



Search Engine Topics

Will Poynter • CLOSER



CLOSER:

Maximise the use, value and impact
of the UK's longitudinal studies, both
at home and abroad

CLOSER Discovery

- Launched January 2017
- DDI-3.2
- Colectica
- 9 longitudinal studies



Searching

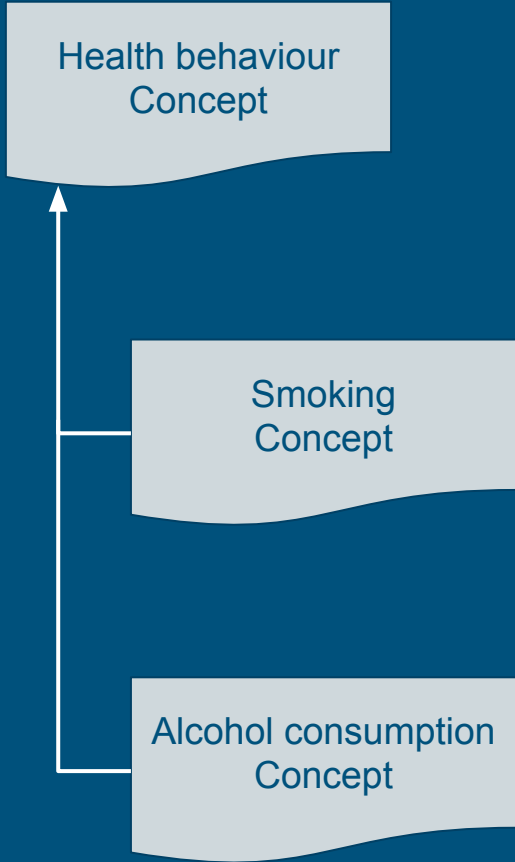
alcohol|

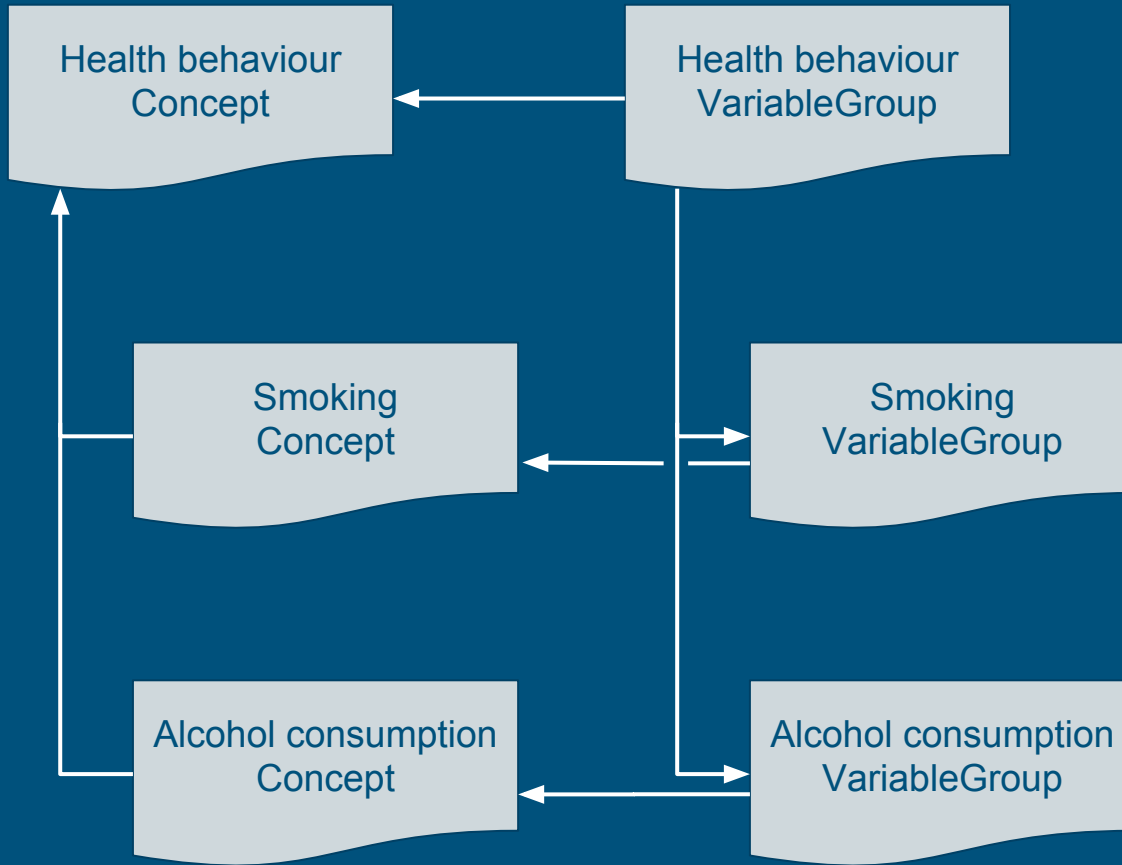
- Health behaviour (Health and lifestyle) (7,443)
 - Diet and nutrition (3,472)
 - Physical activity (888)
 - Sleep (380)
 - Smoking (445)
 - Alcohol consumption (764)
 - Substance abuse (665)
 - Risk taking (88)
-

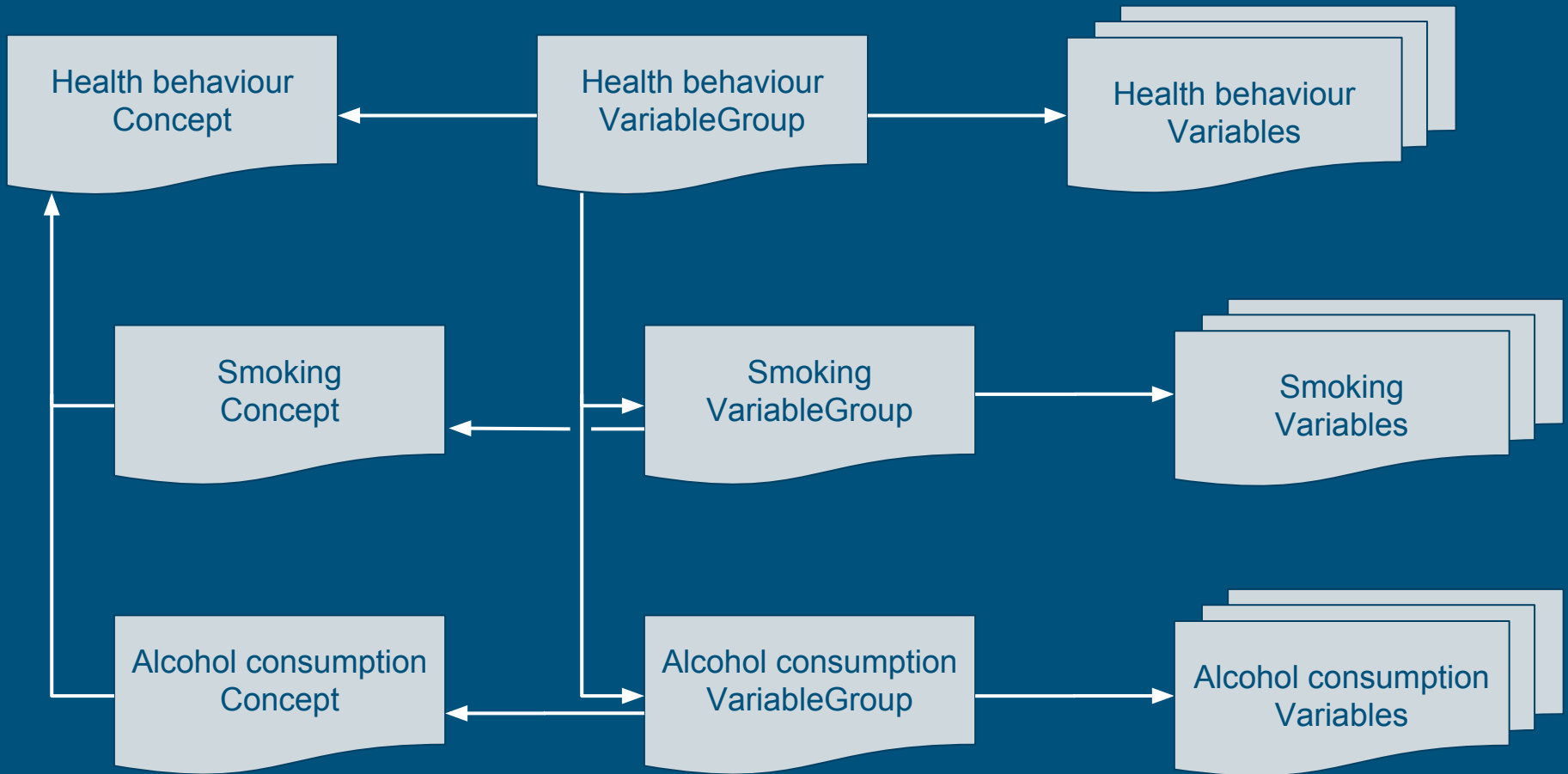
Filtering

DDI Concepts

Health behaviour
Concept







Scale

118

CLOSER Topics

51,000

Questions Constructs

98,000

Variables

How?

Strands

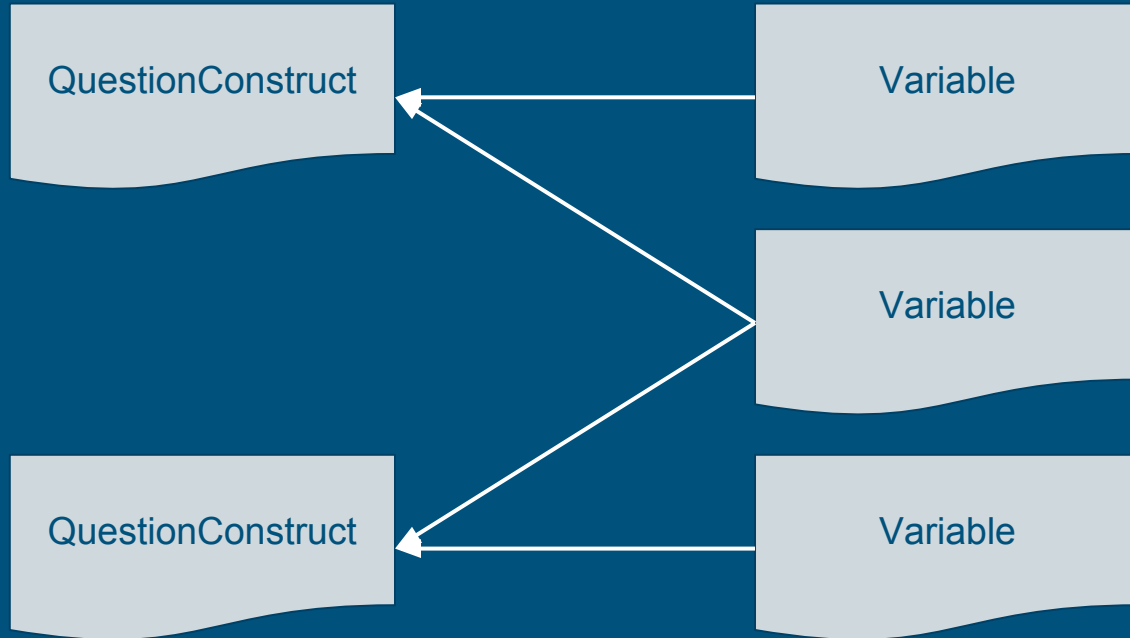
- Collection of QuestionConstructs and Variables
- Joined via **SourceQuestionReference**
- Must have the same topic

Strand



Must have the same topic

Strand

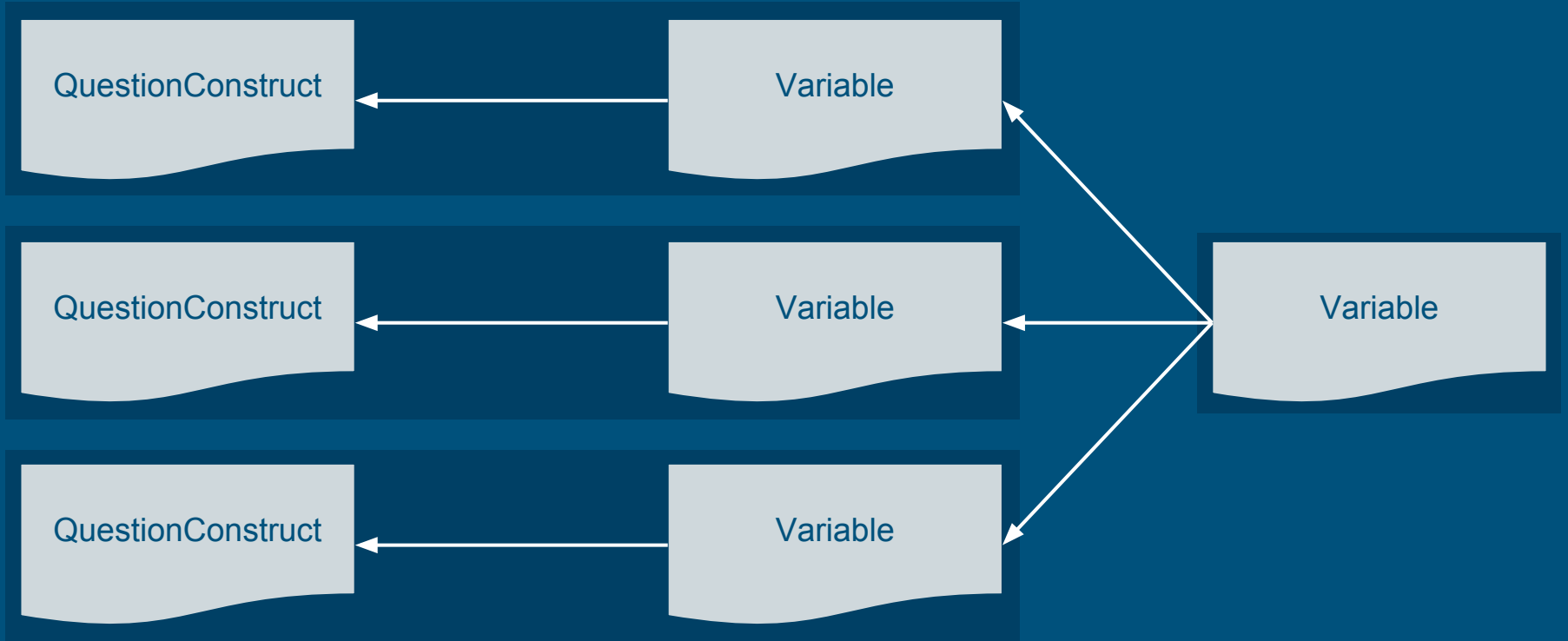


Must have the same topic

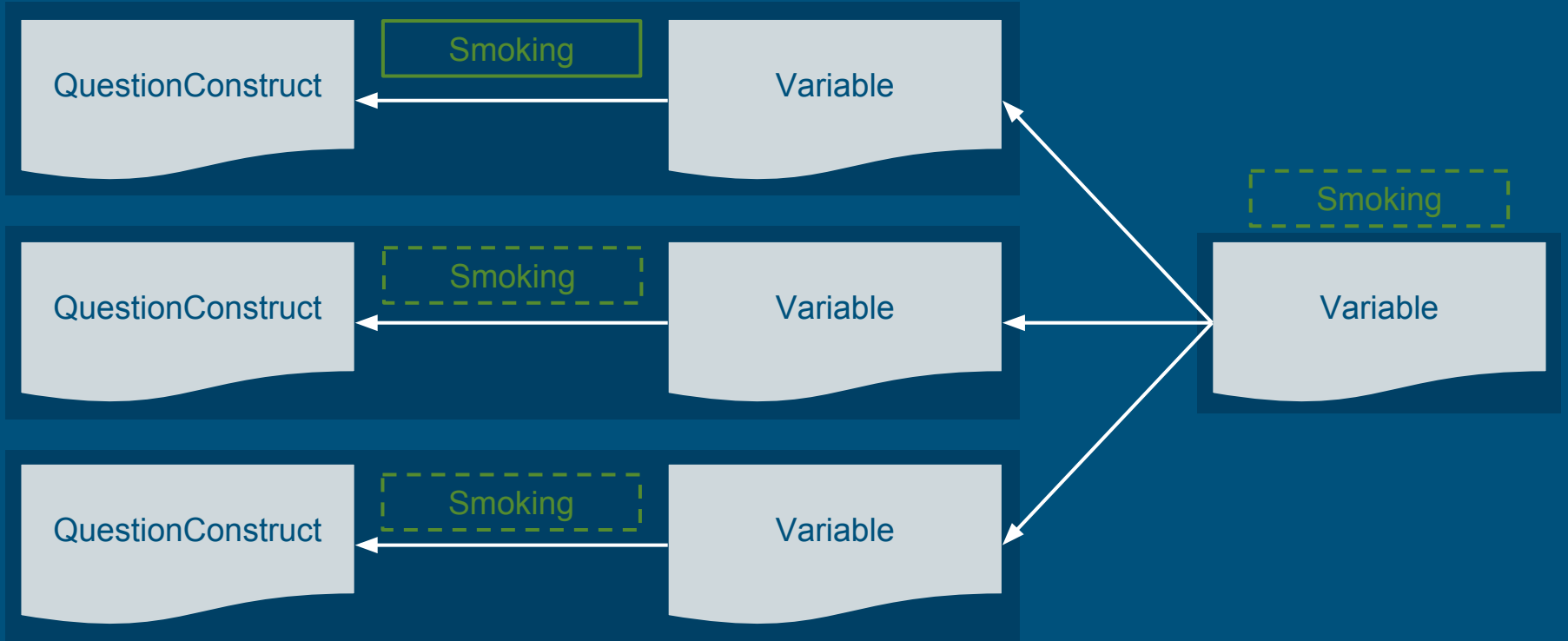
Clusters

- Collection of Strands
- Joined via **SourceVariableReference**
- Suggests the same topic

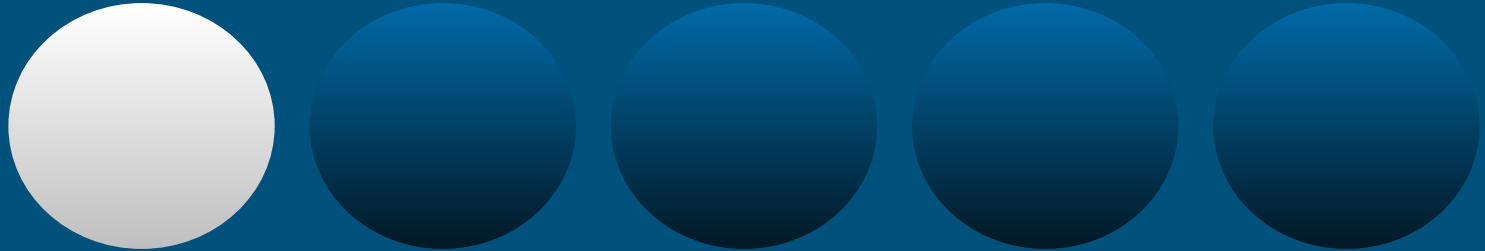
Cluster



Cluster



Savings



30,000

30,000



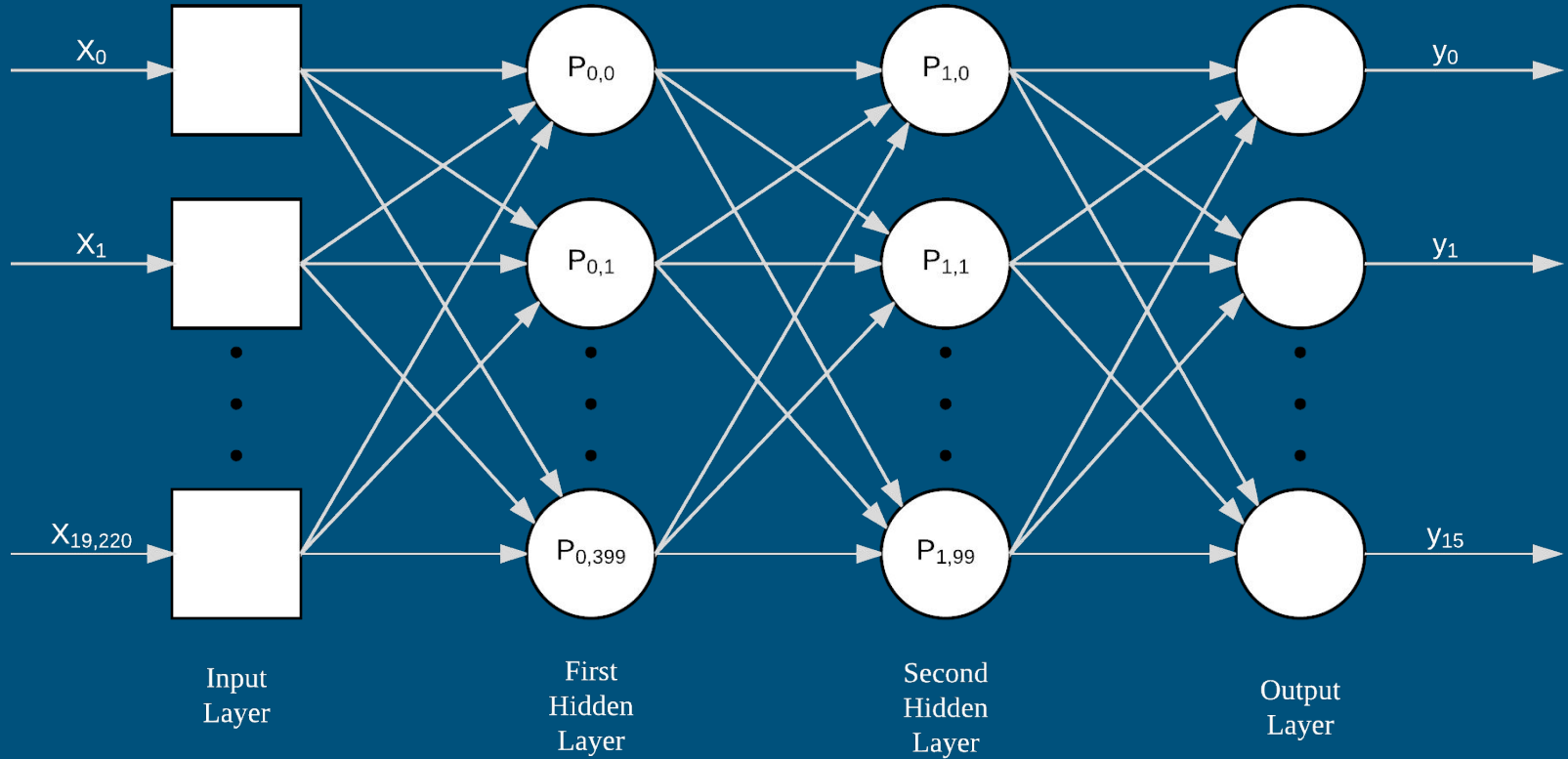
150,000


```
▼<ddi:DDIInstance xmlns:d="ddi:datacollection:3_2" xmlns:ddi="ddi:instance:3_2" xmlns:g="ddi:group:3_2"
xmlns:r="ddi:reusable:3_2" versionDate="2016-12-03T18:38:04+00:00">
  <r:URN>urn:ddi:uk.will:will_three-ddi-000001:1.0.0</r:URN>
  ▼<r:Citation>
    ▼<r:Title>
      <r:String xml:lang="en-GB">will_three instance 01</r:String>
    </r:Title>
    ▼<r:SubTitle>
      <r:String xml:lang="en-GB">Metadata documented by CLOSER using Archivist.</r:String>
    </r:SubTitle>
  </r:Citation>
  ▼<g:ResourcePackage versionDate="2016-12-03T18:38:04+00:00">
    <r:URN>urn:ddi:uk.will:will_three-rp-000001:1.0.0</r:URN>
    ▼<r:Citation>
      ▼<r:Title>
        <r:String xml:lang="en-GB">will_three resource package 01</r:String>
      </r:Title>
    </r:Citation>
    ▼<r:Purpose>
      <r:Content xml:lang="en-GB">not specified</r:Content>
    </r:Purpose>
    ▶<d:InterviewerInstructionScheme versionDate="2016-12-03T18:38:04+00:00">...</d:InterviewerInstructionScheme>
    ▶<d:ControlConstructScheme>...</d:ControlConstructScheme>
    ▶<d:QuestionScheme versionDate="2016-12-03T18:38:04+00:00">...</d:QuestionScheme>
    ▶<d:QuestionScheme versionDate="2016-12-03T18:38:04+00:00">...</d:QuestionScheme>
    ▶<l:CategoryScheme versionDate="2016-12-03T18:38:04+00:00">...</l:CategoryScheme>
    ▶<l:CodeListScheme versionDate="2016-12-03T18:38:04+00:00">...</l:CodeListScheme>
    ▶<d:InstrumentScheme>...</d:InstrumentScheme>
  </g:ResourcePackage>
</ddi:DDIInstance>
```

Machine learning

Model

- Multilayer Perceptron (MLP)
- Trained using variable labels
- Vocabulary of 19,221 words
 - Weighted using term-frequency inverse-document-frequency
- Predicting level-1 CLOSER Topics



Entire Dataset

Study	No. of Variables Used
ALSPAC	46,787
BCS70	6,813
HCS	1,026
MCS	898
NCDS	4,431
NSHD	1,445
SWS	2,511
USoc	0
Total	63,911

Training

- 50,000 labels
- Batches of 200
- Single core i7-5775R Intel & 5GB RAM
- 23 iterations
- Approximately 7 hours

Evaluation

- 13,911 labels
- 93% predicted 'correctly'
- Accuracy greatly depends on topic
- Unweighted random allocation accuracy 6.3%

Evaluation

- Most accurate: Omics
- Least accurate: Expectations, attitudes and beliefs

Future work

- Apply model to entire Strands
- Train for level-1 and level-2 topics
- Use study as a weighted input
- Integrate with current metadata enhancement tools aaS

Thank you

Any questions?

Special thanks you to

- ESRC
- MRC
- JetBrains

Contact me

- @willpoynter
- w.poynter@ucl.ac.uk

Contact CLOSER

- @CLOSER_UK
- closer@ucl.ac.uk