# Spfy: bringing real-time, big data batch analyses of E. coli to SuperPhy

using containerization and graph-based data storage

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)
  - Research: transgenic mouse models of Alzheimer's, molecular neurosci.
  - Cryptography & network security

Software Background:

- Largely Python
- Linux
- Virtualization

Career Goals:

- Big-data companies / grad. studies in comp. sci

# Bioinformatics Co-ops @Lethbridge Winter '17

# Why work on *E. coli*?

- ❖ Within Canada it is estimated that over 4,000 hospitalizations and 100 deaths  per year occur due to bacterial infections
- ❖ E. coli represents 33% of bacterial outbreaks from produce
- ❖ Incidence of ~30.3/100,000 people

# Why work on *E. coli*?

Common Symptoms:

- Stomach cramps
- Diarrhea
- Vomiting

Complications:

- Hemolytic uremic syndrome (HUS)
- Can lead to kidney failure

# Why work on *E. coli*?

❖ Traditional monitoring efforts used wet-lab methods for subtyping

# WGS-Based Predictions

## 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 Previous Version

**What Worked:**

- Pre-computed, predictive genomic analyses for Shiga toxin subtype, AMR genes, virulence factors
- Presence / absence of genomic regions and single nucleotide variants, for bacterial sub-groups

**What Didn't:**

- No batch upload and get results in real-time
- Large scale storage and retrieval of results in real-time
- Easy deployment and replication

Goal: BIG data

# Why big data?

Spfy replicates *E.coli* related functions of traditional reference labs

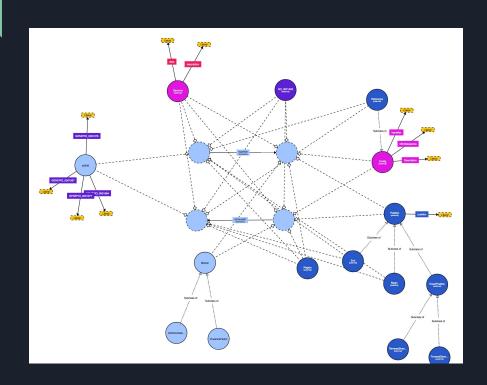- ❖ Performs the same tasks except using whole-genome sequencing (WGS) data

Spfy is also *online* & *in-silico*, incorporating species-specific knowledge for:

- ❖ Population-wide analysis
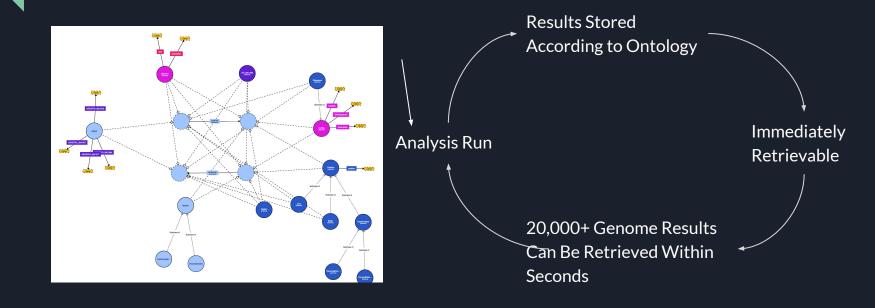- ❖ Historical context
- ❖ Automatically linking new results

# Main Approach

❖ Graph Database          ❖ Parallelization

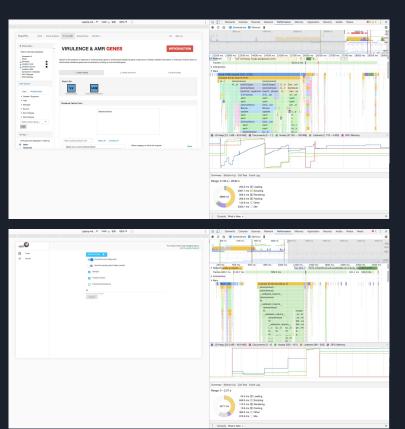# The Graph Database is the Core

# The Graph Database is the Core



Results Stored
According to Ontology

Analysis Run

Immediately
Retrievable

20,000+ Genome Results
Can Be Retrieved Within
Seconds

# The Web Platform

# How do we get to results?

# Major Change:
## 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
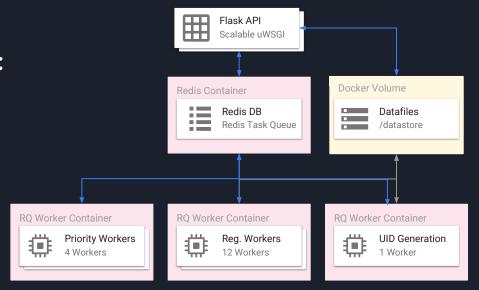
Reg. Workers
12 Workers

RQ Worker Container

UID Generation
1 Worker

Dozen+ Parallel Task Workers

# What can we do with this?

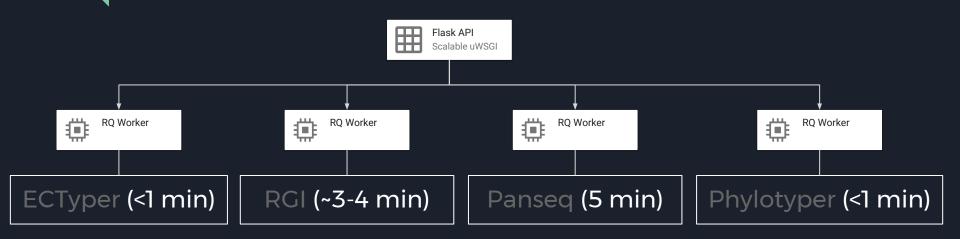given a new whole genome sequencing result...

Serotype, Virulence Factors via ECTyper (<1 min)

Antimicrobial Resistance Genes via RGI (~3-4 mins)

Pangenome Regions via Panseq (5 mins)

Stx Type & Closely Related Strains via Phylotyper (<1 min)

1 + 4 + 5 + 1 = ~11 min

# Recall: Parallel Task Queues

Flask API
Scalable uWSGI

RQ Worker

RQ Worker

RQ Worker

RQ Worker

ECTyper (<1 min)

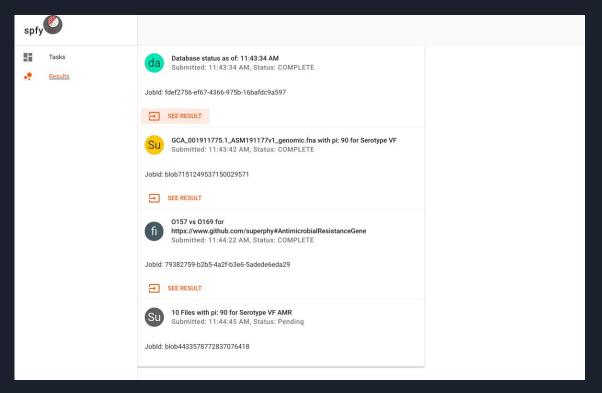RGI (~3-4 min)

Panseq (5 min)

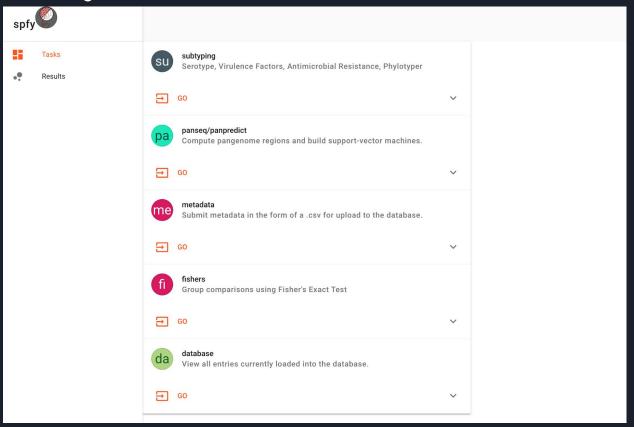Phylotyper (<1 min)

Not affected by the size of the database

# Group Comparisons

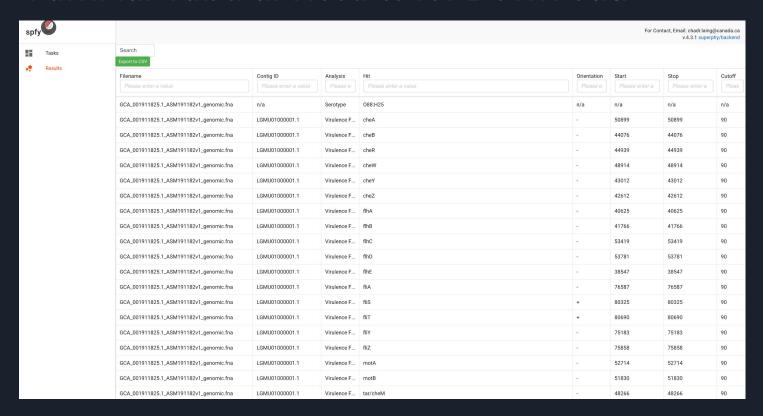Ask: what Pangenome Regions are over-represented among O157 strains over O53 strains?
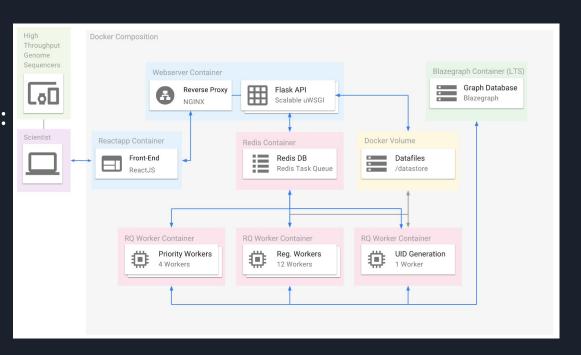
# Simplicity:
## Google's Material Design

Card Based

# Analysis Modules

# Familiar Tabular Results for Download

# Orchestration:
## Docker,
## Docker-Compose



Platform Architecture

Used Spfy to compute & store results for ~20,000 genomes

❖ One-Time Cost: Results are <u>permanently</u> stored for future comparisons - "Big Data"

# What was done

What didn't work in Superphy:

- No batch upload and get results in real-time
- Large scale storage and retrieval of results in real-time
- Easy deployment and replication

How it was solved in Spfy:

- Modern website and use of task queues
- Graph database
- Docker

# Where do we go from here?

❖ More analysis modules

❖ Data Visualization

❖ Machine Learning

# Thanks!

E.coli team @ Lethbridge

- ❖ Chad Laing
- ❖ Matt Whiteside
- ❖ Vic Gannon

Campy team @ Lethbridge

- ❖ Eduardo Taboada
- ❖ Dillon Barker

Questions?