | nidm:nidm_94 | http://purl.obolibrary.org/obo/OBI_0001442 | p-value i.e. nidm:nidm_0011 |
| | p-value adjusted for the search volume, controlling | | nidm:nidm_95 | http://purl.obolibrary.org/obo/OBI_0001265 | p-value i.e. nidm:nidm_0011 |
| pValueUncorrected | Uncorrected p-value | 0.0542 | nidm:nidm_96 | http://purl.obolib | p-value i.e. nidm:nidm_0011 |
| SelectionProcedure | Procedure to select the values that are being reported | | nidm:nidm_97 | | |
| clusterSizeInVoxels | Number of voxels contained in a cluster. | 18 | nidm:nidm_98 | | |
| softwareVersion | Name and Number specifying software version. | SPM99, SPM2, SPM5, SPM8, SPM12b, FSL5.0.0 | nidm:nidm_99 | | nidm:Software |
| softwareRevision | Software revision number. | v5417 | nidm:nidm_100 | | |
| clusterSizeInVertices | Number of vertices contained in a cluster. | 10 | nidm:nidm_101 | | |
| clusterSizeInResels | Number of resels contained in a cluster. | 13 | nidm:nidm_102 | | |
| voxelUnits | Units associated to each dimensions of some N-dimensional data. | {'mm' 'mm' 's'} | nidm:nidm_103 | | |
| ReselSizeInWorldUnits | Volume of a resel, a resolution element, expressed in units. It expresses the smoothness of the noise, with smoother images having larger resels. | | nidm:nidm_104 | http://en.wikipedia.org/wiki/Resel | |
| | | | nidm:nidm_105 | | |
| Map | 2D or 3D structured dataset. | | nidm:nidm_106 | | |
| fileName | Name associated with a file (without path). | | nidm:nidm_107 | | |
| numberOfDimensions | Number of Dimensions of some N-dimentional data. | 3 | nidm:nidm_108 | | |
| dimensions | Dimensions of some N-dimensional data. | [64 64 20] | nidm:nidm_109 | | |
| coordinateSystem | Type of coordinate system. | nidm:mniCoordinateSystem | nidm:nidm_110 | | |
| searchVolumeInVoxels | Total number of voxels within the search volume. | 68656 | nidm:nidm_111 | Synonyms of nidm:volumeInVoxels | |
| searchVolumeInResels | Total number of resels within the search volume. | 151.3 | nidm:nidm_112 | Synonyms of nidm:volumeInResels | |

WARWICK

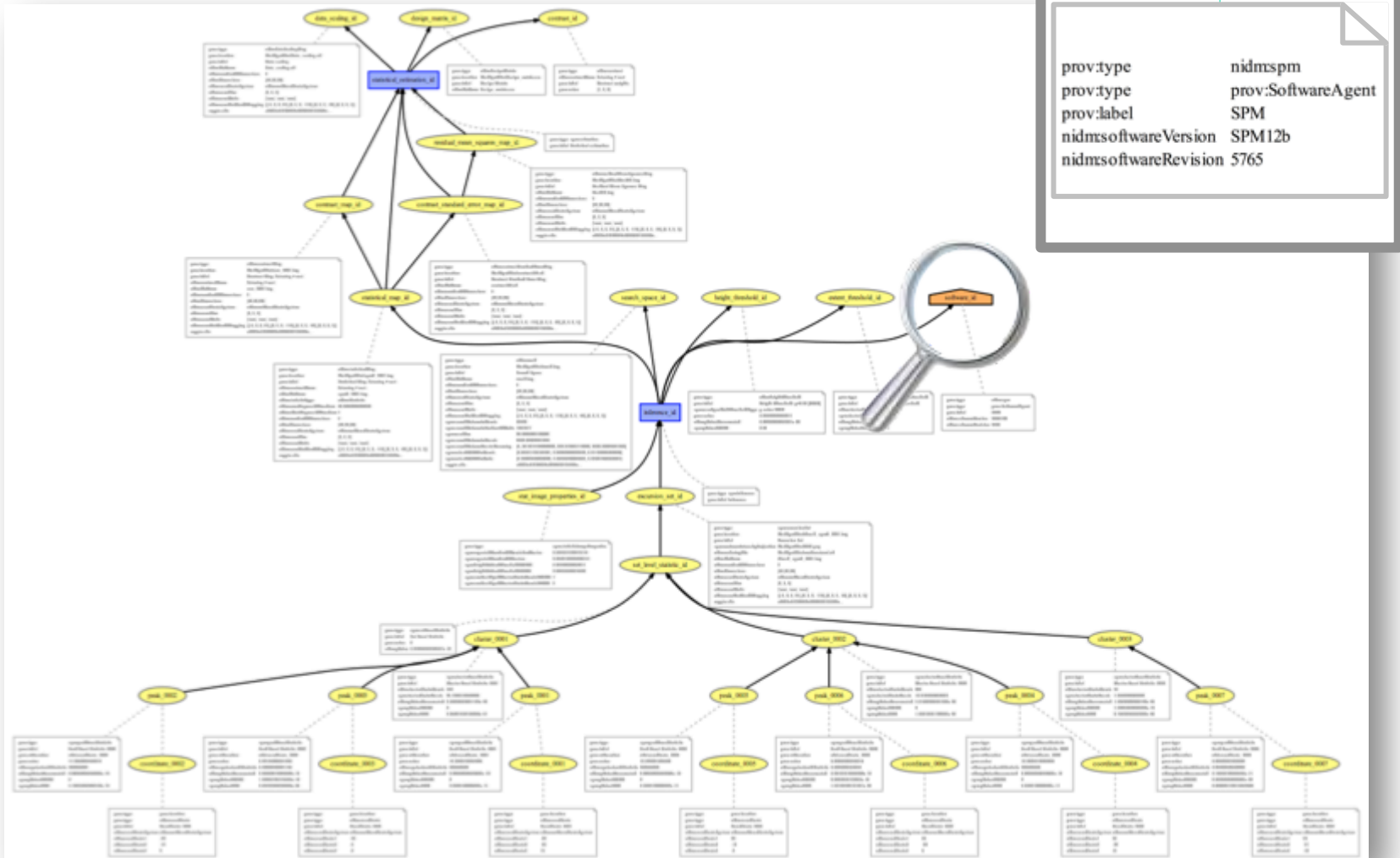# Data model

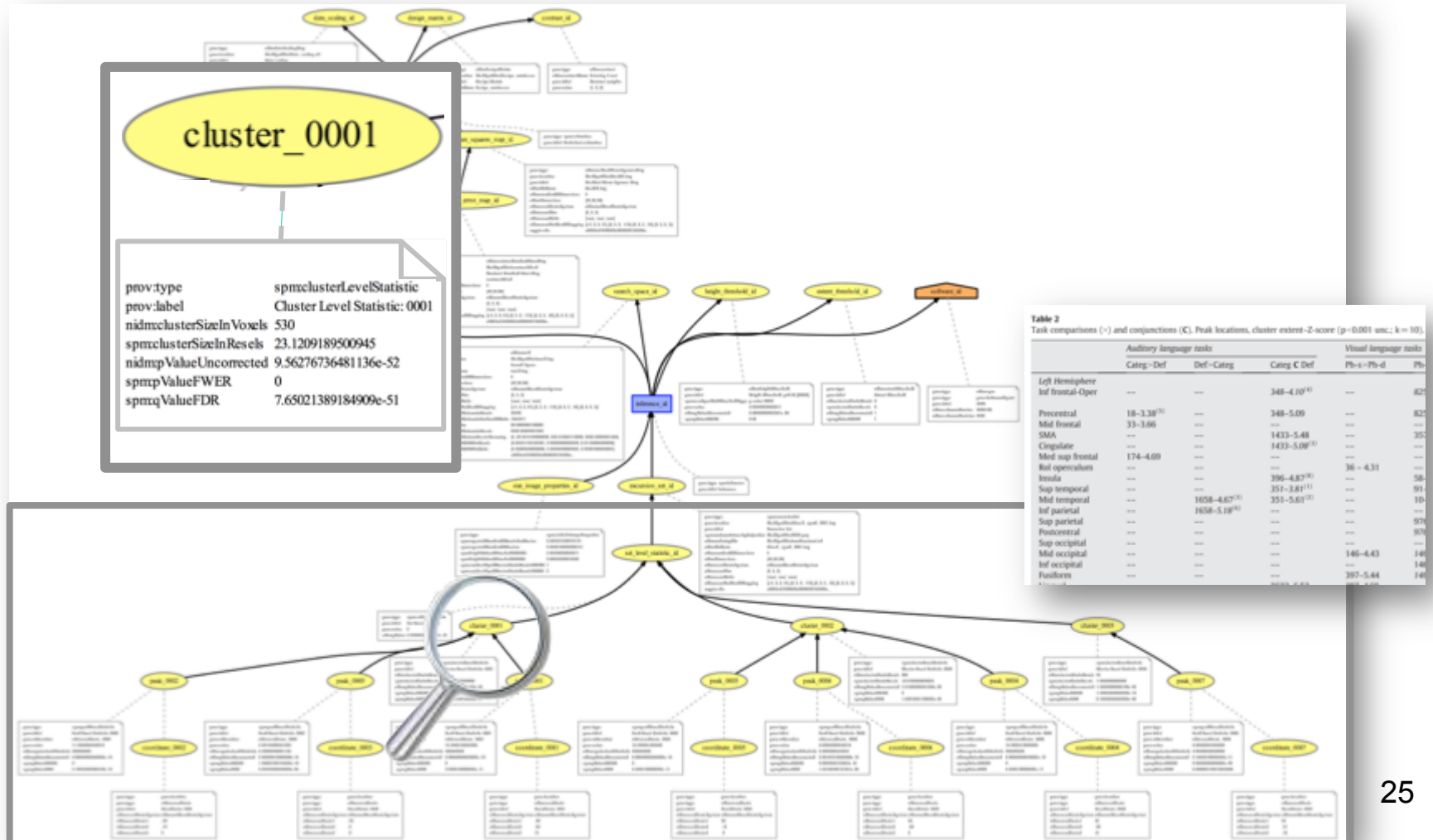- Based on PROV-DM a W3C recommendation to encode provenance.
  www.w3.org/TR/prov-dm/

WARWICK

# Data model

WARWICK

# Data model: activities

WARWICK
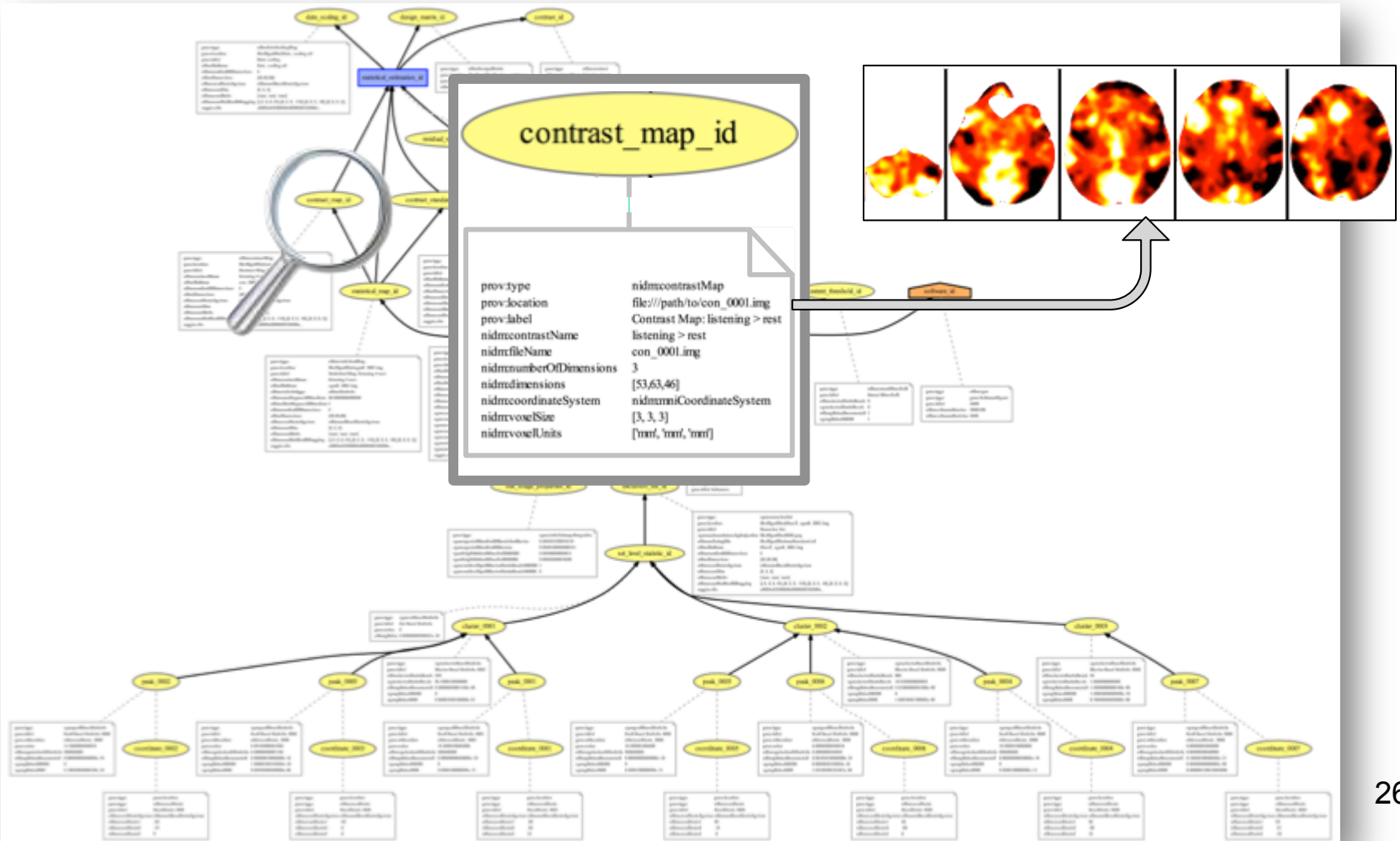
# Data model: agent

WARWICK

# Data model: entities

WARWICK

# Data model: entities

WARWICK

# Conclusion

- Data sharing is one key to reduce skepticism.
- There is already a number of technical solutions for data sharing in neuroimaging.
- A meta-data standard would beneficiate to all of these efforts
  - NI-DM: http://nidm.nidash.org

WARWICK

# Q & A

This work is supported by the **wellcome**trust

WARWICK