# Writing Schemas for Wikidata

Kat Thornton
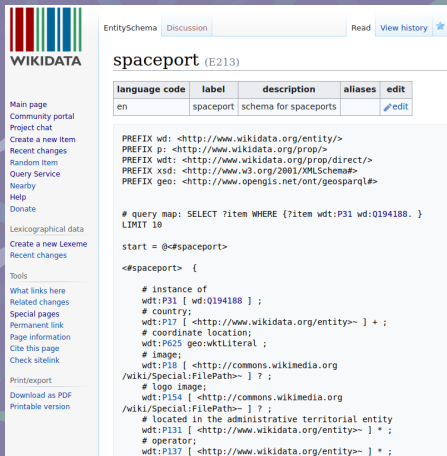
30 October, 2021

# Spaceport



**Figure 1:** Tanegashima Space Center, Photo Credit: NASA/Bill Ingalls, via Wikimedia Commons

# Finding a Schema



**Figure 2:** E213 Schema for Spaceports

# Schemas are data models

# Schema Namespace

- Wikidata's E Namespace is dedicated to schemas.
- 300+ schemas have been contributed since the namespace became available in May, 2019.
- Wikidata supports multiple data models per domain.
- In some cases this is necessary, in other cases it will be possible to build consensus around one shared model.

# ShEx

- ShEx is a formal modeling and validation language for RDF graphs
- Allows humans and machines communicate unambiguously about data assets
- Supports agile development of data models
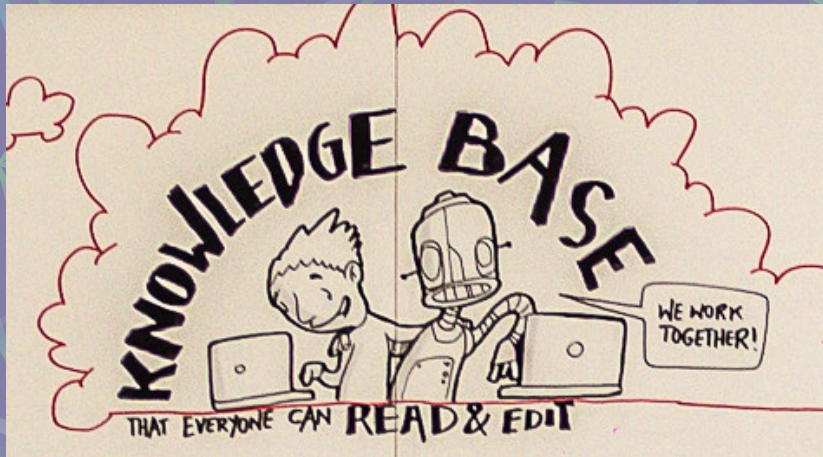- Learn more: http://shex.io

# Humans and Machines Together



**Figure 3:** Detail of mural by Magdalena Wiegner

# Members of National Academies of Science with academic degree PhD

```
SELECT ?item ?itemLabel ?id WHERE
{
?item wdt:P5380 ?id .
?item wdt:P512 wd:Q752297.
}
```

## Members of National Academies of Science, educated at statement, academic degree as a qualifier

```
SELECT ?item ?id
WHERE
{
?item wdt:P5380 ?id .
?item p:P69 [
ps:P69 ?inst ;
pq:P512 wd:Q752297 ].
}
```

# Building Consensus

- Frequency of usage of the 2 patterns?
- Propose a schema
- Discuss on talk page of schema

# A Shared Data Model for Academic Degrees is Desirable

- People who would like to reuse Wikidata data could become frustrated by needing to identify multiple related modeling patterns.
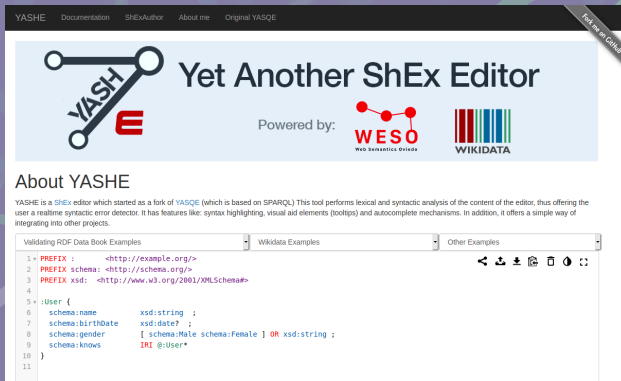
# Schemas: Internal or External

# Bringing an External Schema into Wikidata

- Expressing a model from a database, ontology, system external to Wikidata
- PRONOM technical registry of file formats
- file format with PRONOM id (E79) and
- File format with identification pattern (E237)

# Extracting an Internal Schema from Wikidata

- What if there is no schema yet?
- Then we can extract a schema
- Tools to try:
- sheXer
- Shape Designer

# Extract Schemas from Wikibase



**Figure 4:** Try YASHE out!

# ShEx Implementations

- shex.js (runs on n3.js)
- SHACLex (Scala)
- ShexJava
- PyShEx
- Ruby ShEx

# ShEx Test Suite

- 1103 validation tests
- 99 negative syntax tests
- 14 negative structure tests
- 431 schema conversion tests between ShExC, ShExJ and ShExR

# Online Validators for ShEx

- ShEx2 Simple Online Validator (JavaScript)
- RDF Shape (Scala)

# Interlinked ecosystem of schemas

# Schema Gallery



**Figure 5:** Chemistry schemas available in Wikidata

# Multilingual Schemas

# Importing One Schema into Another

# What's Next?

- More editors will contribute schemas
- More schemas will be interlinked via IMPORT
- This ecosystem of interlinking schemas will support the work of more editors
- Wikidatans will create more tooling that leverages schemas

# Find ShEx Papers

- Scholia is a Wikidata-based service that generates scholarly profiles
- Scholia for shex



Scholia  Author  Work ▾  Organization ▾  Location ▾  Event ▾  Project ▾  Award  Topic ▾  Tools ▾  Help ▾

topic

## ShEx (Q29377880)

Shape Expressions (ShEx) is a language for validating and describing RDF. It was proposed at the 2012 RDF Validation Workshop as a high-level, concise language for RDF validation. The shapes can be defined in a human-friendly compact syntax called ShExC or using any Resource Description Framework (RDF) serialization formats like JSON-LD or Turtle. ShEx expressions can be used both to describe RDF and to automatically check the conformance of RDF data. ... (from the English Wikipedia)

## Recently published works on the topic 🔲RSS

Show  10 ⌄  entries                                                          Search: [          ]

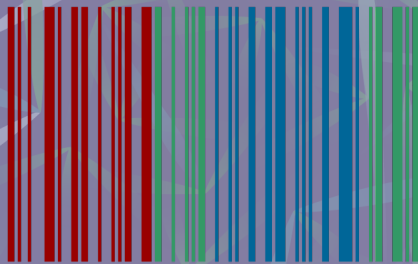| Date | Work | Topics |
|------|------|--------|
| 2019-05-25 | Using Shape Expressions (ShEx) to Share RDF Data Models and to Guide Curation with Rigorous Validation | ShEx |
| 2018-09-14 | XMLSchema2ShEx: Converting XML validation to RDF validation | Extensible Markup Language // ShEx |

# Quick Start

- Jose Labra's Wikipedia Weekly Entity Schemas and Shape Expressions session
- Primer: http://shex.io/shex-primer/
- ISWC 2020 Tutorial Shapes applications and tools
- Read the book Validating RDF Data

# Getting started with ShEx on Wikidata

- Wikidata Wikiproject ShEx
- Schema namespace
- Browse a gallery of schemas

# Wikidata's E Namespace: Where the Ecosystem of Schemas Thrives

Thank you!

WIKIDATA