

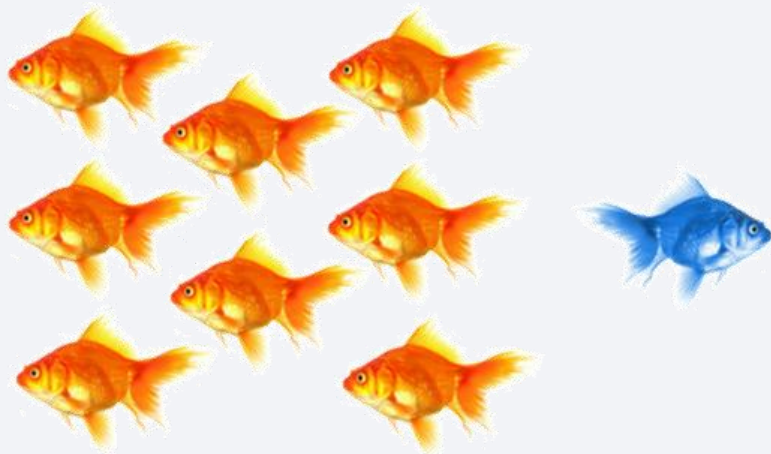
Use of Knowledge Graphs and Relational Machine Learning



Frederic Clarke, Director MINDS
Australian Bureau of Statistics

Australian Bureau of Statistics
Informing Australia's important decisions





Multi-disciplinary team – maths, CS, KM

- Focus on hard problems, learn by doing
- Demonstrate solutions through prototypes

Investigate, experiment, evaluate and inform

- Methods, technologies and models
- New data sources and applications
- ABS strategic capabilities and priorities
- Environmental trends, opportunities, threats

Agenda

The strategic context

A motivating example

ABS use of KL and RL

GLIDE and current work



Strategic context

Inform Australia's important decisions

- Producing new and relevant statistical insights
- Enabling effective and safe use of data
- Building national information capability for the future

On public policy, services and investment


Disruptive change

Driven by powerful global megatrends

- Intelligent Machines
- Digital Connectedness
- Data-driven World

Impact on government is profound

- Overturns extant business models
- Challenges traditional decision-making processes



Lee Sedol and AlphaGo Zero

Complex systems

Most economic, social and environmental systems are complex

- Many interacting entities of different types
- Dynamic and non-linear relationships
- Emergent system structure and behaviour

Connectedness is an essential condition for complexity

Complex systems underlie most 'wicked' problems

Big data



Personal
Devices



Imagery
Systems



Smart
Meters



Product
Scanners



Environmental
Monitors



Telematics
Devices

Diverse new sources of human-generated and machine-generated data

Can be used for statistical purposes

- Improve quality, time-to-market and cost-effectiveness of mainstream statistical products
- Create new statistical products that fill information gaps

Issues: heterogeneity, bias and dimensionality

ABS operates in an increasingly congested and contested environment

Web
Applications



Social Media
Services



Logistics
Systems



Accounting
Systems



Administrative
Collections

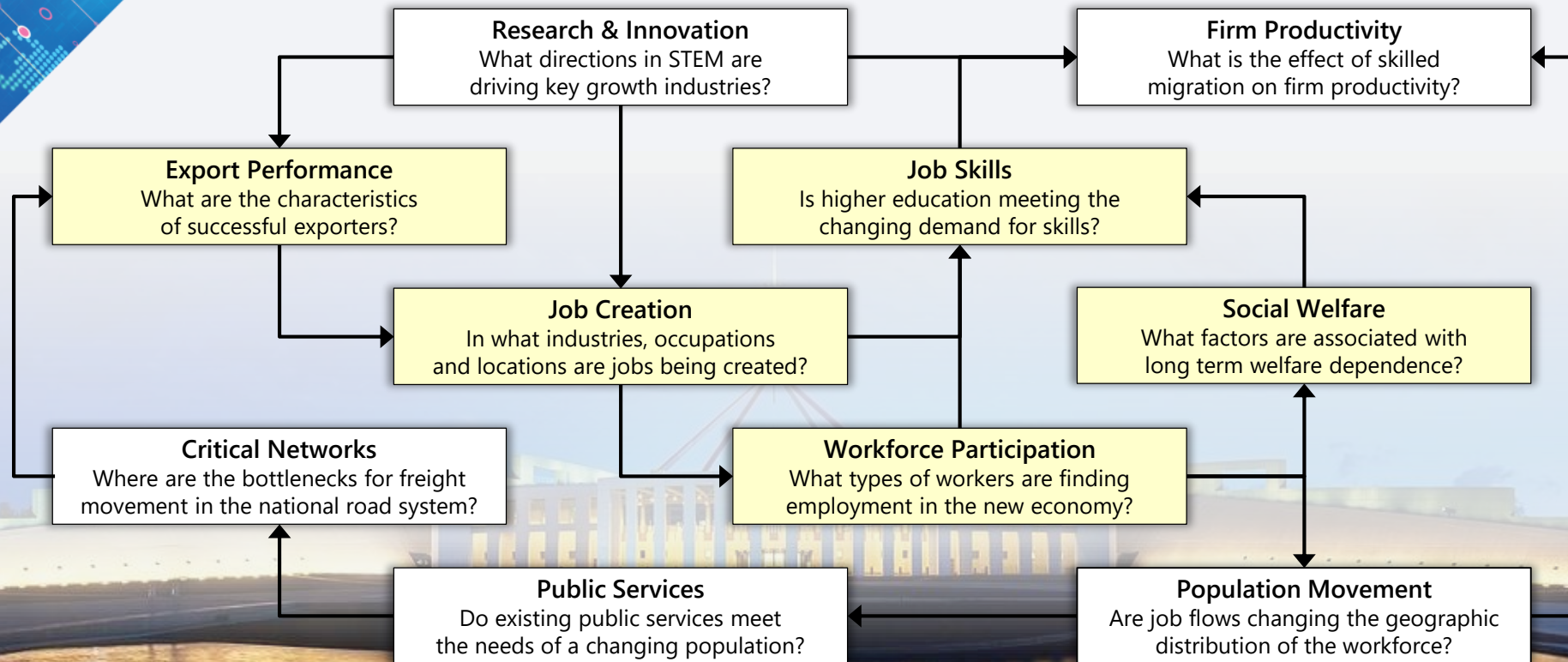


Surveys and
Censuses



Motivating example

Connectedness of policy concerns



Next generation analytical capability

Built on system-centric information models

- Composable, interpretable and semantically precise
- Join up interrelated concept and data spaces
- Connect individuals and groups in multiple ways

Dynamically integrates heterogeneous multisource data

- Structured and unstructured
- Cross-sectional and longitudinal temporal linkages

Next generation analytical capability

Enables multiple analytical perspectives and objectives

- Exploration (pattern sensing) – finding statistical features and correlations
- Explanation (model building) – testing hypotheses about the observed data
- Extrapolation (system simulation) – projecting beyond known cases

ABS use of KG and RL

System is depicted as a graph of entities and relationships

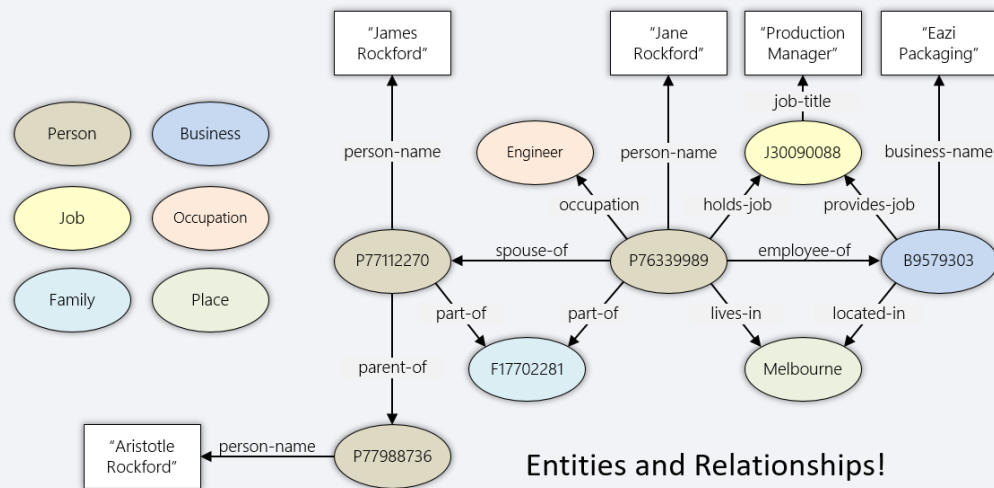
- Entity – individual thing or group of things
- Relationship – association between entities
- Entities interact through relationships of analytical interest

Use W3C Semantic Web formalism for knowledge graphs

- Graph composition (standard: RDF)
- Semantic modelling (standard: OWL)
- Knowledge discovery (standard: SPARQL)



Knowledge graphs – simple example



Systems are partitioned the into context-specific analytical domains

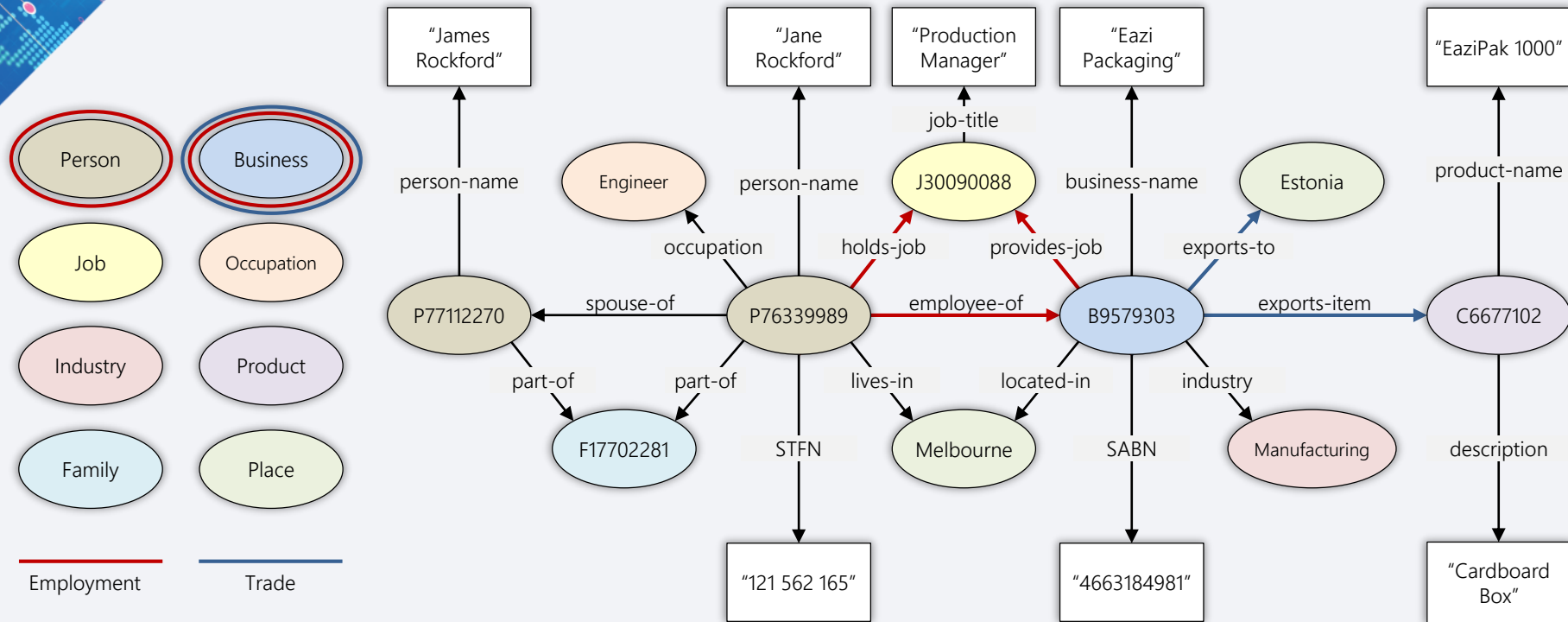
- Example: Trade, Employment, Production, Education, Welfare, etc

Each domain has set of associated entity types and relationships

- Basic entity types can exist in multiple domains
- Relationships are usually context-specific and so bound to one domain

Domains are connected through common entity types

Analytical domains – simple example

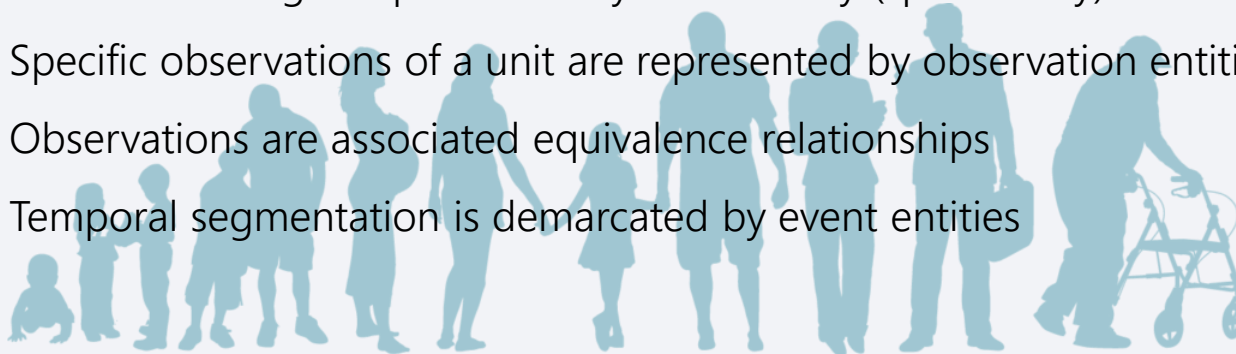


Observable characteristics of an entity can change over time

- Example: person name, address, marital status, and employment details

Entities are associated with their respective observations in data

- Real world thing is represented by a unit entity (spine entity)
- Specific observations of a unit are represented by observation entities
- Observations are associated equivalence relationships
- Temporal segmentation is demarcated by event entities



Associating observations – example

Example: observations of the same person in different data sets

- Record-3 is much later than Record-1 and Record-2
- Significant events: Change of Residence, Marriage, Graduation

Record-1

Family Name	Given Name	Address	DOB	Country of Birth	Sex	Marital Status	Occupation	STFN
Smith	Jane	1 Long Street Broadford VIC	05-08-1985	Australia	F	Single	Student	

Observation date

10-02-2005

Record-2

Family Name	Given Name	Address	DOB	Country of Birth	Sex	Marital Status	Occupation	STFN
Smyth	Jane	1 Long Street Broadford VIC	05-08-1985		F	Single	Student	121 562 165

17-05-2005

Record-3

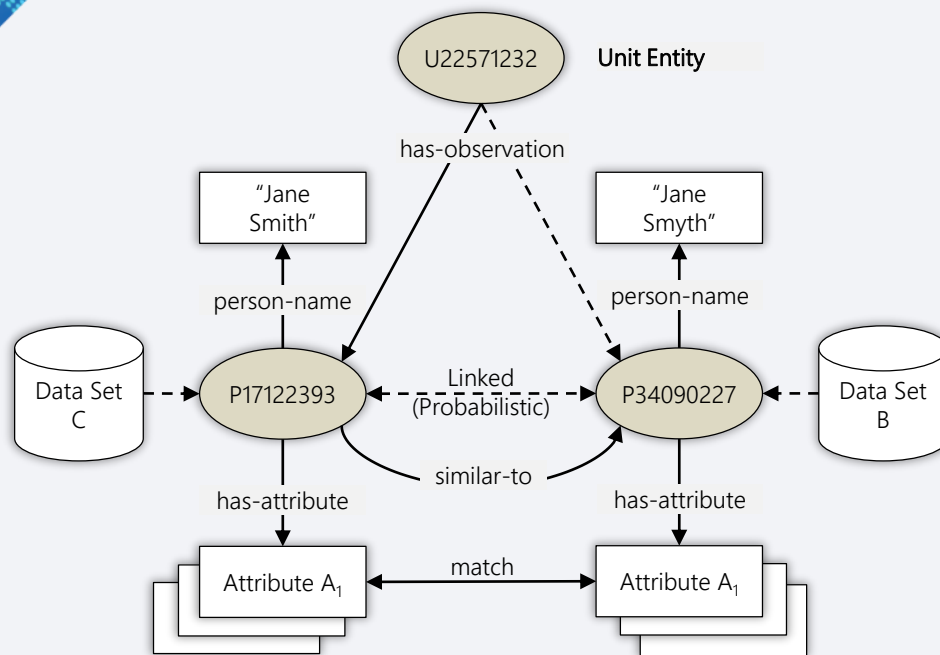
Family Name	Given Name	Address	DOB	Country of Birth	Sex	Marital Status	Occupation	STFN
Rockford	Jane	32 King Street Lalor VIC	05-08-1985	Australia	F	Married	Engineer	121 562 165

15-07-2015



- Current: identifier match using common unique key
- Future: fact match using deductive rules (FOL)
- Multiple class inheritance

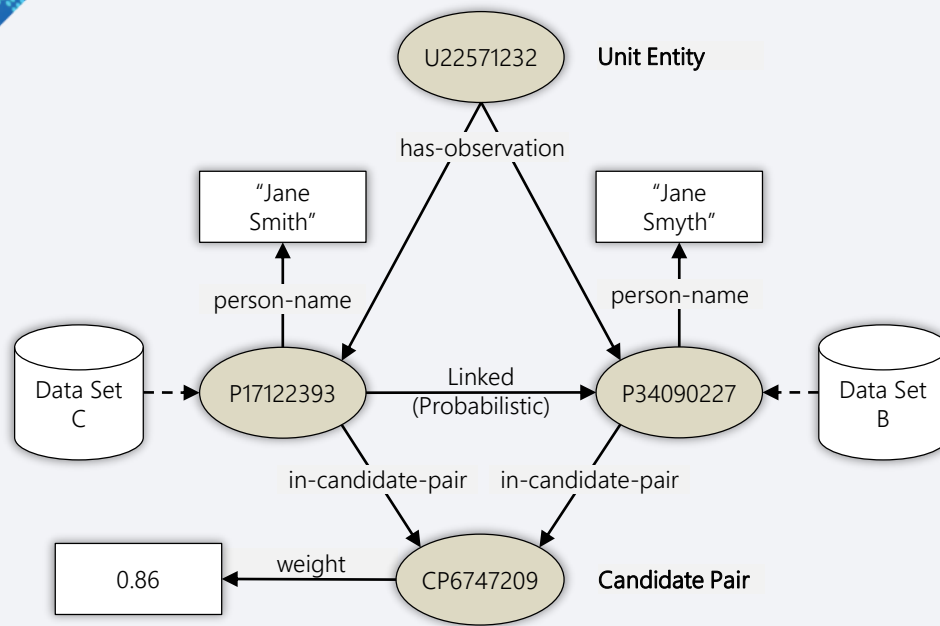
Associating observations – example



Probabilistic association

- Current: similarity match on entity characteristics using Sunter-Fellegi model
- future: similarity match on entity characteristics evaluation or relational characteristics using machine learning
- Multiple class inheritance

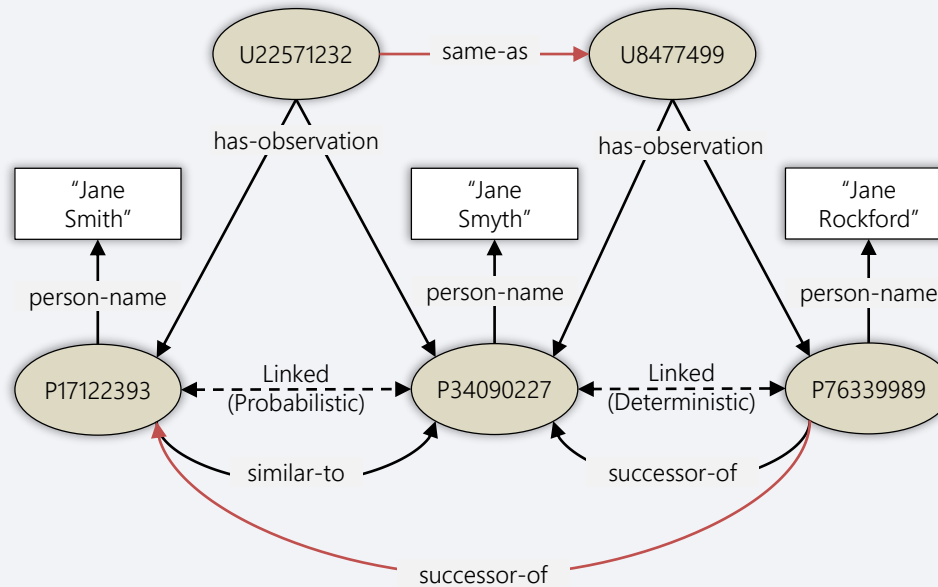
Representing weights/probabilities



Accommodates SF pairwise linking

- Candidate pairs with associated weights
- An observation can be linked to more than one CP
- Much better represented in a property graph construct
- Can we combine the two paradigms?

Associating observations – example



Combining data

Extract entities and relationships from source data

- Automatically by content analysis (fact extraction) tools

Create the knowledge graph in SW form

- Asserted facts from source data

Associate observation entities in the graph

- Inferred facts from reasoning processes



Based on pattern of connections among entities

- Extend scope of probabilistic linking beyond IID assumption

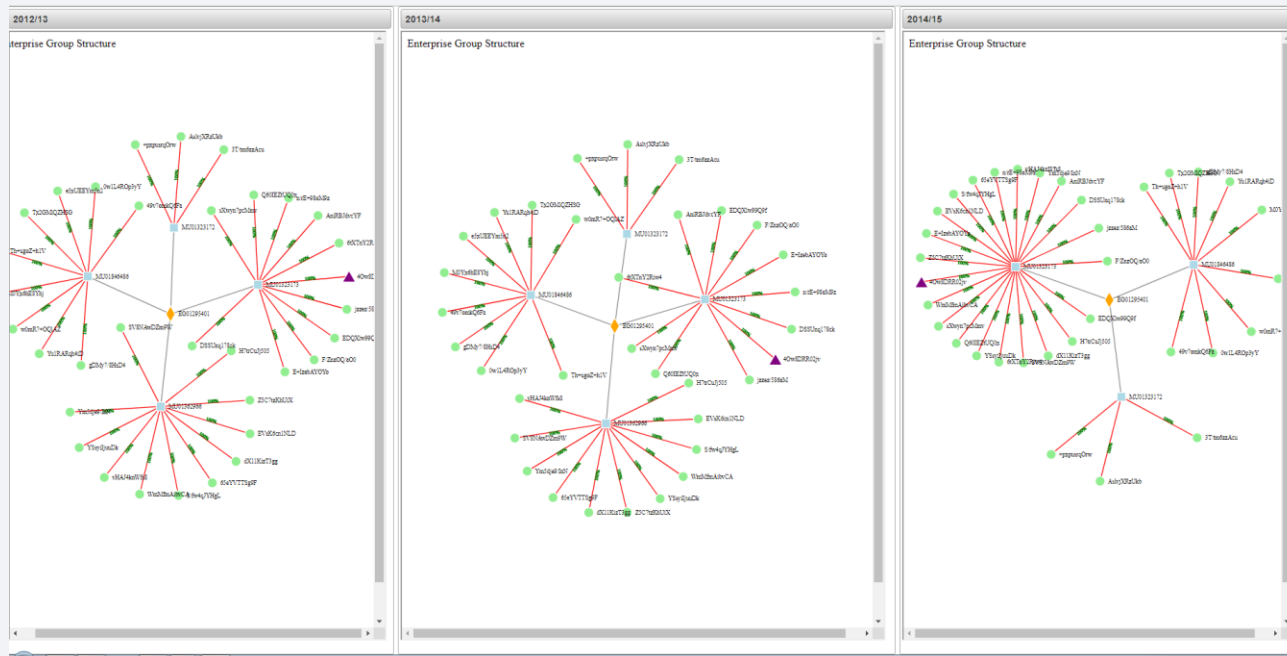
Associate entities across time and in disparate data sources

- Needed when there are no reliable common identifiers
- Example: persons by family and household relationships

Detect events that involve changes in the structure of groups

- Example: business reorganisation, closure, takeover

Detecting change over time



RL approaches

Graph kernel learning

- Represent the structural form of a graph for use in kernel learning algorithms
- Kernels: Weisfeiler-Lehman (WL), Intersection Tree Path (ITP)

Tensor factorisation

- Manipulates 3-order adjacency tensor of the knowledge graph
- Estimate probability distribution over possible states of graph
- Algorithms: RESCAL, Complex Embedding, HOLE, TransE

Overview of GLIDE

Capability vision for enabling informed decision making

- About policy, services, investment and (possibly) regulation

Based on insights derived from different types of analysis

- Exploratory analysis, hypothesis testing, system simulation

Using a dynamic evidence base from diverse data sources

- Surveys, admin collections, sensors, transactions, web content, ...

Graph-Linked Information Discovery Environment

What GLIDE will provide

One platform for data analysis and linking

- Browser interface – rich, interactive, navigable, context-sensitive visualisation
- Program interface (R, Python) – statistical and econometric modelling
- Spatial, temporal and compositional perspectives of problem
- Deterministic and probabilistic linkage methods – common key, Sunter-Fellegi
- Extraction of entities and facts from heterogeneous data
- Reusable, plug-and-play models and components

Graph-Linked Information Discovery Environment

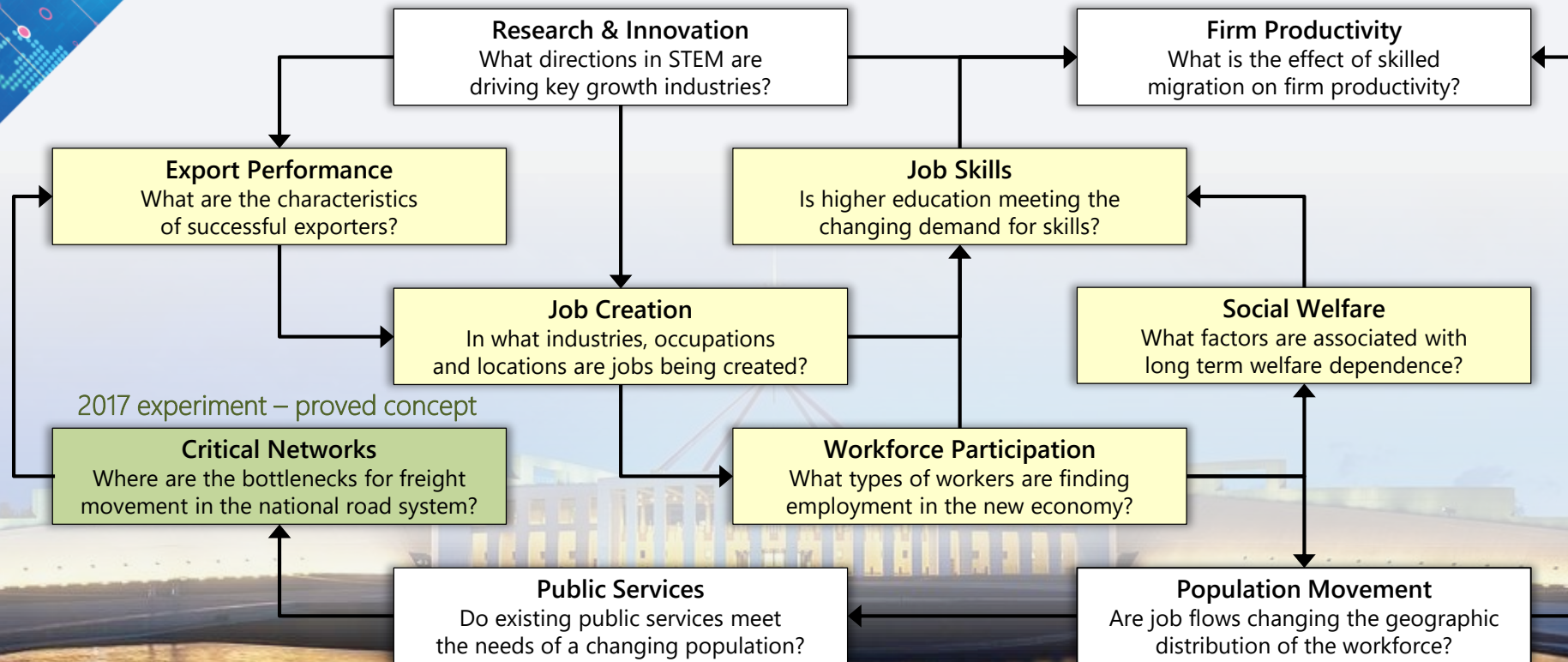
What GLIDE will provide

A set of extensible, interoperable 'data pools'

- Multiple structured or unstructured data sets
- Longitudinally linked data (given social license)
- Different entity and relationship types in the form of a knowledge graph

Graph-Linked Information Discovery Environment

Pathway initiatives



Pathway initiatives – freight movement

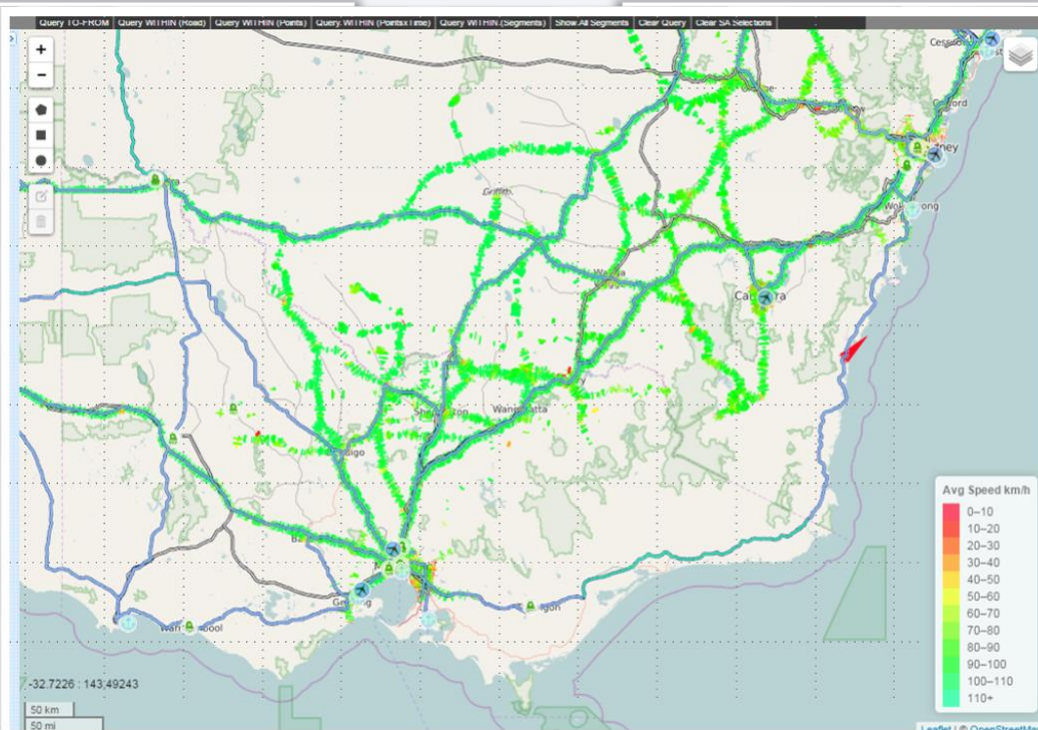
Research
What are the
characteristics
of successful
exporters?

Export Performance
What are the characteristics
of successful exporters?

In what
ways and locations
are freight
movements
occurring?

Critical Networks
Where are the bottlenecks for freight
movement in the national road system?

Do existing
infrastructure
meet the needs
of a changing
population?



What is the
distribution of
the workforce?

Pathway initiatives – freight movement

Research & Innovation

What do
driving k

Firm Productivity

Export Performance
What are the characteristics
of successful exporters?

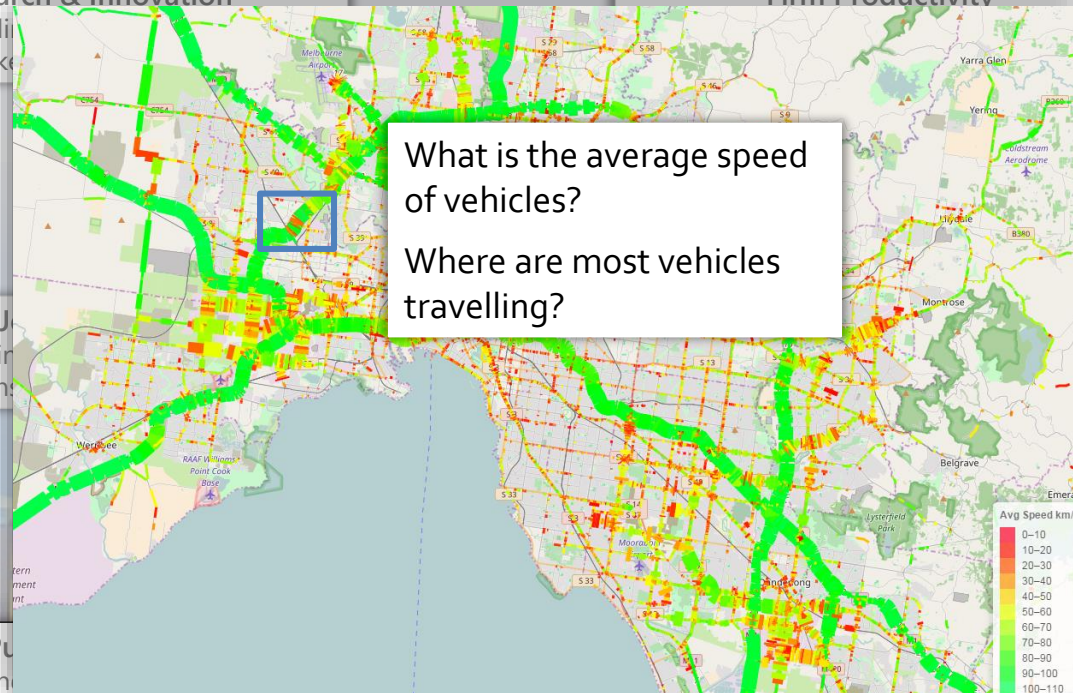
In what in
and location

What is the average speed
of vehicles?
Where are most vehicles
travelling?

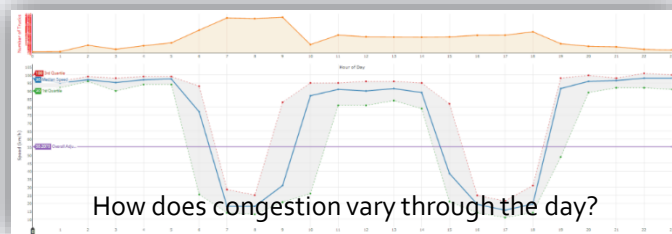
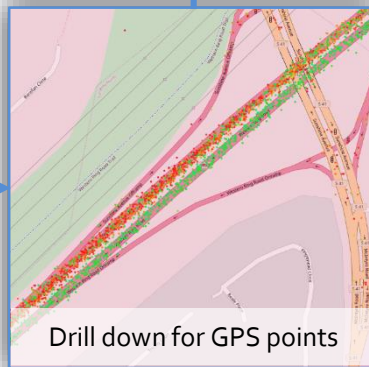
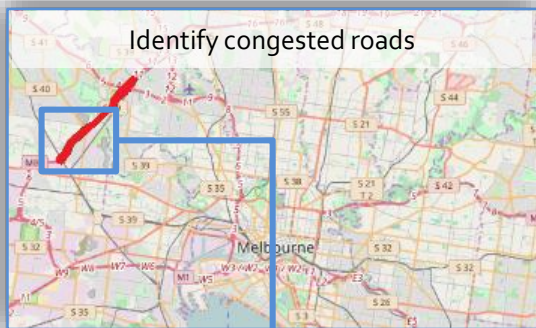
Critical Networks
Where are the bottlenecks for freight
movement in the national road system?

Do existin
the needs of a changing population?

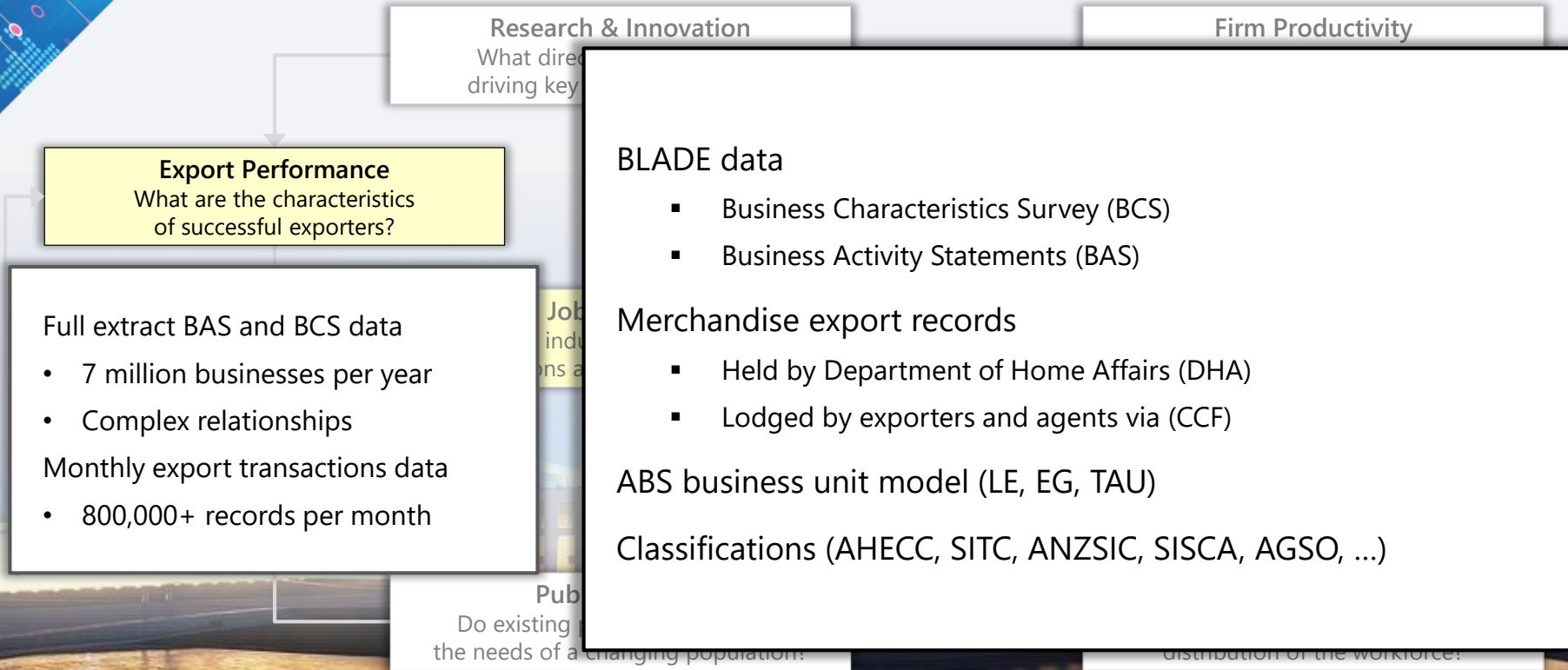
distribution of the workforce?



Freight movement – drilling down



Pathway initiatives – successful exporters





Questions?