# The Linkable Open Data Environment - LODE:
## Overview of Data, Tools, and Collaborations

Data Exploration and Integration Lab (DEIL)
Centre for Special Business Projects (CSBP)

May 2, 2019

Delivering insight through data, for a better Canada

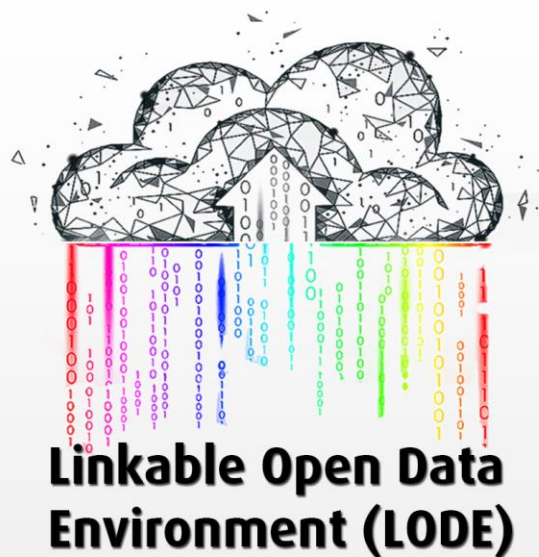# Open microdata: a vast, growing but still underutilized type of data

- **Microdata** – non sensitive and non-personal information on buildings, businesses, addresses, property values, infrastructure assets, and much more

- **From authoritative sources** – municipal, regional, provincial governments and, increasingly, also private sector stakeholders

- **Released with an open data license** – that encourages the use of the data

- **Rapidly expanding**

Delivering insight through data, for a better Canada

Canada

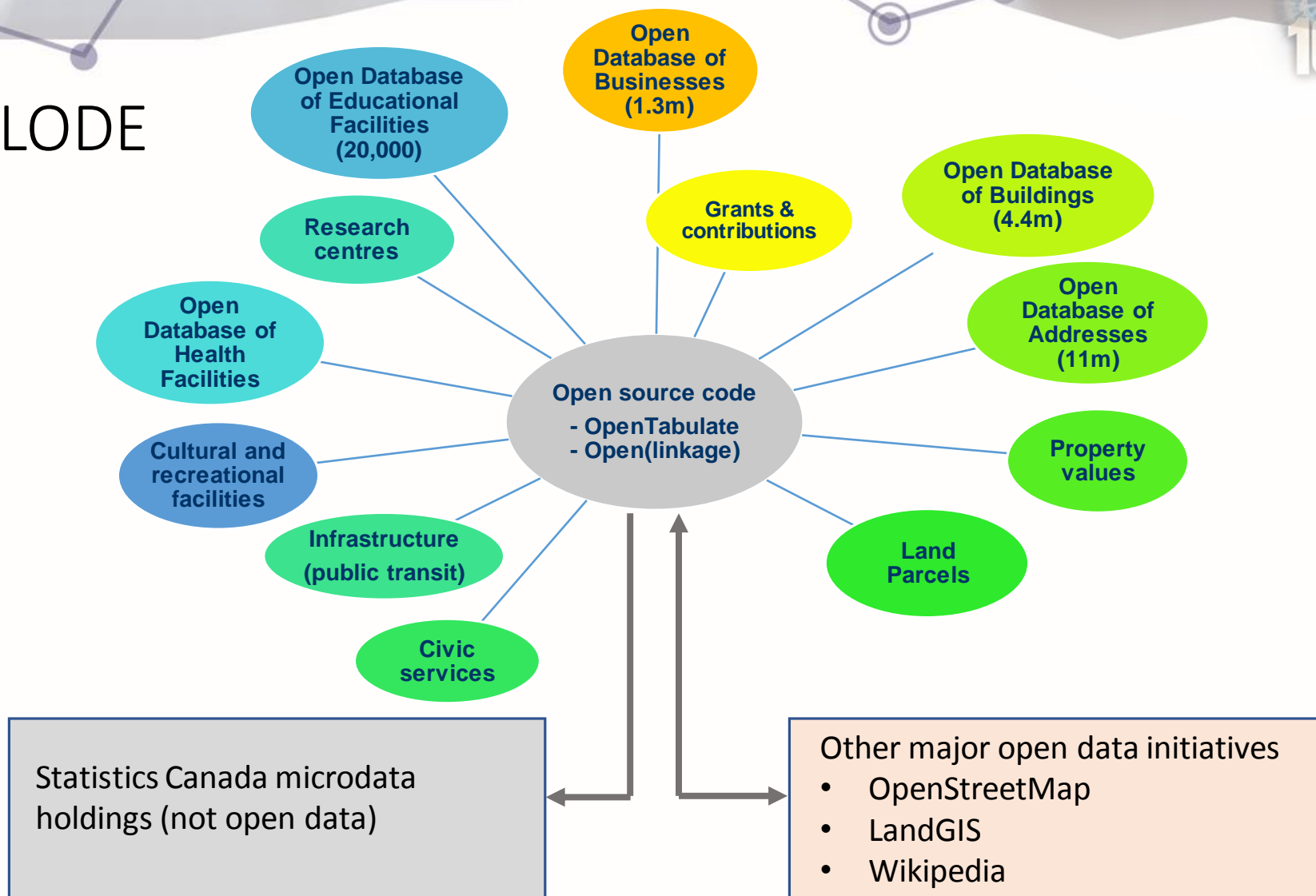# Linkable Open Data Environment (LODE)

- Open microdata from authoritative sources that have been brought into an environment that is suitable for data linkage

- Goal: harmonised and standardised datasets made available under the [Open Government License - Canada](#)

- Vast majority of datasets are from governmental sources (municipal, regional, provincial or federal)

- LODE is in development and you can use/contribute to it

2

## Draft logo



Linkable Open Data
Environment (LODE)

- Experimental and work in progress

- Much of LODE development is expected to be completed between April 2019 and March 2020

- Funded through external cost-recovery projects and internal projects

# LODE



**Open Database of Businesses (1.3m)**

**Open Database of Educational Facilities (20,000)**

**Grants & contributions**

**Open Database of Buildings (4.4m)**

**Research centres**

**Open Database of Health Facilities**

**Open Database of Addresses (11m)**

**Open source code**
- **OpenTabulate**
- **Open(linkage)**

**Cultural and recreational facilities**

**Property values**

**Infrastructure (public transit)**

**Land Parcels**

**Civic services**

Statistics Canada microdata holdings (not open data)

Other major open data initiatives
- OpenStreetMap
- LandGIS
- Wikipedia

13

# LODE Databases
## open for you to be used and shared

Statistics Canada / Statistique Canada

Canada

# Open Database of Buildings – version 2, March 1st, 2019

https://www.statcan.gc.ca/eng/open-building-data/index



- A compilation of 65 datasets originating from various government sources of open data (provincial, municipal)

- 4.4 million records of building footprints and variables calculated and standardized across all data providers

- Harmonised and standardised dataset made available under the Open Government License - Canada

3

# The ODB: example of the data

- Ex: Footprints for Richmond Hill, Toronto

- Quality is generally high, buildings are tightly knit



0    250    500 Metres

| OBJECTID* | Shape* | Longitude | Latitude | CSDUID | CSDNAME | Data_prov | Build_ID | Shape_Length | Shape_Area |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Polygon | -115.561757 | 51.18907 | 4815035 | Banff | Banff | 48150350000001 | 16.560241 | 16.963528 |
| 2 | Polygon | -115.569331 | 51.171372 | 4815035 | Banff | Banff | 48150350000002 | 87.531972 | 330.625531 |
| 3 | Polygon | -115.569616 | 51.178173 | 4815035 | Banff | Banff | 48150350000003 | 104.044015 | 573.938947 |

**Open Database of Educational Facilities (ODEF)**

- Completed with cost-recovery project funded by Crown-Indigenous Relations and Northern Affairs Canada (CIRNAC)

- 20,000 records (open data + publicly available data)

- Currently finalizing the review of licenses and approvals to release (public data) with GoC open license

- Harmonised and standardised dataset will be made available under the Open Government License – Canada

- Expected release: May/June 2019 – CSV file

9

**Open Database of Addresses (ODA)**

- Preliminary version is completed, work in collaboration with OpenAddresses

- About 11 million records from municipal and provincial open data sources

- Harmonised and standardised dataset will be made available under the Open Government License – Canada

- It will not include Postal Codes

- Expected release: June 2019

10

Statistics Canada   Statistique Canada

Canadä

- Preliminary version is available in LODE (GitHub)

- 26 data providers (federal, provincial, municipal)

- Approximately 1.4 million records (to be cleaned, de-duplicated, harmonised, etc.)

11

## Draft logo



Open Database of Health Facilities (ODHF)

- Work in progress

- Preliminary compilation of datasets available in LODE (GitHub)

# LODE resources and processing tools
## for open collaboration, use and development

# LODE on GitHub: repository of open microdata data sources

See: https://github.com/CSBP-CPSE/LODE-ECDO

# OpenTabulate 1.0 - a Python package, on Pypi, for LODE compilation and data cleaning

See: https://pypi.org/project/opentabulate/

# Open microdata is enabling collaboration

## Opportunity: collaborative data creation and analysis with many stakeholders **outside Statistics Canada**

- Open data as enabler and data value multipliers

- Reduces barriers to data sharing and costs (administrative, production, management). No need to have complex agreements and administrative burden

- Enriching the open data ecosystems may be one way we can unlock more data and more of their value

- Examples:
  - **Fleming College** collaborative project (ODB analytics and web mapping)
  - **UBC Master's of Data Science** collaborative project (ODB analytics)
  - **Canadian Read Cross** and **Digital Academy (CSPS)** on data analytics
  - **OSM communities across Canada** data imports into OSM and data cleaning and improvements

14

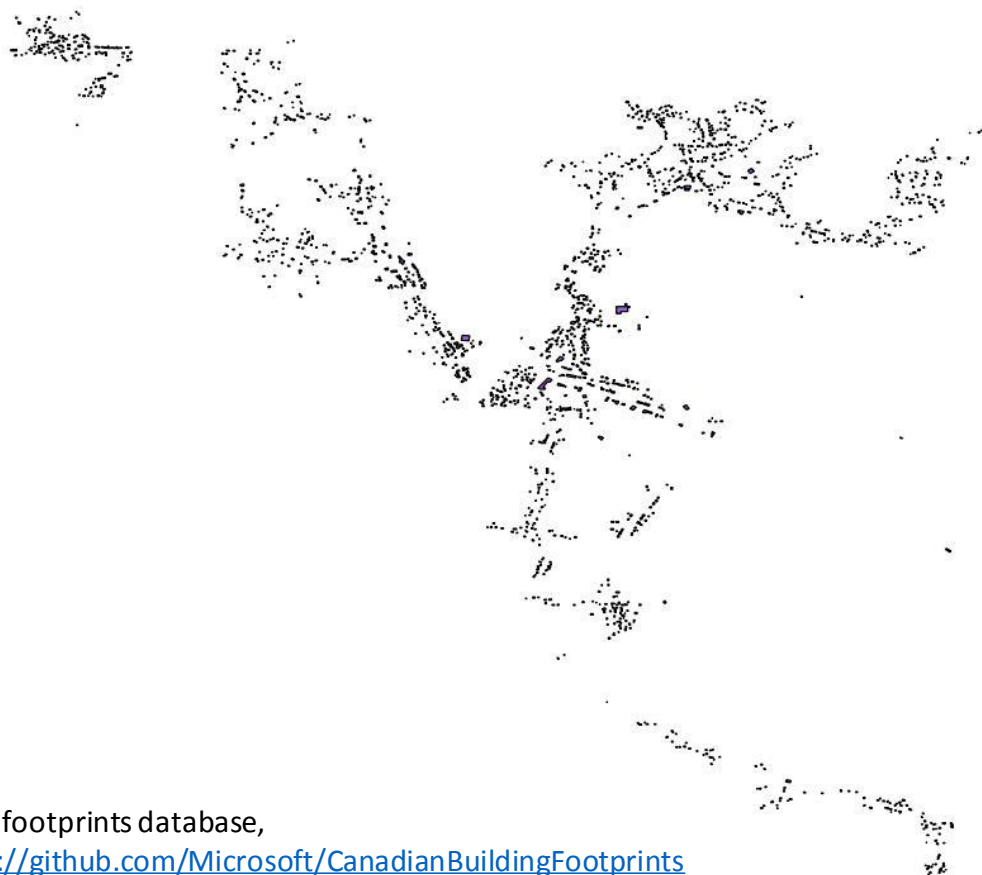# Enabling collaboration: Microsoft building footprints

- July 2018, StatCan and Microsoft started a collaboration to complete the mapping of building footprints across Canada

- Microsoft had released an open database of 125 million footprints for the U.S. based on satellite imagery extraction

- Microsoft used the Open Database of Building (version 1.0) to train a neural network model to extract building footprints from satellite imagery. The Microsoft database (about 12 million footprints) is available at:
  - Microsoft blog post (link)
  - Bing blogs (link)
  - https://github.com/Microsoft/CanadianBuildingFootprints

- Open data is a collaboration enabler and value multiplier

5

# Newfoundland and Labrador … as never seen before

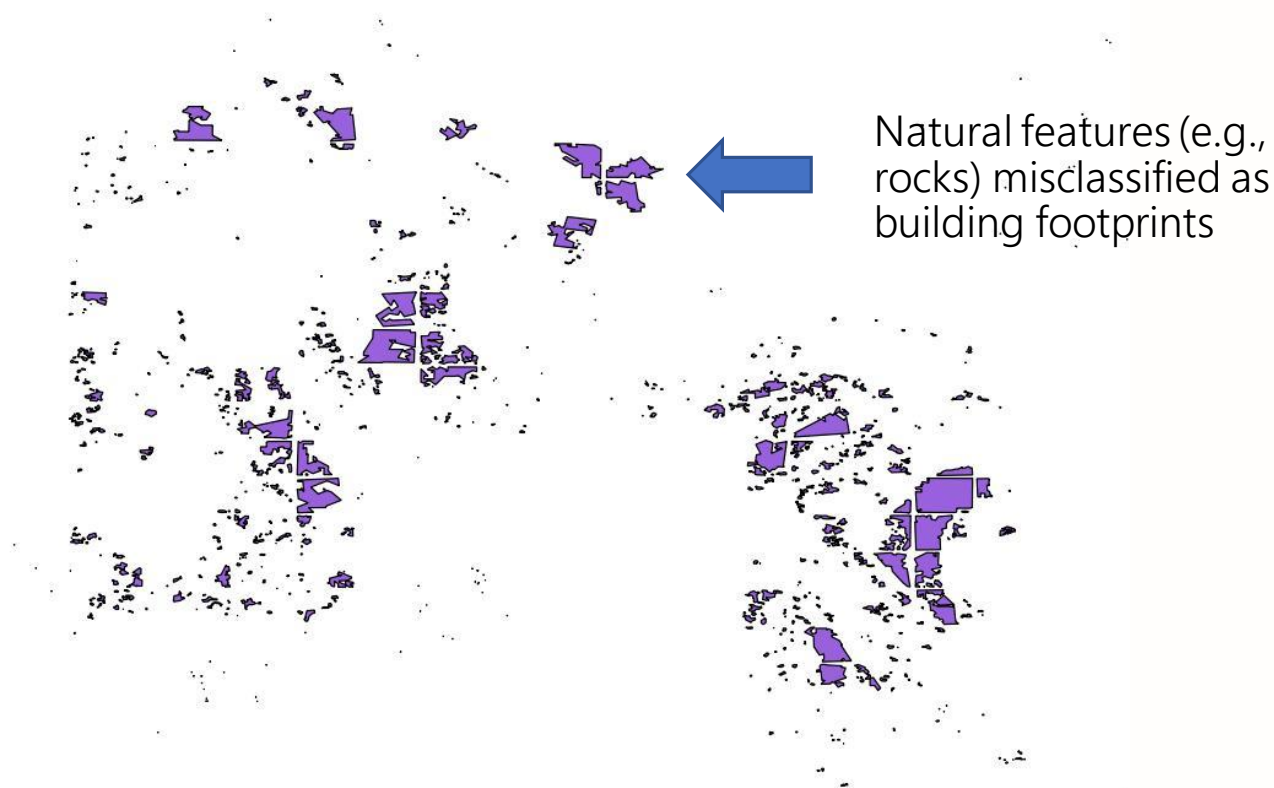Source: Microsoft building footprints database,
     available at: https://github.com/Microsoft/CanadianBuildingFootprints

# Twillingate region (NFL)

Source: Microsoft building footprints database,
    available at: https://github.com/Microsoft/CanadianBuildingFootprints

# ...but still a lot of false positives (mostly in remote regions)



Natural features (e.g., rocks) misclassified as building footprints

8

Source: Microsoft building footprints database, examples from Labrador region, available at: https://github.com/Microsoft/CanadianBuildingFootprints

# Conclusions

- You can use LODE (data and code)

- You can contribute to data development

- You can contribute to code development

- LODE is open and is there to be used!

15

# THANK YOU!

**For more information,**
alessandro.alasia@canada.ca
haaris.jafri@canada.ca

# MERCI!

**Pour de plus amples renseignements,**
alessandro.alasia@canada.ca
haaris.jafri@canada.ca

#StatCan100