# Spfy: speedy predictive genomics

Supervisor: Chad Laing
Student: Kevin Le

# Kevin Le

Currently:

- BASc Software Engineering (uOttawa)
    - Focus on applied math
- Co-op @NML Lethbridge (8-months)

Previously Completed:

- BSc Neuroscience (Dalhousie)
    - Transgenic mouse models of Alzheimer's
    - Cryptography & network security

Software Background:

- Largely Python
- Linux
- Virtualization

Career Goals:

- Big-data companies

# Background

**SuperPhy: predictive genomics for the bacterial pathogen Escherichia coli.**

Whiteside MD[1], Laing CR[2], Manji A[1], Kruczkiewicz P[1], Taboada EN[1], Gannon VP[1].

⊖ **Author information**

1    National Microbiology Laboratory @ Lethbridge, Public Health Agency of Canada, Lethbridge, AB, T1J 3Z4, Canada.
2    National Microbiology Laboratory @ Lethbridge, Public Health Agency of Canada, Lethbridge, AB, T1J 3Z4, Canada. chad.laing@canada.ca.

**Abstract**
**BACKGROUND:** Predictive genomics is the translation of raw genome sequence data into a phenotypic assessment of the organism. For bacterial pathogens, these phenotypes can range from environmental survivability, to the severity of human disease. Significant progress has been made in the development of generic tools for genomic analyses that are broadly applicable to all microorganisms; however, a fundamental missing component is the ability to analyze genomic data in the context of organism-specific phenotypic knowledge, which has been accumulated from decades of research and can provide a meaningful interpretation of genome sequence data.

# The problem

**Speed**

- Predict serotype, VF, AMR within a few minutes
- More responsive user-interface

**Simplicity**

- Easier to upload genomes and get results
- Intuitive navigation of the website

**Scalability**

- Quickly integrate new analyses & results
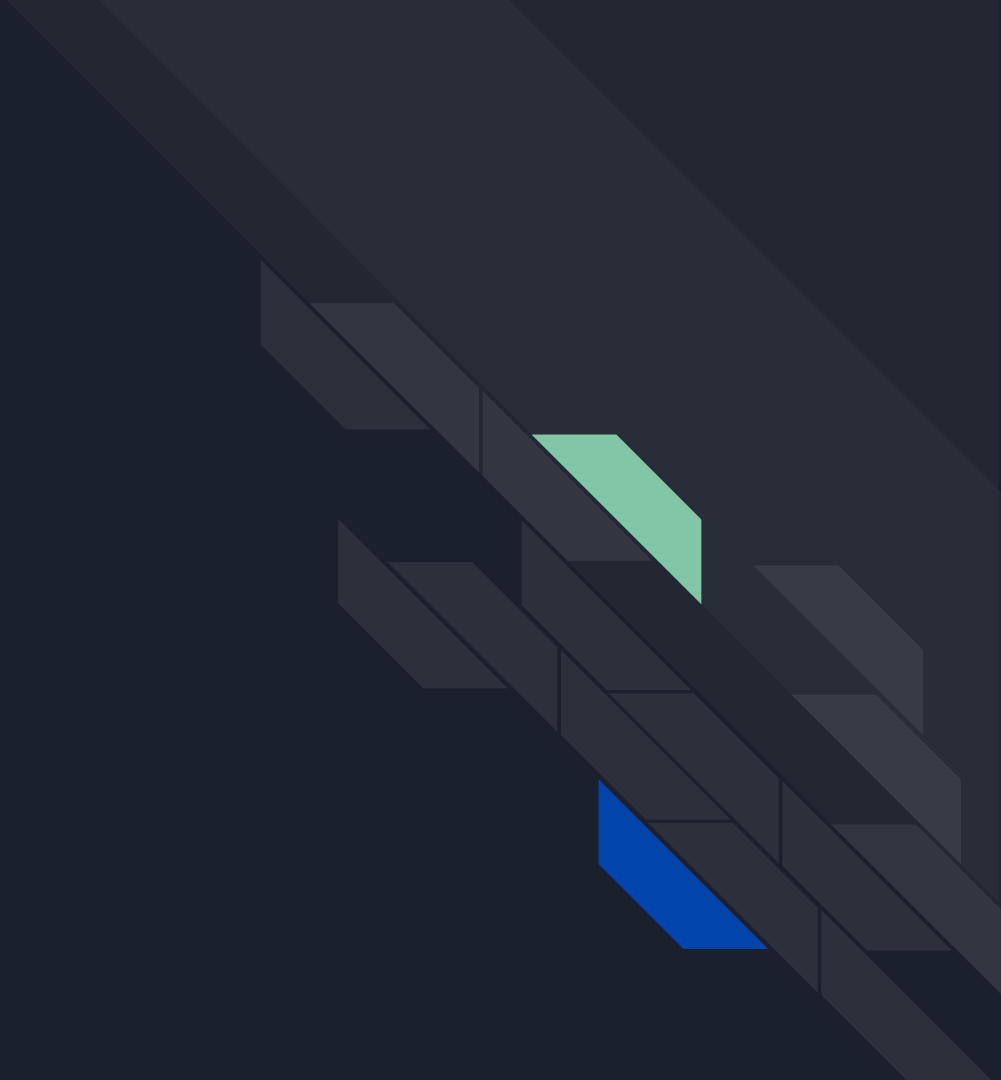- Ability to perform group comparisons across thousands of genomes

# Challenges deep-dive

Goal: Speed

- Predict serotype, VF, AMR within a few minutes
- More responsive user-interface

Solution: Modernize

- Docker
- Task queues
- ReactJS

Docker

# Docker

# Solution

Speed: task queues



Flask API
Scalable uWSGI

Redis Container
Redis DB
Redis Task Queue

Docker Volume
Datafiles
/datastore

RQ Worker Container
Priority Workers
4 Workers

RQ Worker Container
Reg. Workers
12 Workers

RQ Worker Container
UID Generation
1 Worker

Immediate Results

# Solution



Speed: ReactJS

# Solution



Speed: ReactJS

Webserver Container

Reverse Proxy
NGINX

Flask API
Scalable uWSGI

Scientist

Reactapp Container

Front-End
ReactJS

Docker

# Challenges deep-dive

**Goal: Simplicity**

- Easier to upload genomes and get results
- Intuitive navigation of the website

**Solution: Familiar Design**

- One click results
- Material design (Google)

# Solution

Simplicity: one-click results

# Solution

Simplicity: one-click results

# Solution

Simplicity: material design



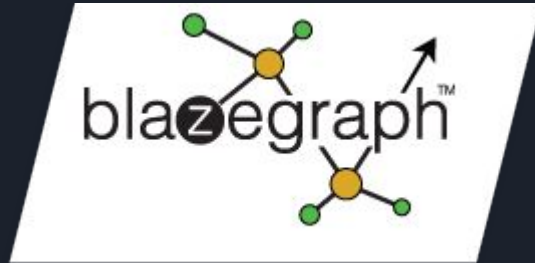Card Based

# Challenges deep-dive

Goal: Scalability

- Quickly integrate new analyses & results
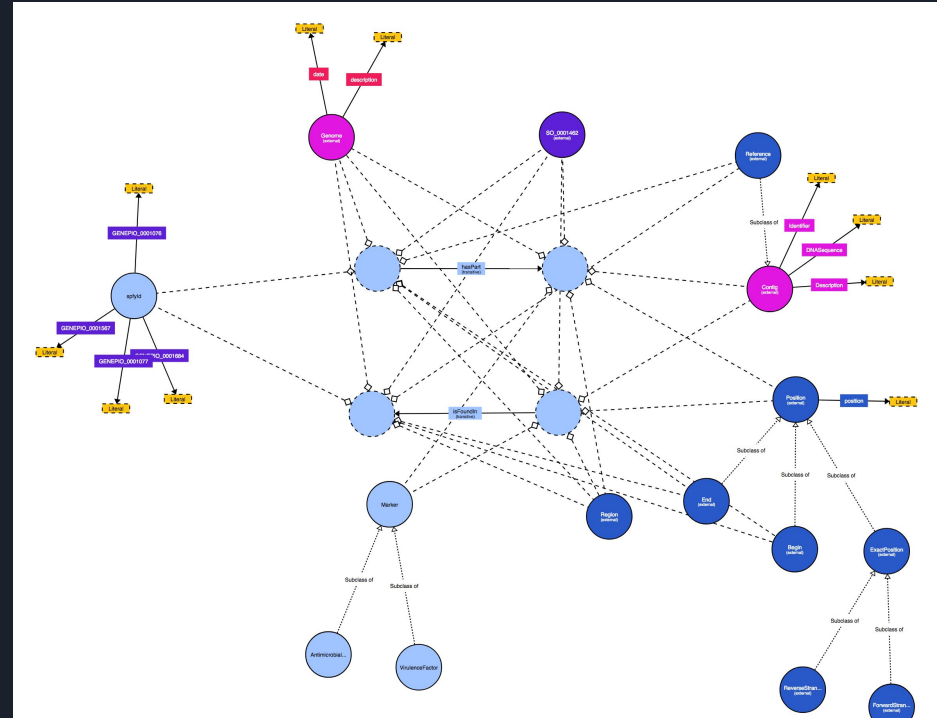- Ability to perform group comparisons across thousands of genomes

Solution: Technology

- Graph database
- Pandas, SciPy

# Solution



Scalability: graph database



Semantic Web

# Solution

Scalability: graph database

Blazegraph Container (LTS)

Graph Database
Blazegraph

RQ Worker Container

Priority Workers
4 Workers

RQ Worker Container

Reg. Workers
12 Workers

RQ Worker Container

UID Generation
1 Worker

Docker

# Solution



$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Scalability: pandas, scipy

Python Wrappers to C Code

# Challenges deep-dive

Speed

Simplicity

Scalability

**Modernize**

- Task queues
- ReactJS

**Familiar Design**

- One click results
- Material design (Google)

**Technology**

- Graph database
- Pandas, SciPy

# Timeline

- Wrapper for serotype, VF, AMR
- Upload to graph database

- Docker
- QC / CI
- Initial Deployment

| January | February | March | April |
|---------|----------|-------|-------|

- Task queues
- API (Flask)
- Basic website (AngularJS)

- Graph traversals
- Group comparisons: backend code (Flask)

# End Result

Platform Architecture

Live Demo

# End Result:

- Real-time serotype, VF, AMR prediction
- Within 2-3 minutes

- Storage and retrieval of genomes & results in a Graph Database
- Test set: 5353 GenBank genomes

- Live group comparisons of database entries within seconds
- Ex. O157 vs O53 for all known VFs