



Connectedness and Meaning: A New Analytical Platform

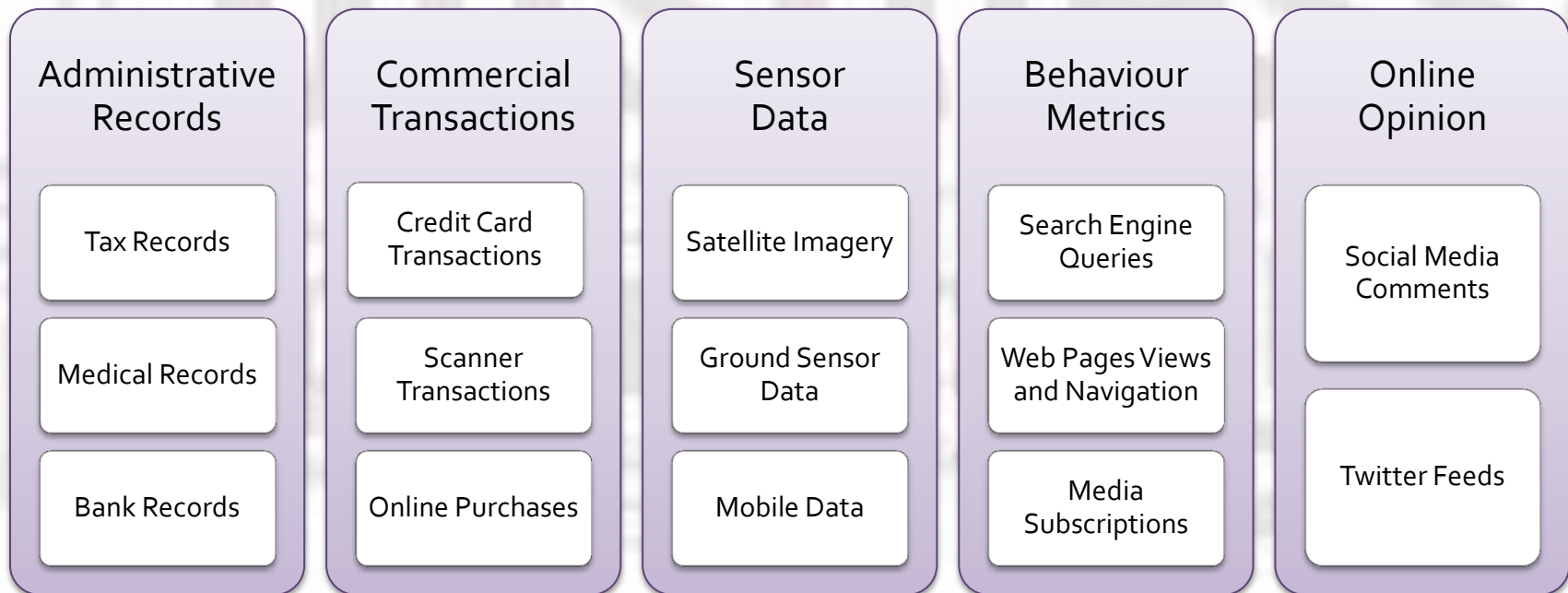
Joseph (Chien-Hung) Chien
Australian Bureau of Statistics

Disclaimer:

The information presented is not for official statistics and only synthetic data is used for demonstration purposes. The opinions expressed in this presentation are those of the presenter, not the ABS

Big Data = Big Data sources

- Administrative records and Scanner Data are not new.
- Sensor Data has potential – but has significant challenge
- Behaviour metrics and online opinion – potentially large inherent statistical biases.



Increasing methodological challenges

Analytical challenges of Big Data



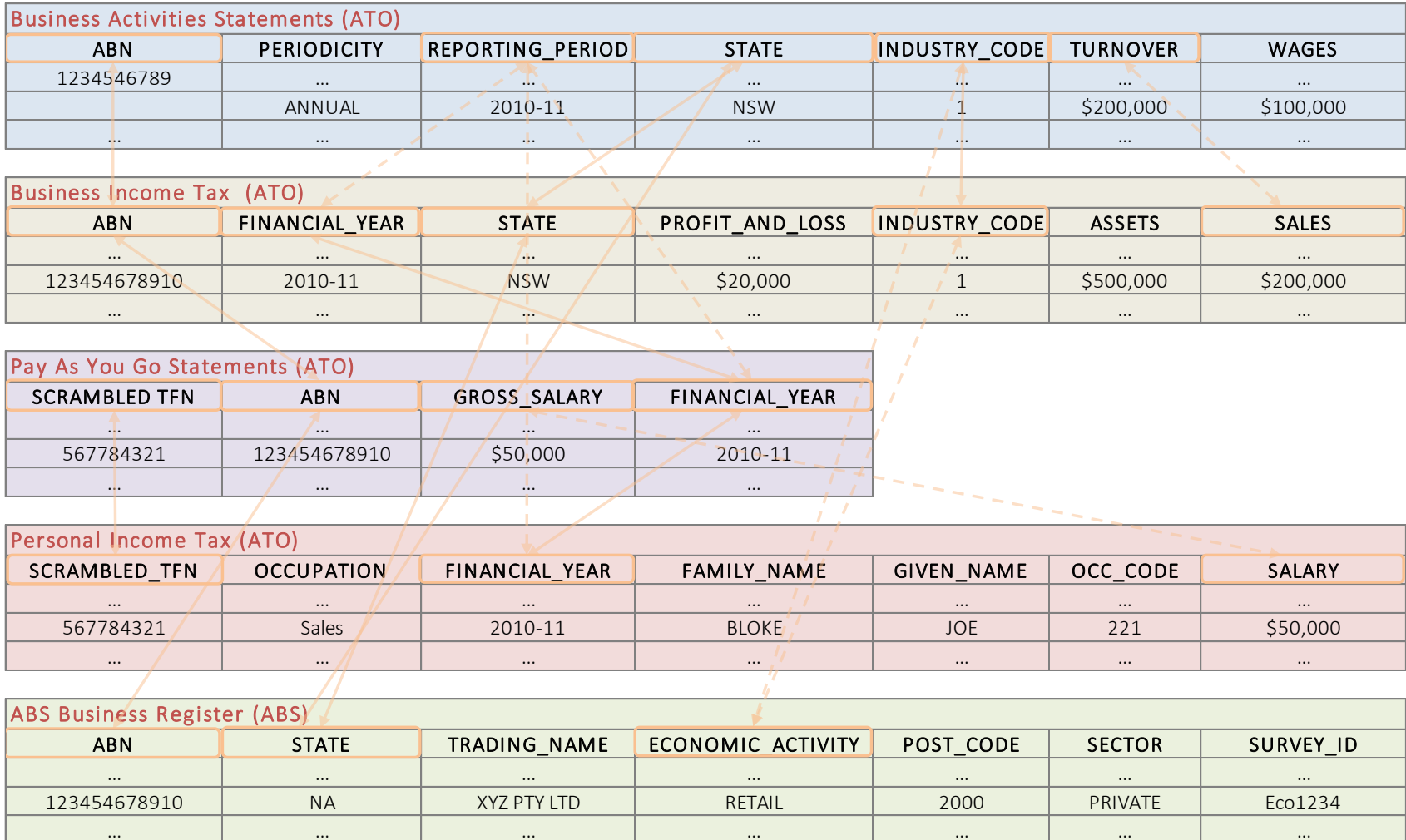
Multisource – conceptual coherence is difficult without a common frame of reference

Multiconnected – a strong need to expand existing record linking methods

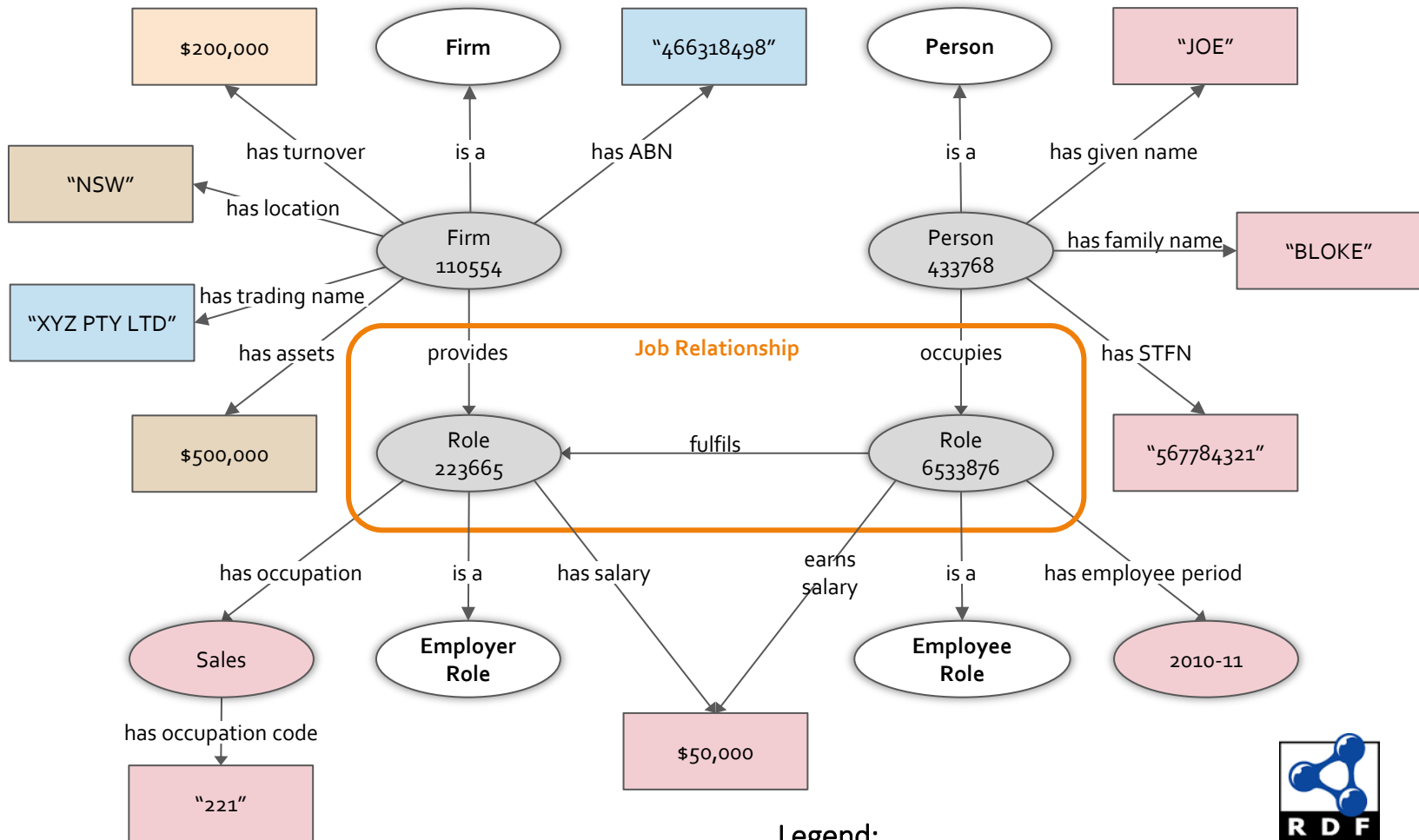
Multistructured – ABS lacks dynamically configurable, schema-last data repositories

Multidimensional – a mismatch between the pattern of data access and how the data is represented and stored

From structured data sets ...



... to a network of entities and relationships



Legend:



W3C Linked Data



Application 1: Research question and Methodology

Research question: Can we distinguish true and spurious firm death events from analysing the network connections in the prototype semantic LEED?

The statistical importance of accounting for true firm deaths.

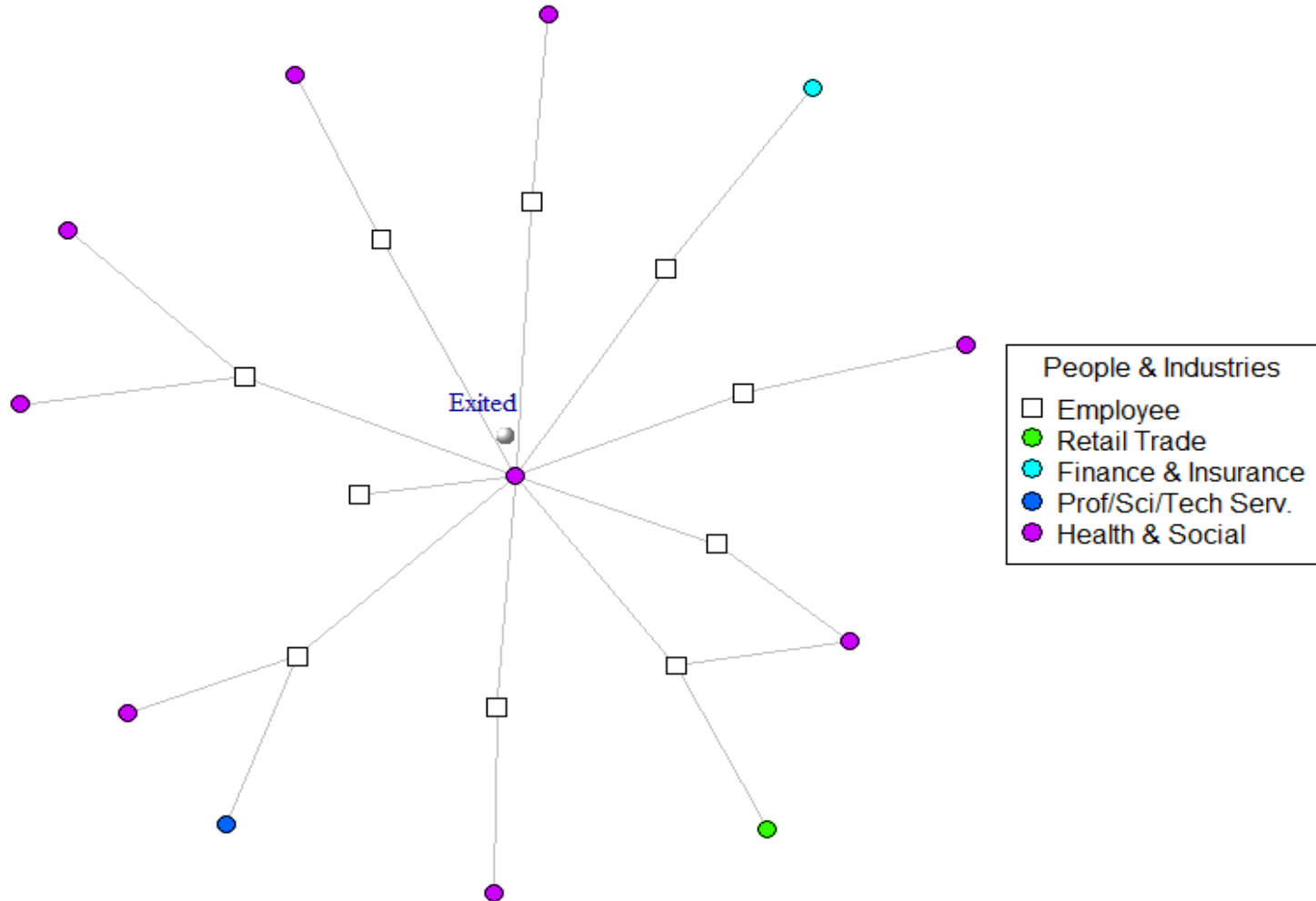
- Exits \neq Deaths (OECD)
- Statistical bias if we don't correct them.

Methodology

- Derive network statistics using semantic LEED (Important).
- Combining Multilevel modelling and Bayesian network.

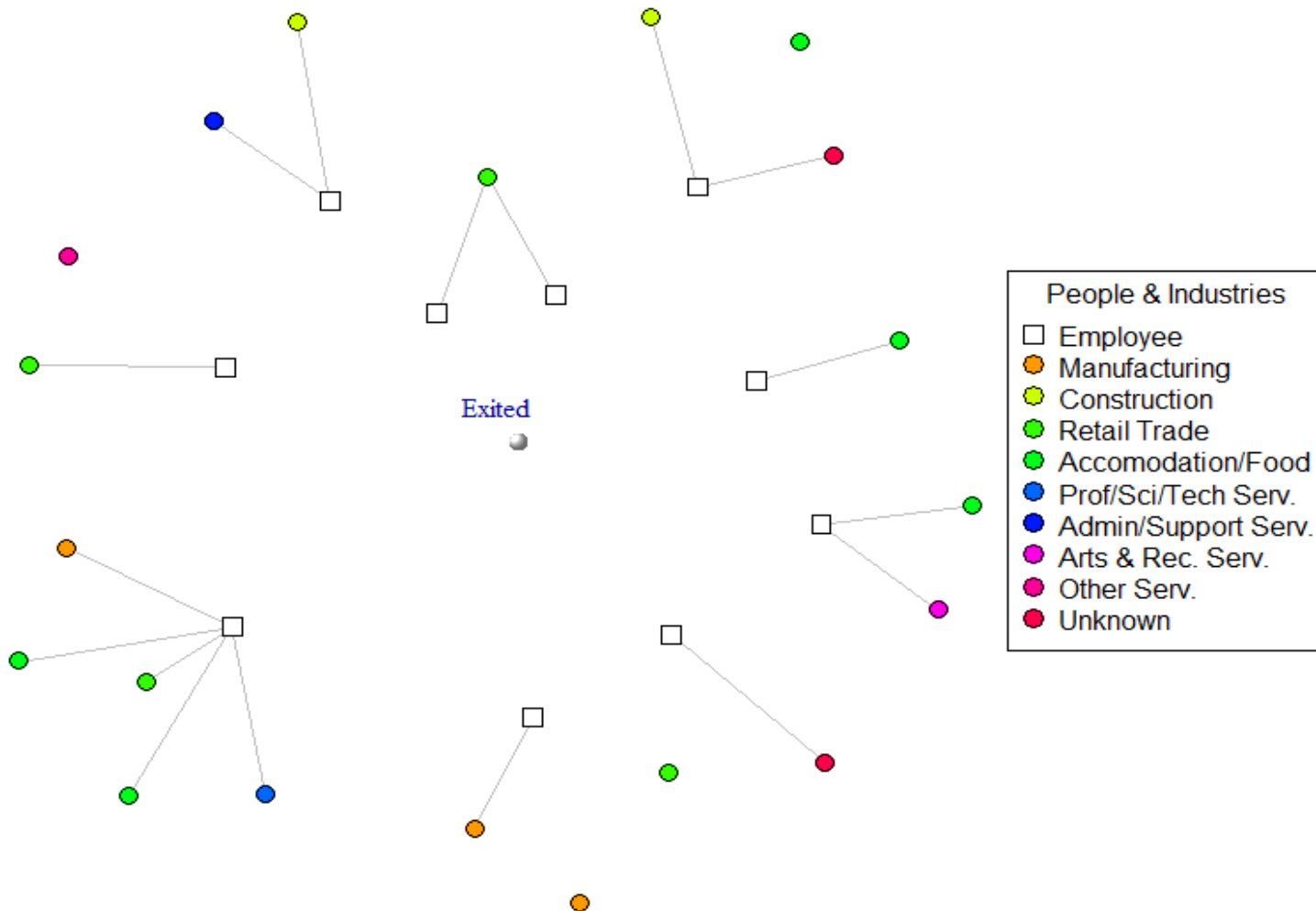
Firm exit: A continuing example

After firm exits

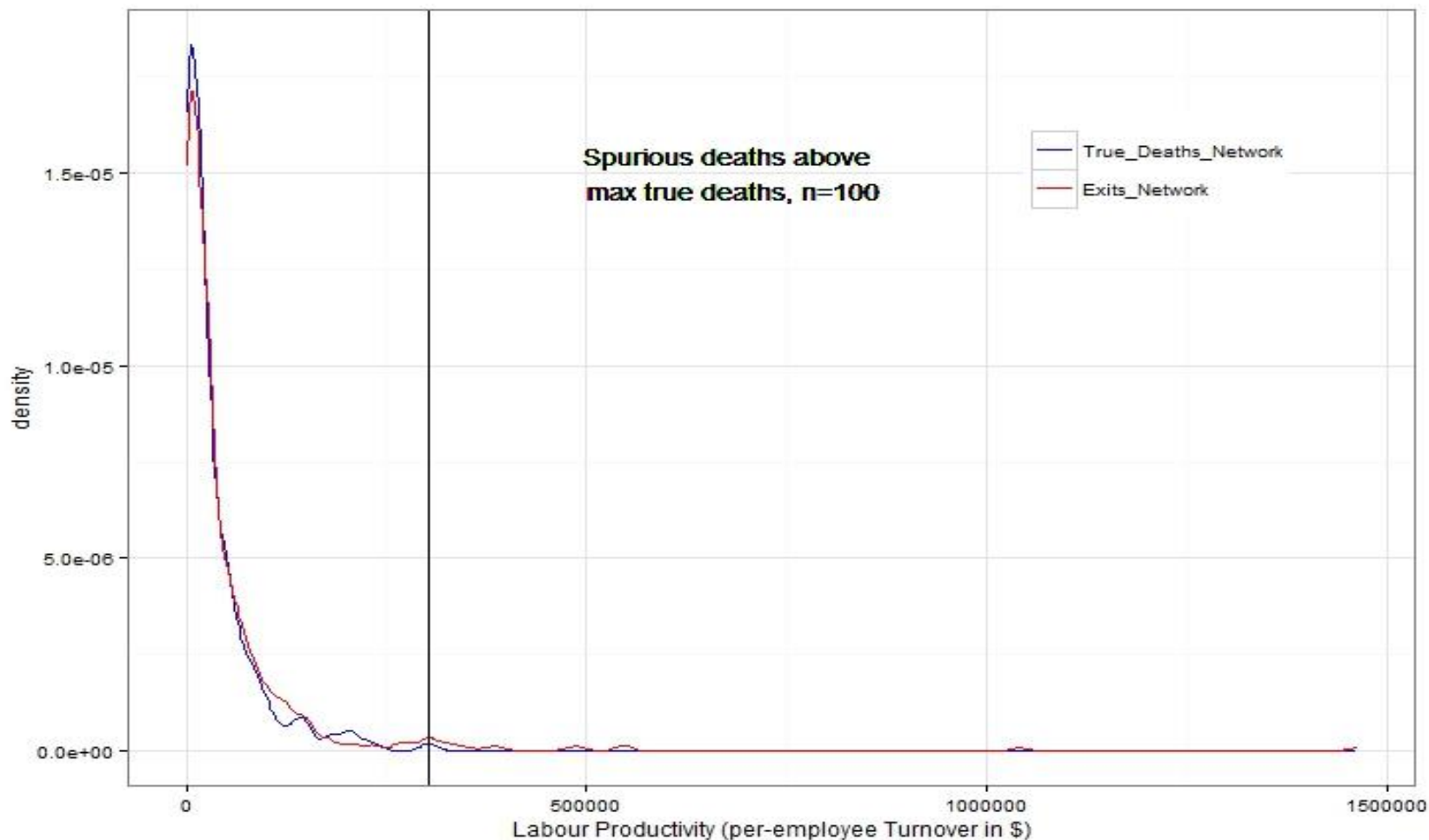


Firm exit: a death example

After firm exits



Summaries and results – bias without correction



Application 2: DEMO

Advantages for Official Statistics

- ✓ Consistency in the use of statistical concepts across data collections;
- ✓ Linking multiple disparate datasets for different analytical perspectives
- ✓ Integrating dynamic structured and unstructured content
- ✓ Manipulating highly multidimensional data in statistical computation
- ✓ Fast and adaptive information discovery on the scale of Big Data
- ? **Research direction** – how can we use this approach to enhance existing data linking methods in the ABS?



Questions?

Email: joseph.chien@abs.gov.au
chien-hung.chien@anu.edu.au

