# Open-source, transparent, and reproducible bioinformatics

PUBS2015
Franklin Bristow and Eric Marinier
National Microbiology Laboratory
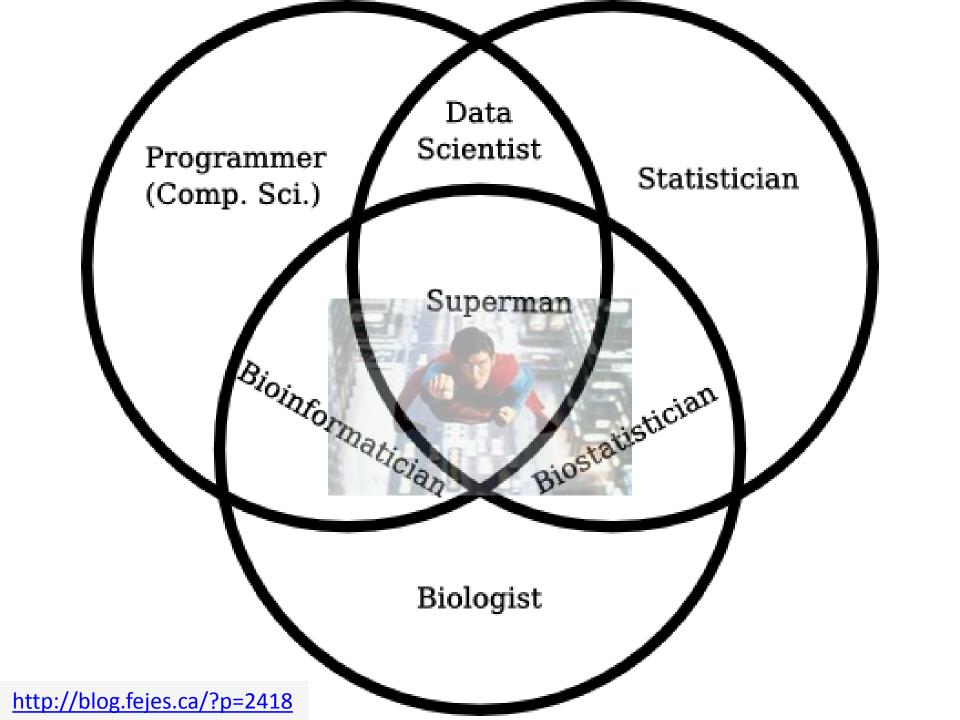Public Health Agency of Canada

# Overview

- [Who are we?](#)
- [Data](#)
- [Bioinformatics](#)
- [Open-source software](#)
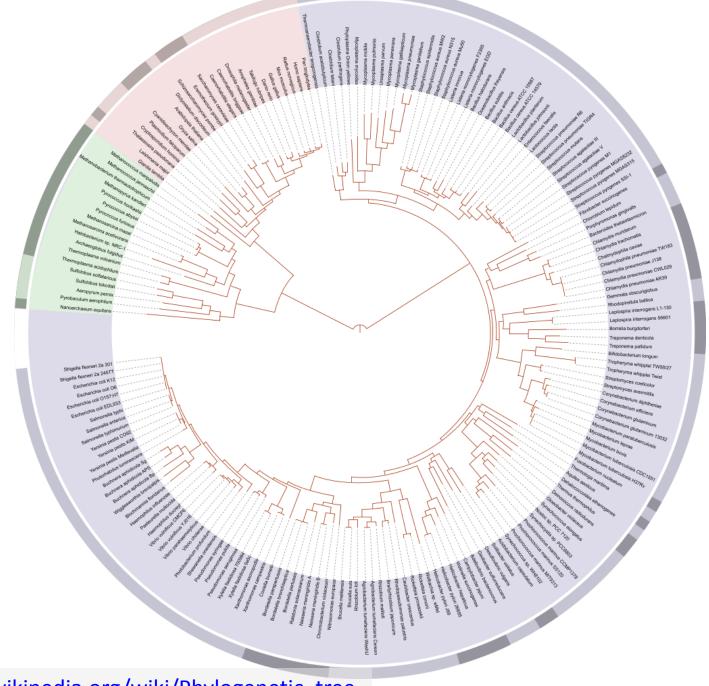- [Reproducibility](#)
- [Galaxy](#)

# Who Are We?

- Franklin Bristow

  ✉ franklin.bristow@phac-aspc.gc.ca

  🐙 https://github.com/fbristow

- Eric Marinier

  ✉ eric.marinier@phac-aspc.gc.ca

  🐙 https://github.com/emarinier
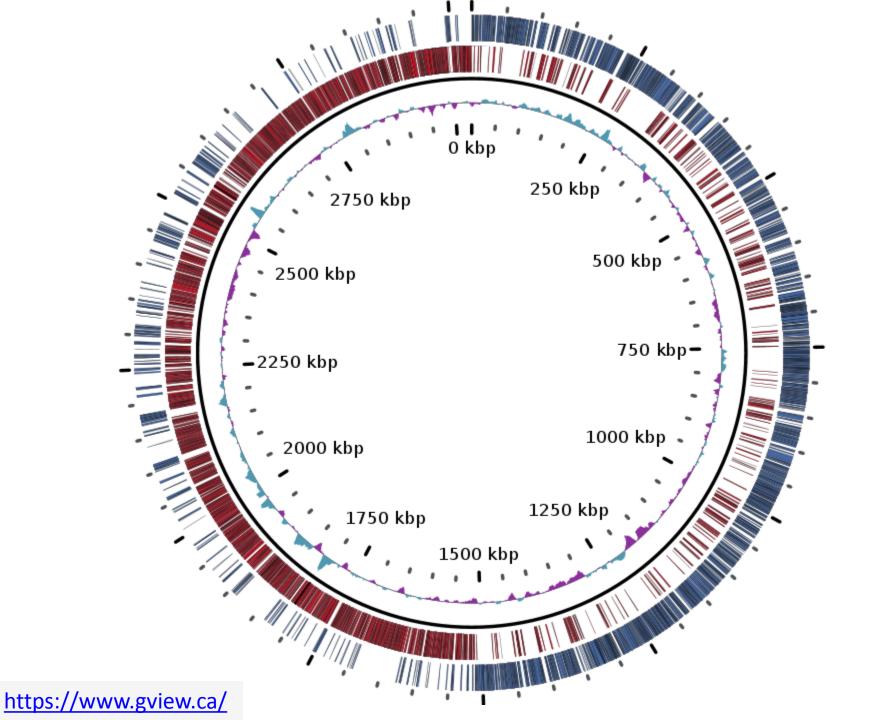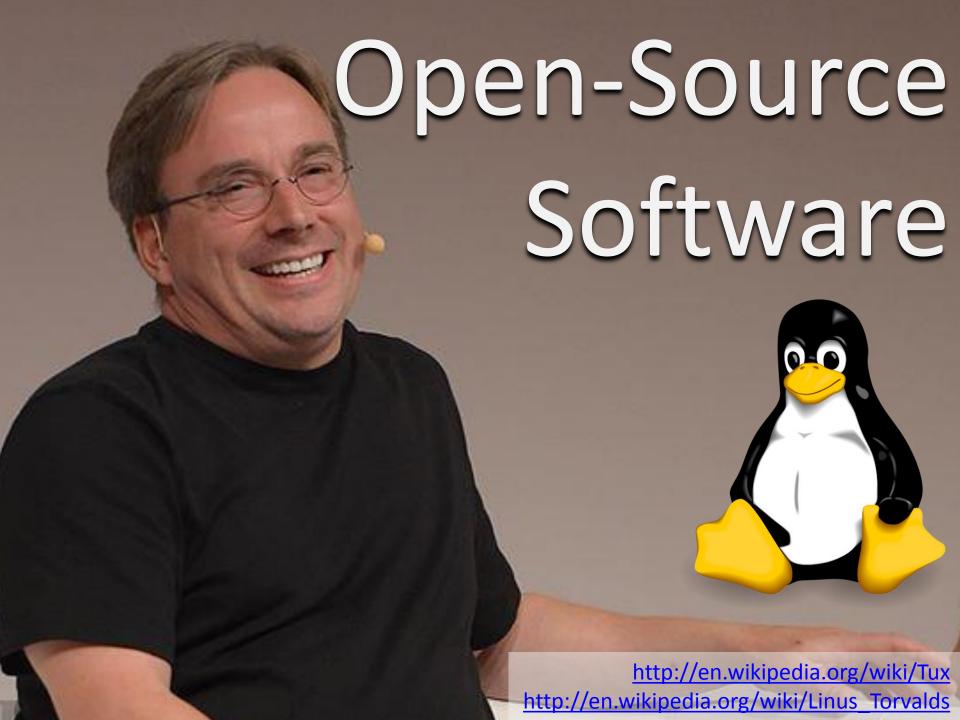
**Data**

# Data

Open-Source Software

Open-Source Software

# Command-line Software

# Reproducibility

Reproducibility

# Links

- Lecture notes and exercises:
  - https://github.com/emarinier/
- Public Galaxy instances:
  - https://usegalaxy.org
  - https://orione.crs4.it