# Accuracy of natural language processing-based classifiers for automated identification of studies on humanistic and economic burden of disease

Jean-Baptiste Krohn

# heoro.com

An effortless and comprehensive approach to burden of illness reviews

About Us   Subscriptions   Blog   User Guide   Use Cases   Contact Us   Dashboard                    Sign in

# Burden of Illness Database

The heoro.com™ database saves you time and money by pre-screening thousands of abstracts and indexing them by disease, type of intervention, study methodology and geographical setting.

**Patient- and clinician-reported outcome studies**: identify all instruments and utility measurements used in a particular disease, and shortlist validation studies with a single click.

**Economic burden studies**: instantly find data on direct or indirect costs, resource use and treatment patterns.

**Economic evaluations**: rapidly filter cost-effectiveness or cost-utility analyses from other economic evaluations.

**Mortality trend studies**: efficiently identify studies reporting relative mortality and trends in survival.

*Legal terms and conditions*

➡ Sign in

| Email address | Password | Sign in | Go to dashboard |

# The problem

- >100,000 abstracts from PubMed search for quality of life, economic burden, economic evaluations and mortality since 2005

- >180,000 abstracts since 1960

- Need a quick, affordable and accurate way to index these

Crystallise

heoro.com

# The data

# The data

>4,700 diseases

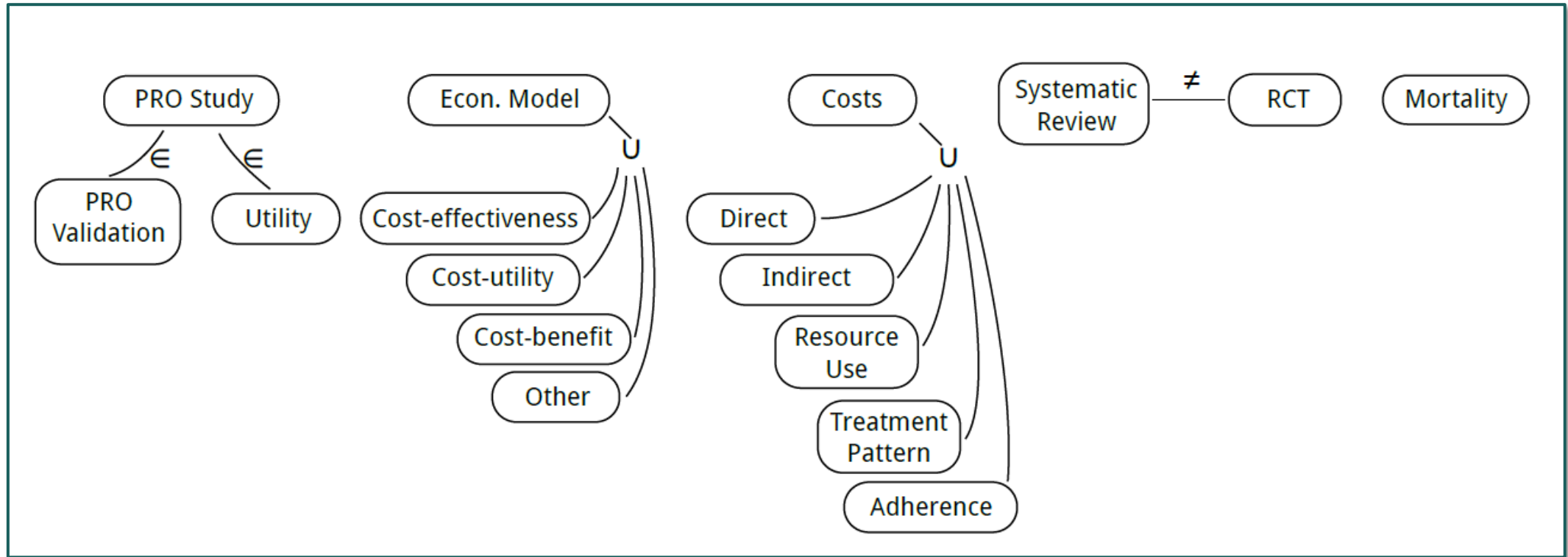4 study types,
11 sub-types



170 geographical
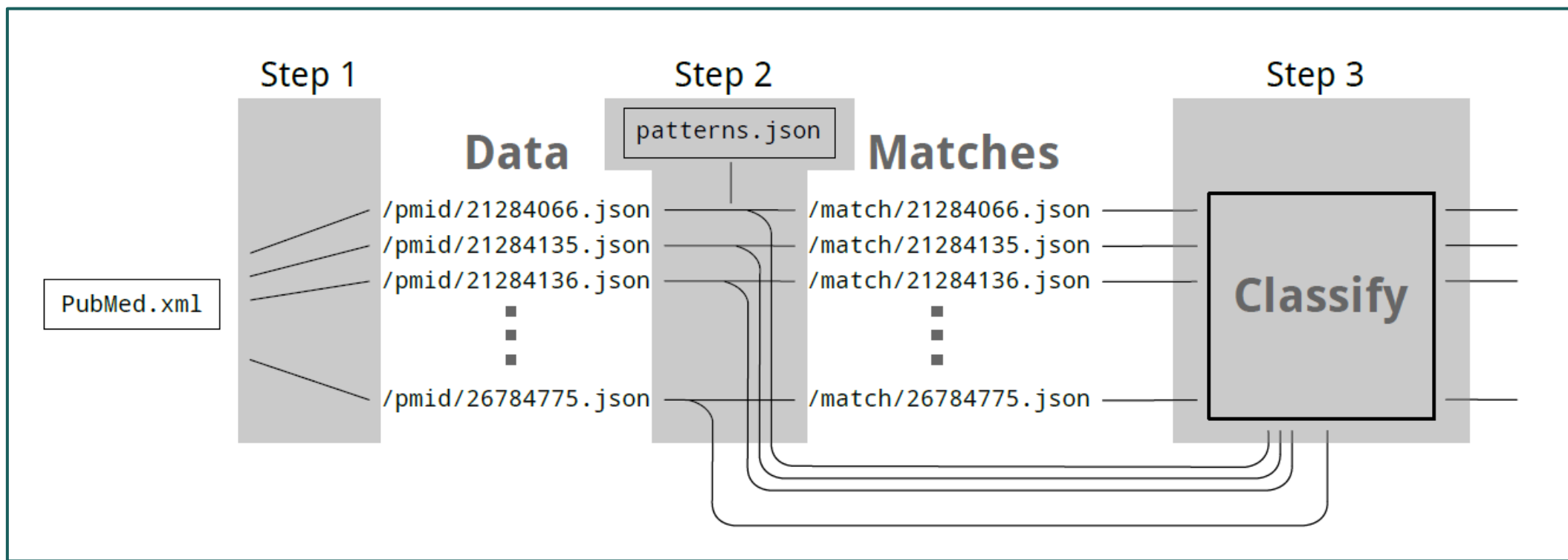locations

2 study methodologies

>8,400 interventions

>5,000 PRO instruments

# Type constraints

# The approach

# Step 1: PubMed data extraction

- Text: title and abstract.

- Structure: abstract paragraphs (e.g. OBJECTIVES, METHODS, CONCLUSIONS)

- Metadata: MeSH headings, journal, keywords, authors, etc

# Step 2: Pattern matching

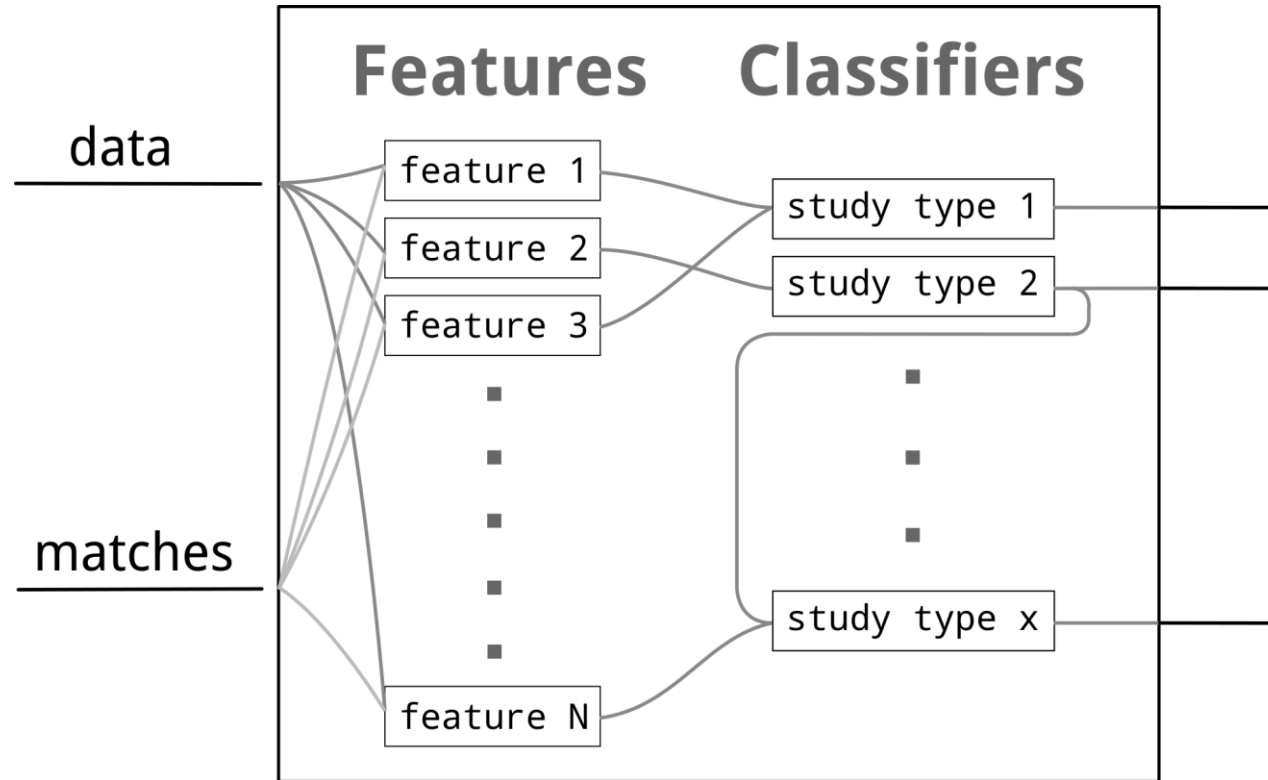- Example pattern (JSON):

```
"sf6d": [
  "/\\b[sS]hort(\\s+|-)[Ff]orm((\\s+|)[sS]urvey)?(\\s+|)6[Dd]\\b/",
  "/\\bSF-6D\\b/i"
]
```
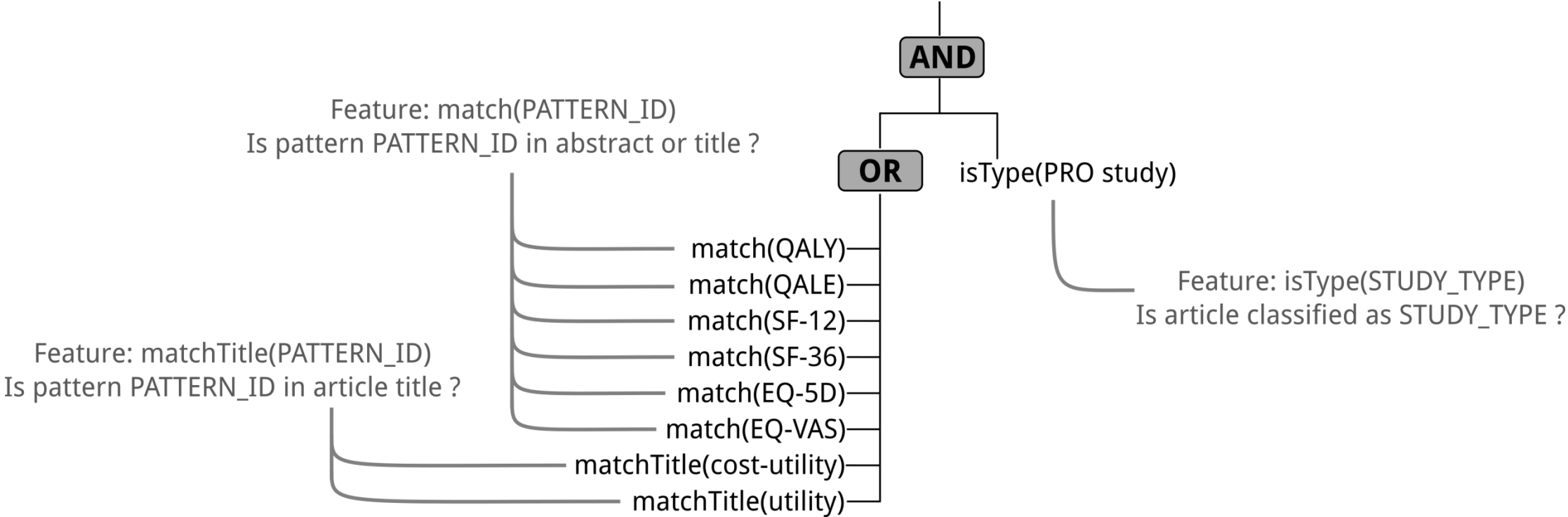
- Example matches (JSON):

```
{"name":"sf6d","path":["abstract",1,"text"],"details":[292,"Short Form-6D"]}
{"name":"sf6d","path":["abstract",5,"text"],"details":[ 80,        "SF-6D"]}
```

Crystallise

heoro.com

# Step 3: Classification

# Step 3: Example classifier

# Learning process

- Human expert assisted by Machine Learning system suggesting improvements based on learning set.

- Human expert intervenes on:
  - Pattern matching: define expressions to detect
  - Feature selection: determine features to use in classifier
  - Classifier structure: define boolean expression combining features

Crystallise

heoro.com

# Classifier accuracy: publications from 2005-2016

| Study type | Sensitivity | Specificity |
|---|---|---|
| PRO studies | 96% | 96% |
| Economic models | 95% | 97% |
| Costs and resource use | 92% | 95% |
| Mortality | 82% | 97% |
| RCT | 95% | 99% |
| Systematic review | 93% | 99% |
| Geographical location | 97% | |

# Classifier accuracy: publications from 2005-2016

| Study type | Sensitivity | Specificity |
|---|---|---|
| PRO studies | 96% | 96% |
| PRO validation | 99% | 97% |
| Utilities | 100% | 98% |

Crystallise

heoro.com

# Classifier accuracy: publications from 2005-2016

| Study type | Sensitivity | Specificity |
|---|---|---|
| Economic models | 95% | 97% |
| Cost-effectiveness | 87% | 100% |
| Cost-utility | 99% | 98% |
| Cost-benefit | 99% | 97% |
| Other | 93% | 93% |

Crystallise

heoro.com

# Classifier accuracy: publications from 2005-2016

| Study type | Sensitivity | Specificity |
|---|---|---|
| Costs and resource use | 92% | 95% |
| Direct costs | 79% | 98% |
| Indirect costs | 97% | 97% |
| Resource use | 90% | 97% |
| Treatment patterns | 99% | 93% |
| Adherence | 99% | 92% |

Crystallise

heoro.com

# Classifier accuracy: publications from 1960-2004

| Study type | Sensitivity | Specificity |
|---|---|---|
| PRO studies | 90% | 90% |
| PRO validation | 98% | 91% |
| Utilities | 99% | 99% |
| Cost-effectiveness models | 96% | 99% |
| Cost-utility models | 100% | 99% |
| Cost-benefit models | 100% | 100% |
| Other models | 86% | 94% |

# Classifier accuracy: publications from 1960-2004

| Study type | Sensitivity | Specificity |
|---|---|---|
| Direct costs | 80% | 93% |
| Indirect costs | 94% | 84% |
| Resource use | 76% | 91% |
| Treatment patterns | 98% | 82% |
| Adherence | 96% | 83% |
| RCTs | 97% | 98% |
| Systematic reviews | 95% | 98% |
| Mortality | 96% | 91% |

# Questions?

- **Contact:**
- Jean-Baptiste Krohn
- mail@jbk.io

- Alison Martin
- alison.martin@crystallise.com
- www.heoro.com

Crystallise

heoro.com