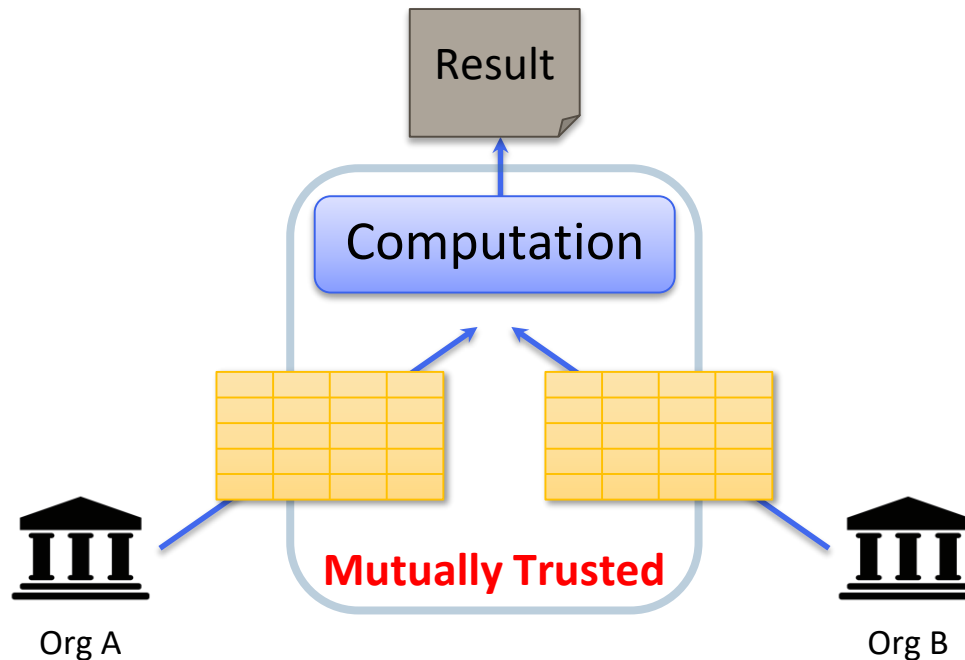# Confidential Computing

**Analytics with data privacy and control**

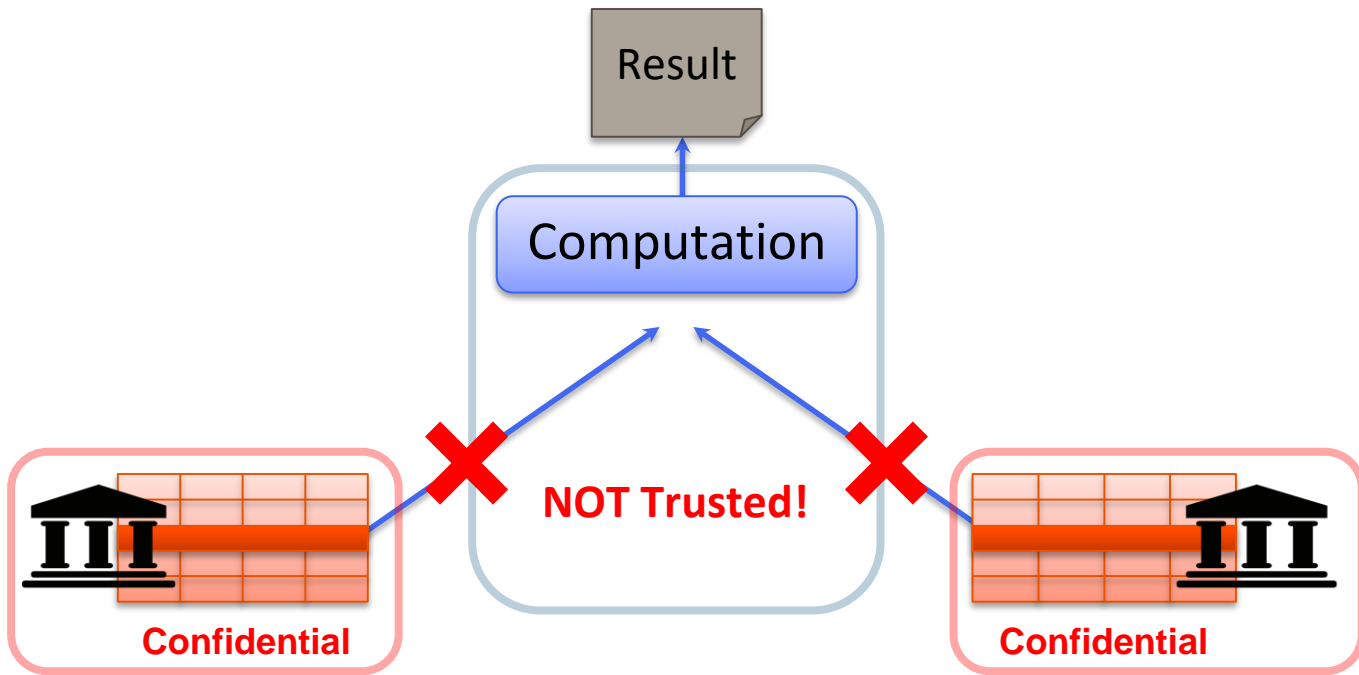Brian Thorne

# Privacy Preserving Linkage Motivation
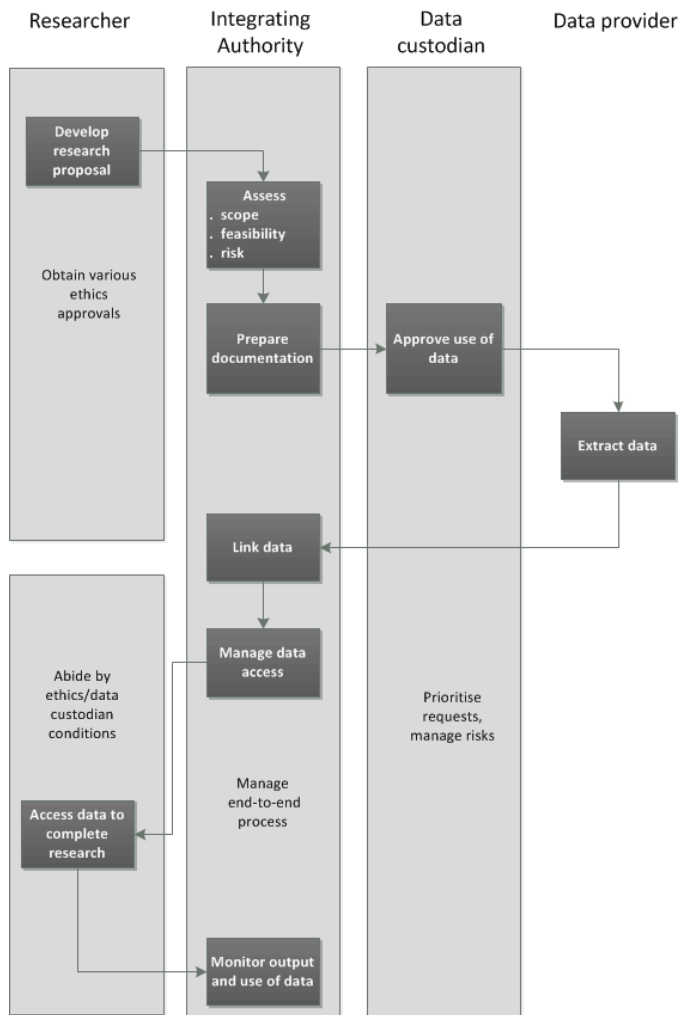
# Multi-Organisation Analytics Today

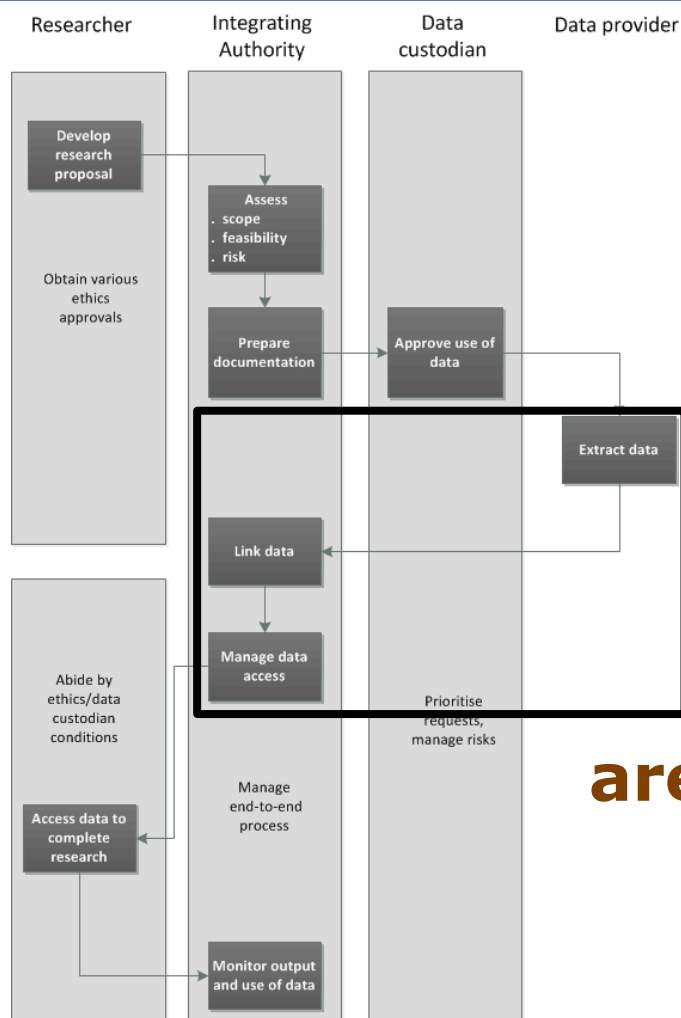# But many Opportunities are Blocked



Confidential Computing | brian.thorne@data61.csiro.au

# Entity Matching

# Overview of a typical data integration project within GOV

**area I'm covering today**

Confidential Computing | Brian.Thorne@data61.csiro.au

# Privacy-preserving entity resolution

- **Goal**: match *corresponding* rows in two distinct databases

| Name | DOB | ... |
|------|-----|-----|
| Klara Jovel | 07/09/1942 | ... |
| Scott Redo | 04/08/1923 | ... |
| Tori Mckone | 07/06/1921 | ... |
| Rusty Brod | 25/07/2014 | ... |

**?**

| Name | DOB | ... |
|------|-----|-----|
| Tori Mckone | 07/06/1921 | ... |
| Scotty Undo | 24/01/1965 | ... |
| Scott Redo | 04/08/1923 | ... |
| Clara Jovel | 07/09/1942 | ... |

- **Constraint**: can't share Personally Identifiable Information (PII)
- **Solution**: *fuzzy & private* matching

# Privacy-preserving entity resolution



One way hash functions      One way hash functions

# How?

For every record we process the PII into a **Cryptographic Longterm Key** or (CLK)

Briefly, we hash the bi-grams for each PII feature into a bloom filter.

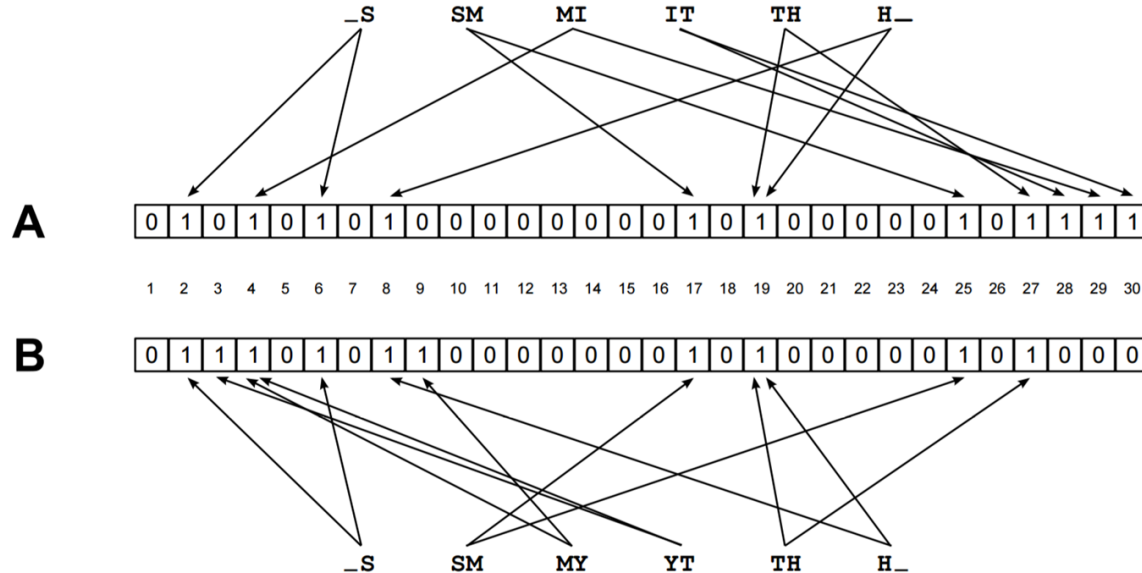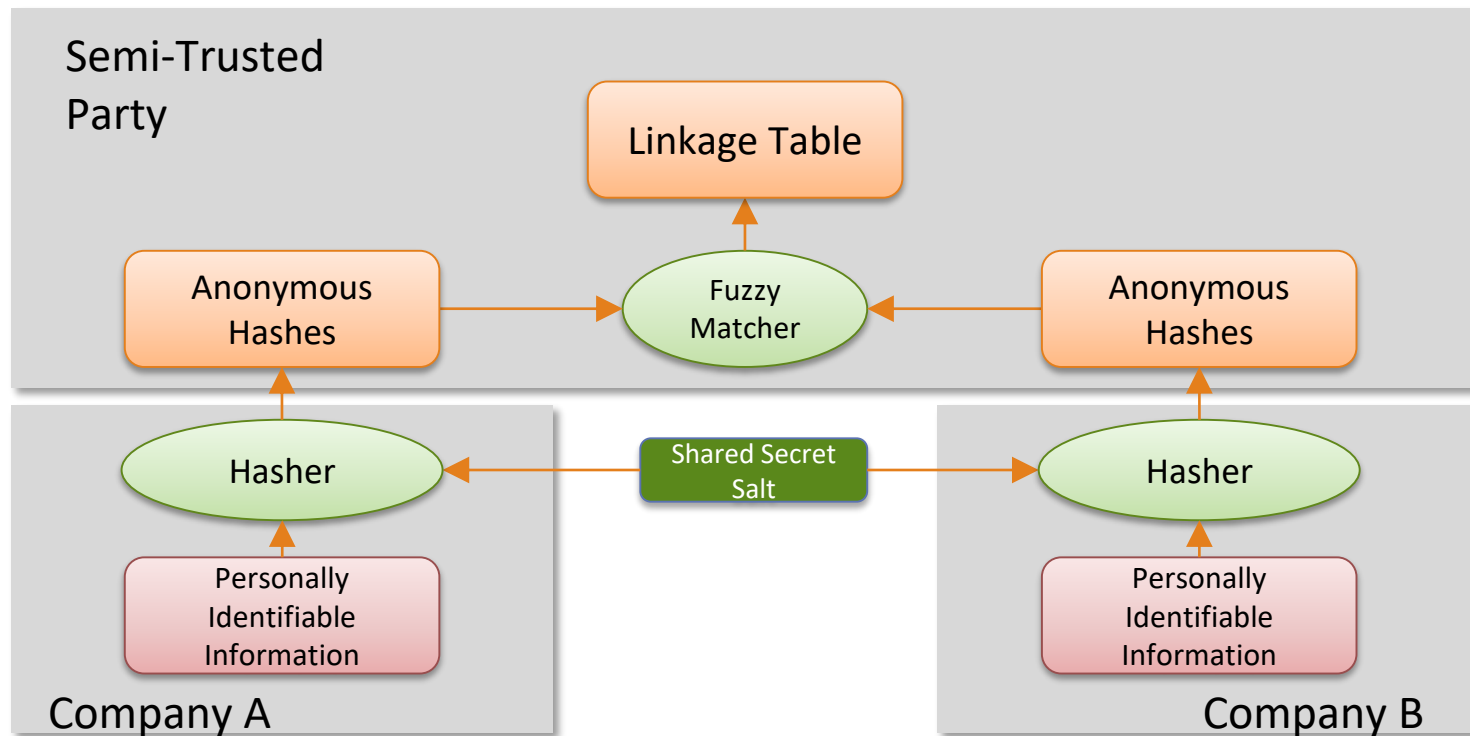https://github.com/n1analytics/clkhash/

# Cryptographic Longterm Key



**Figure 1**: Example for the mapping of two names (SMITH, SMYTH) using bigrams and two hash functions to two Bloom-Filters (A, B) with 30 bits each.

# Private Record Linkage



PII cannot be recovered from the hashes
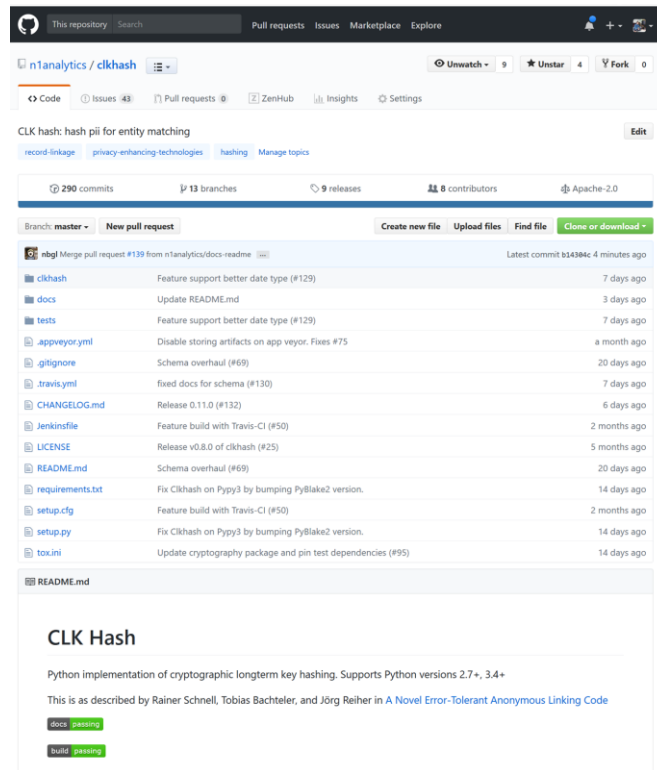
# anonlink

# Semi-trusted Third Party

- Only hashed data is uploaded to the entity resolution service

- Hash security relies on a shared secret between parties

- Implemented the service with a simple JSON + REST API

- All communication is secured with HTTPS

- Authentication tokens created for each job

- Result type and agreed schema is set at beginning

# Client side: Command Line Utility

- Locally hashes PII data

- Creates new mapping jobs on the server

- Uploads hash data

- Retrieves results

```
In [7]: !clkutil hash --help

        Usage: clkutil hash [OPTIONS] INPUT OUTPUT

          Process data to create CLKs

          Given a file containing csv data as INPUT, and optionally a json document
          defining the expected schema, verify the schema, then hash the data to
          create CLKs writing to OUTPUT.

          Use "-" to output to stdout.

        Options:
          -k, --keys <TEXT TEXT>...
          -s, --schema FILENAME
          --help                          Show this message and exit.

In [8]: %%time
        # Hash the data using the secret keys that the linkage authority doesn't know
        !clkutil hash --keys smooth oreo alice.txt alice-hashed.json

        Assuming default schema
        Hashing data
        CLK data written to alice-hashed.json
        CPU times: user 53.3 ms, sys: 16.7 ms, total: 70 ms
        Wall time: 2.23 s
```
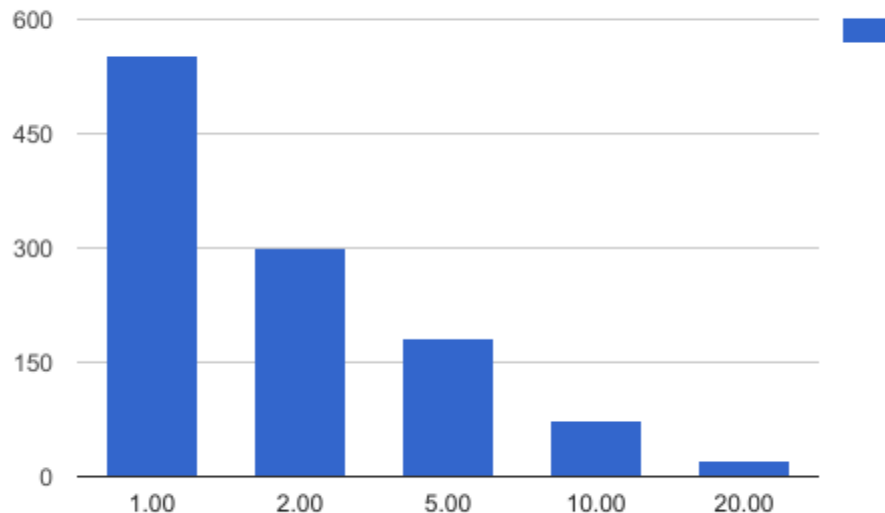
# Performance & Case Study

# Speed and Scale

- 1.3B hash comparisons/s
- Handle uploads of 35M hashes
- 1M x 1M match takes around 5 hours

Running on four r4.4xlarge instances on AWS

**100K Match - Time taken with more workers**

Computing similarity between CLKs is a very parallel problem. Our implementation utilizes multiple workers to carry out comparisons using a kubernetes cluster
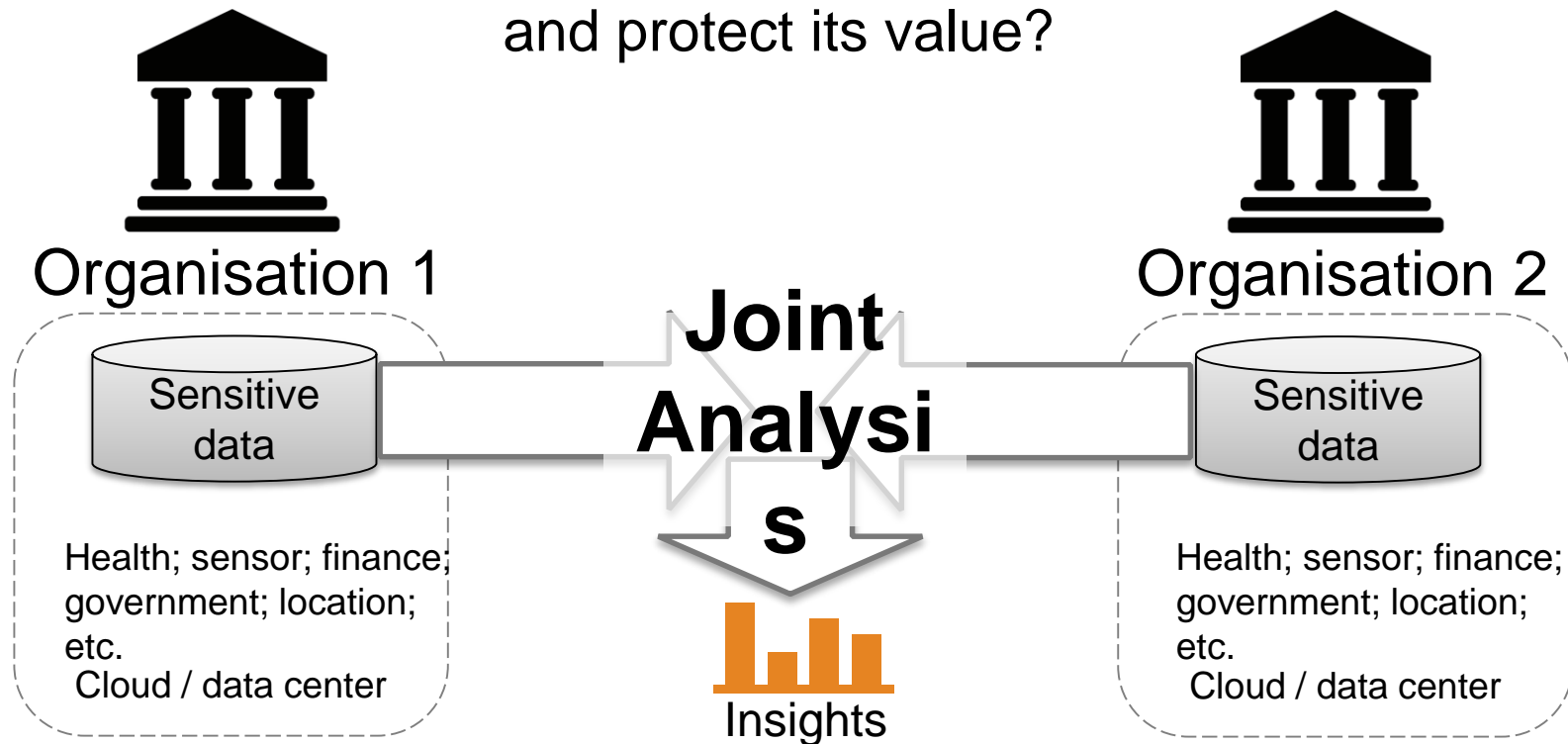
# Data61 Privacy Projects

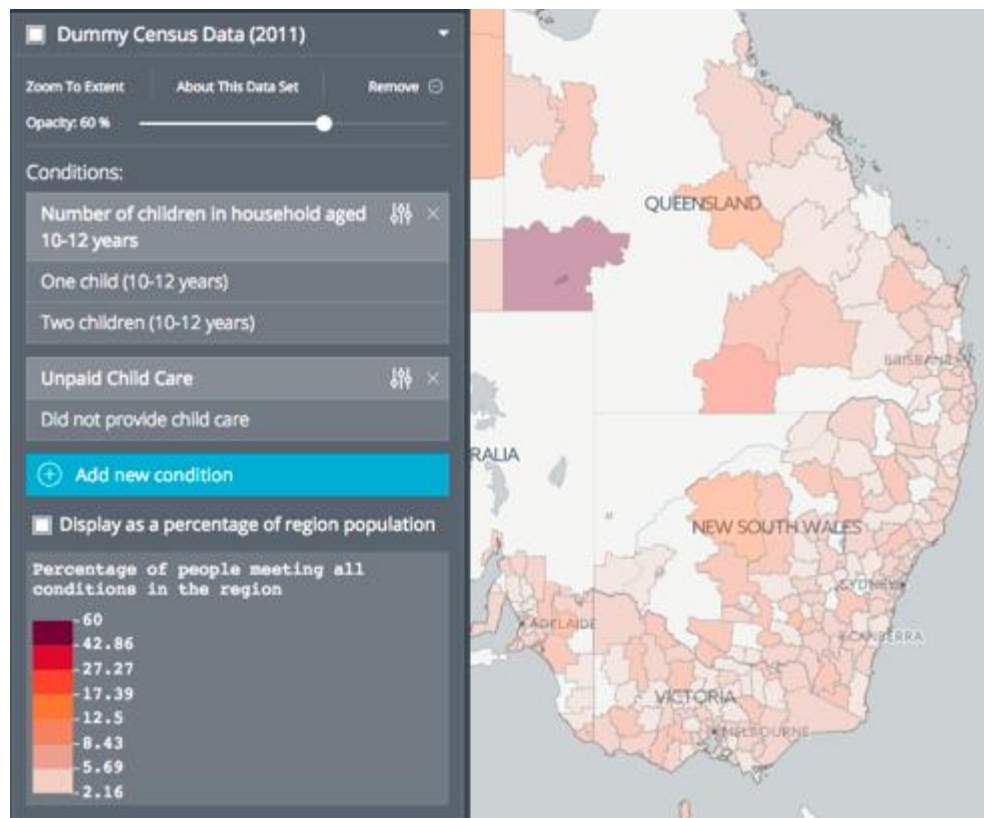**Protari, SENDA, Risk Identification, N1, Private Linkage**

# Confidential Computing

How can we learn insights from data from multiple sources and protect its value?



Organisation 1

Organisation 2

Sensitive data

Sensitive data

**Joint Analysis**

Health; sensor; finance; government; location; etc.
 Cloud / data center

Health; sensor; finance; government; location; etc.
 Cloud / data center

Insights

# Protari

# N1 Analytics

## Goals

:

| | |
|---|---|
| Release your data without losing control | Access data that is currently too sensitive |

## Technologies

:

| | | |
|---|---|---|
| Fully, Somewhat, Partially Homomorphic encryption | Secure Multiparty Compute | Learning from Aggregates |

## Capabilities

:

| | | |
|---|---|---|
| Learn and deploy models | Secure aggregation of data | Clustering/Anomaly Detection |

# Confidential Computing

**Analytics with data privacy and control**

www.n1analytics.com