# ProkEvo: an automated, reproducible, and scalable framework for high-throughput bacterial population genomics analyses
—

**Natasha Pavlovikj**

PhD student
Department of Computer Science and Engineering
University of Nebraska-Lincoln

February 25th , 2021

# Bacterial Population Genomics

- Group of individuals of the same bacterial species

- Populations evolve

- Differences between individuals are tiny

- Conclusions about populations

# Bacterial Population Genomics
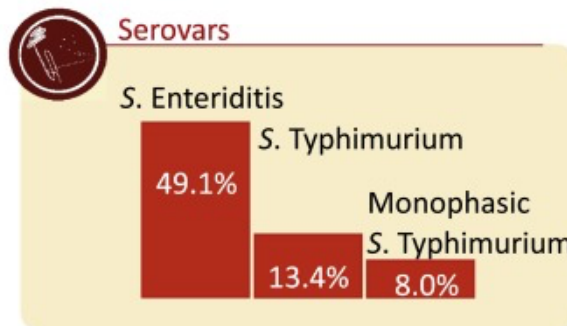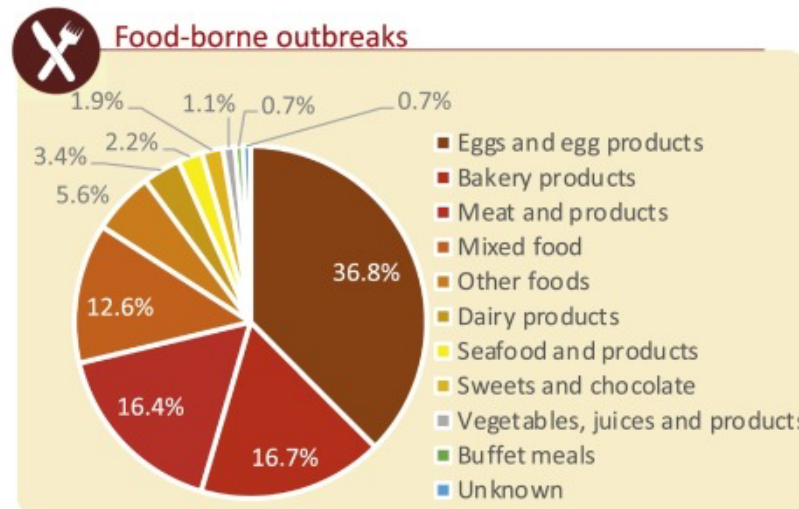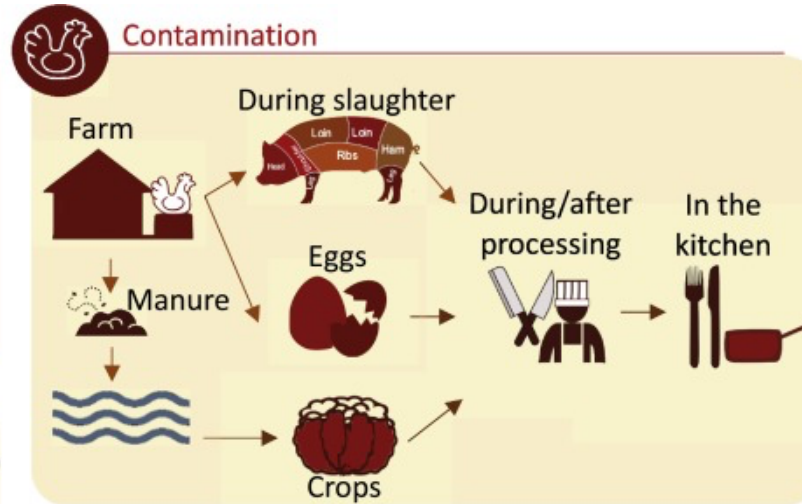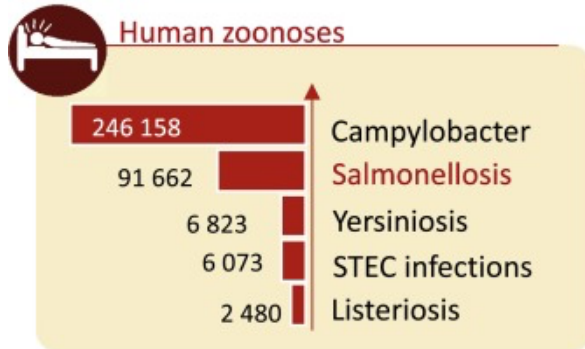
- Group of individuals of the same bacterial species

- Populations evolve

- Differences between individuals are tiny

- Conclusions about populations

We need a comprehensive pipeline that allows us to do large-scale population genomics, because only then we can understand something about the population structure, distribution of genotypes, evolutionary history, accessory genome distribution that COUPLED with ecological and epidemiological information can bring better understanding of how populations diversify and evolve over time.

One application of ProkEvo is in the field of Public Health and Food Safety: *Salmonella enterica*

# *Salmonella enterica*

**Goals:**

1. Determine the most Important genotypes
2. Infer the genomic events important for survival across the food chain in order to mitigate the risk of infection



### Human zoonoses
- 246 158 — Campylobacter
- 91 662 — Salmonellosis
- 6 823 — Yersiniosis
- 6 073 — STEC infections
- 2 480 — Listeriosis

### Salmonellosis in numbers
- Notification rate: 19.7/100 000
- Hospitalization: 42.5%
- Case-fatality rate: 0.25%
- € 3 billion

### Serovars
- *S.* Enteriditis — 49.1%
- *S.* Typhimurium — 13.4%
- Monophasic *S.* Typhimurium — 8.0%

### Contamination
Farm → During slaughter → During/after processing → In the kitchen
Manure, Eggs, Crops

### Food-borne outbreaks
- 36.8% Eggs and egg products
- 16.7% Bakery products
- 16.4% Meat and products
- 12.6% Mixed food
- 5.6% Other foods
- 3.4% Dairy products
- 2.2% Seafood and products
- 1.9% Sweets and chocolate
- 1.1% Vegetables, juices and products
- 0.7% Buffet meals
- 0.7% Unknown

**Example:**

**sugE** – resistance to quaternary ammonium
If that is a driver for a serovar to dominate, then it indicates a change in hygiene practice is needed, perhaps by switching disinfectants

Salmonella is one example only

# Applications

- Outbreak detection
- Source tracking
- Understanding epidemics
- Public Health Surveillance
- Pathogen transmission
- Discover ecological properties
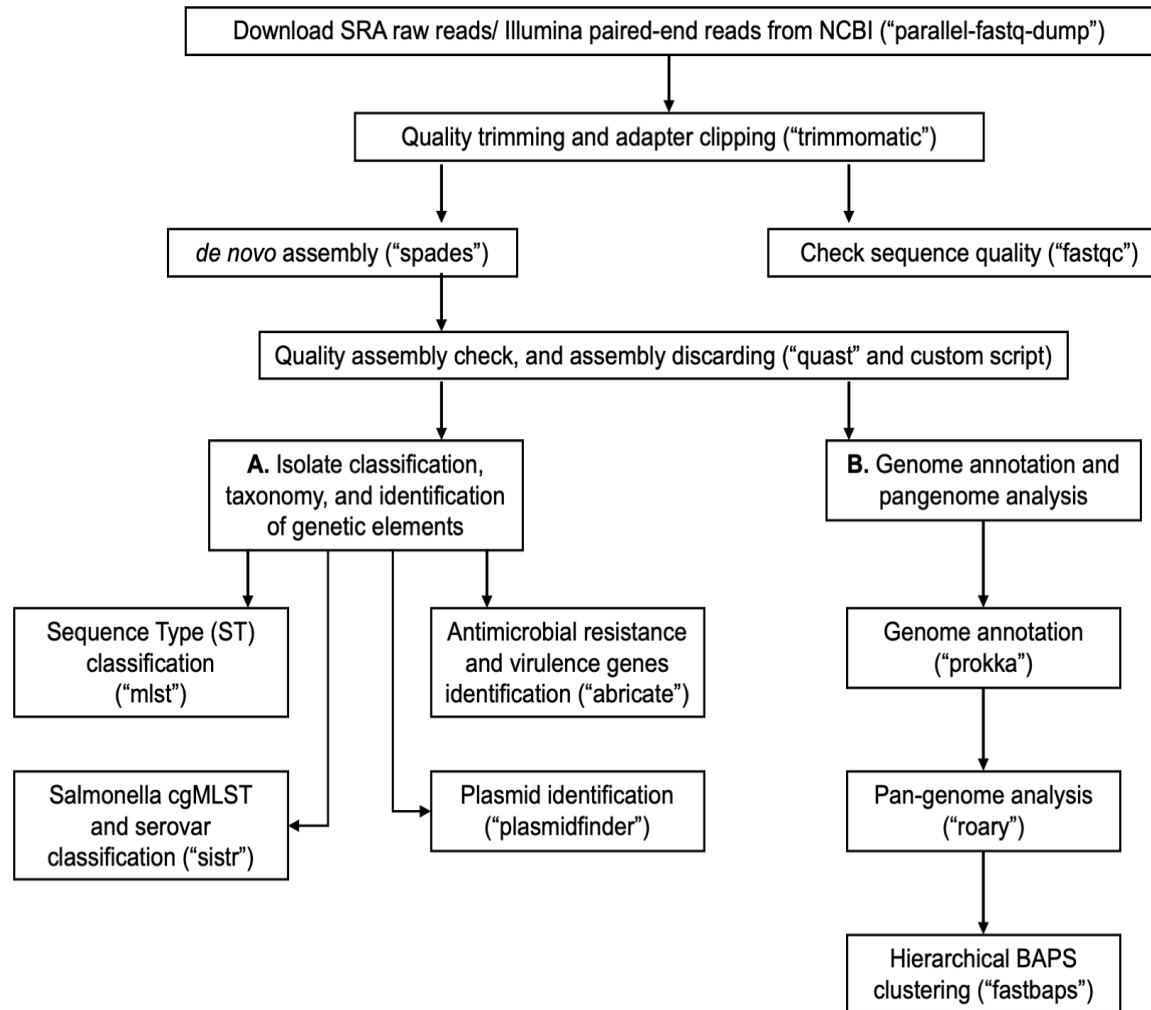
# Challenges

- Scaling and automating WGS analyses can be a challenging task that comes with its own costs and benefits

- Several existing automated pipelines for analyses of bacterial genomes

- No usage of workflow management system
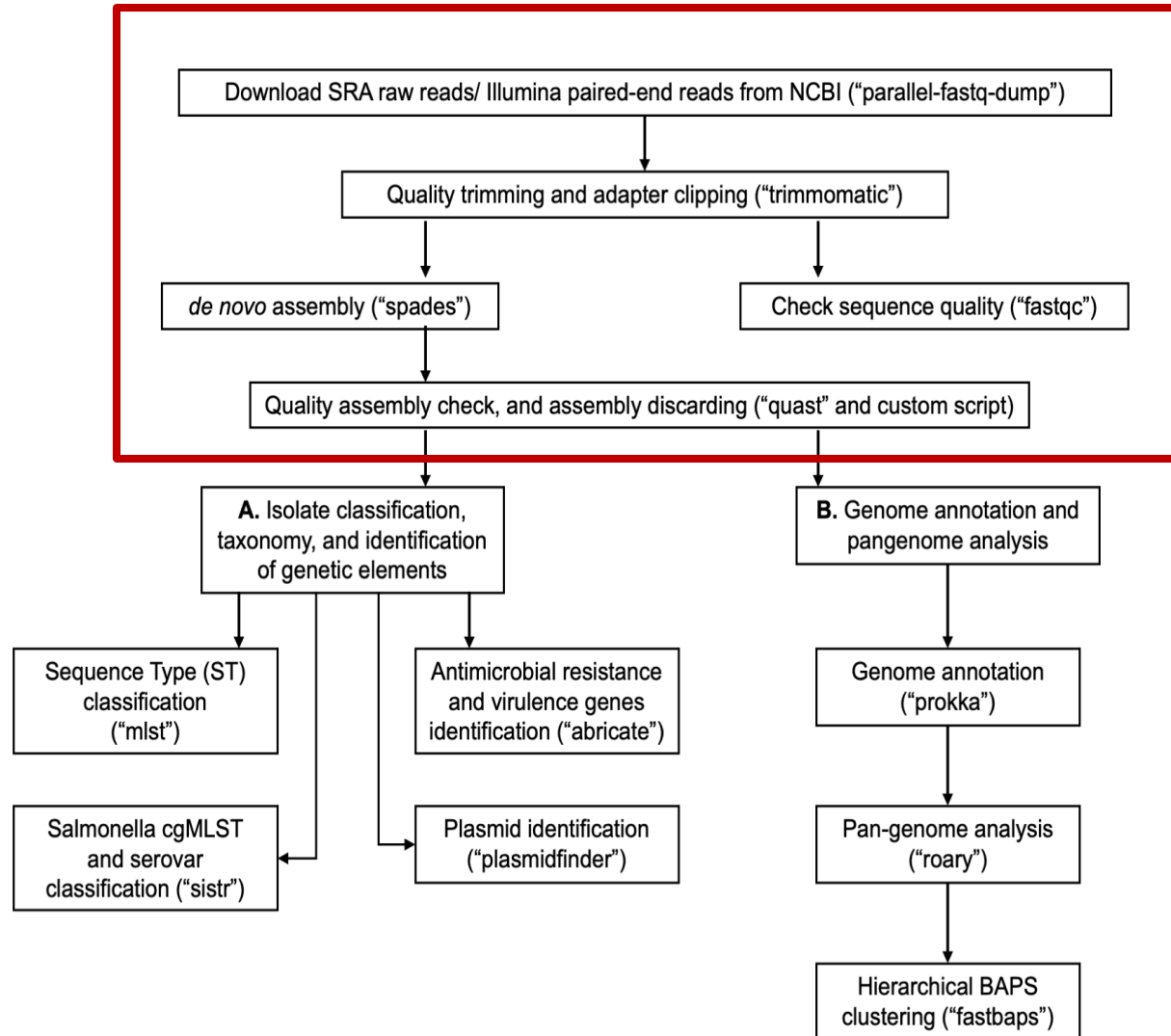
- Small number of used genomes

# ProkEvo

- Automated, open-source pipeline for population genomics analyses

- Uses Pegasus Workflow Management System

- Analyses of few genomes, as well as tens of thousands using high-performance and high-throughput computational resources

- Easily modifiable and expandable pipeline to include additional steps, custom scripts, user databases, and species

- Modular pipeline that can run thousands of analyses concurrently if the resources are available

- Distributed with conda environment and Docker image for all bioinformatics tools and databases needed to perform population genomics analyses

- The output of ProkEvo is used to provide some guidance on how to perform population-based analyses with reproducible Jupyter Notebooks
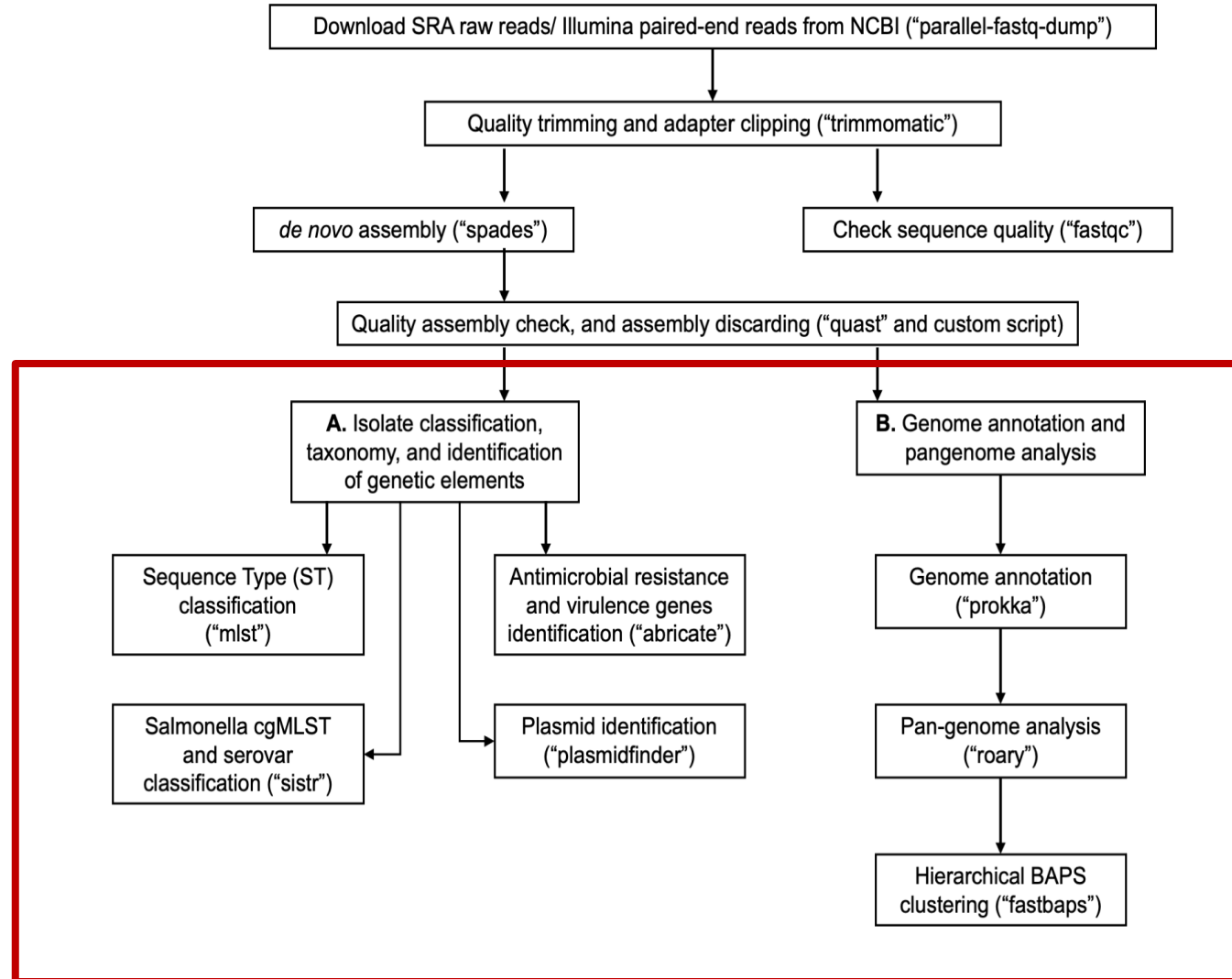
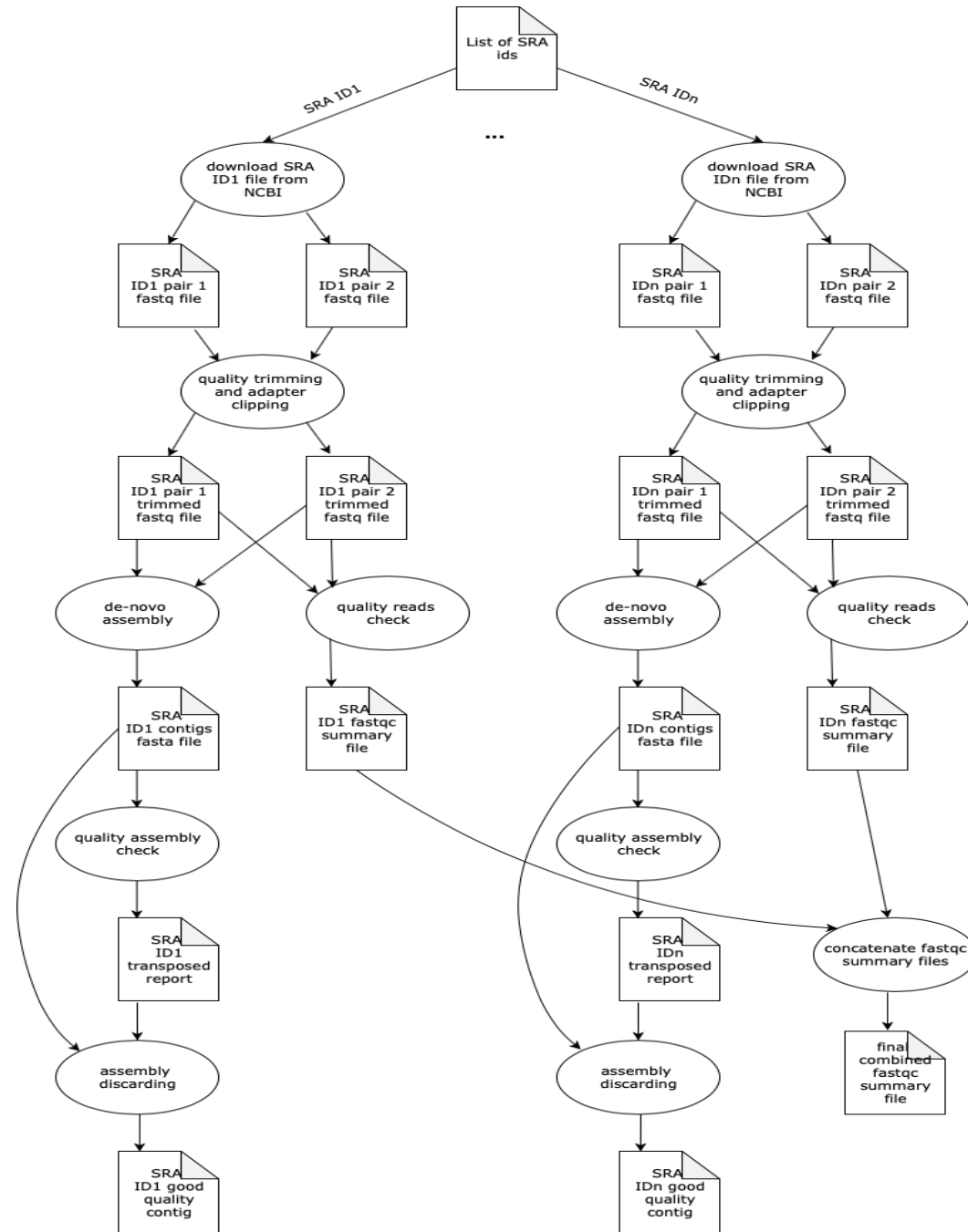- https://github.com/npavlovikj/ProkEvo
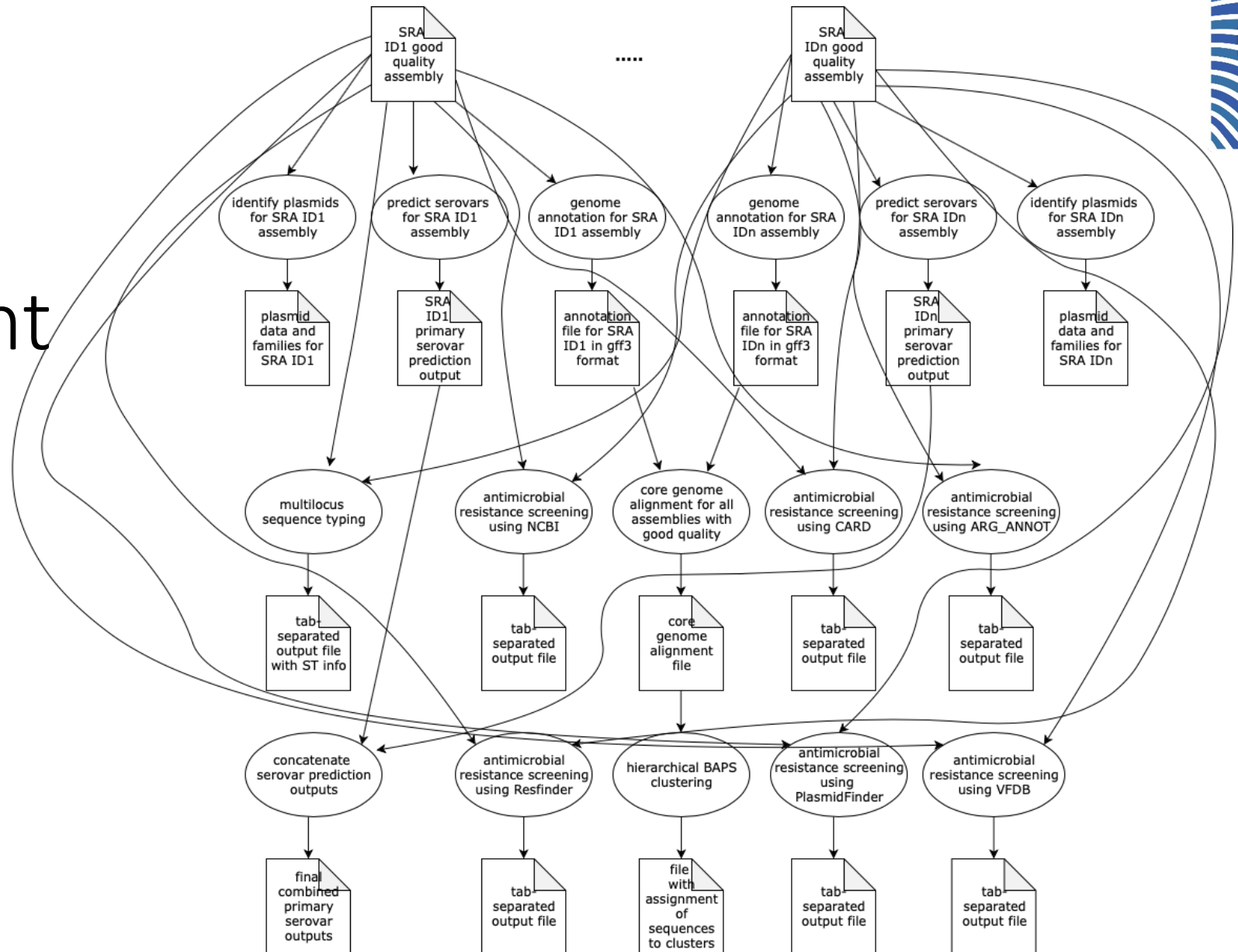
# ProkEvo

# ProkEvo

# ProkEvo

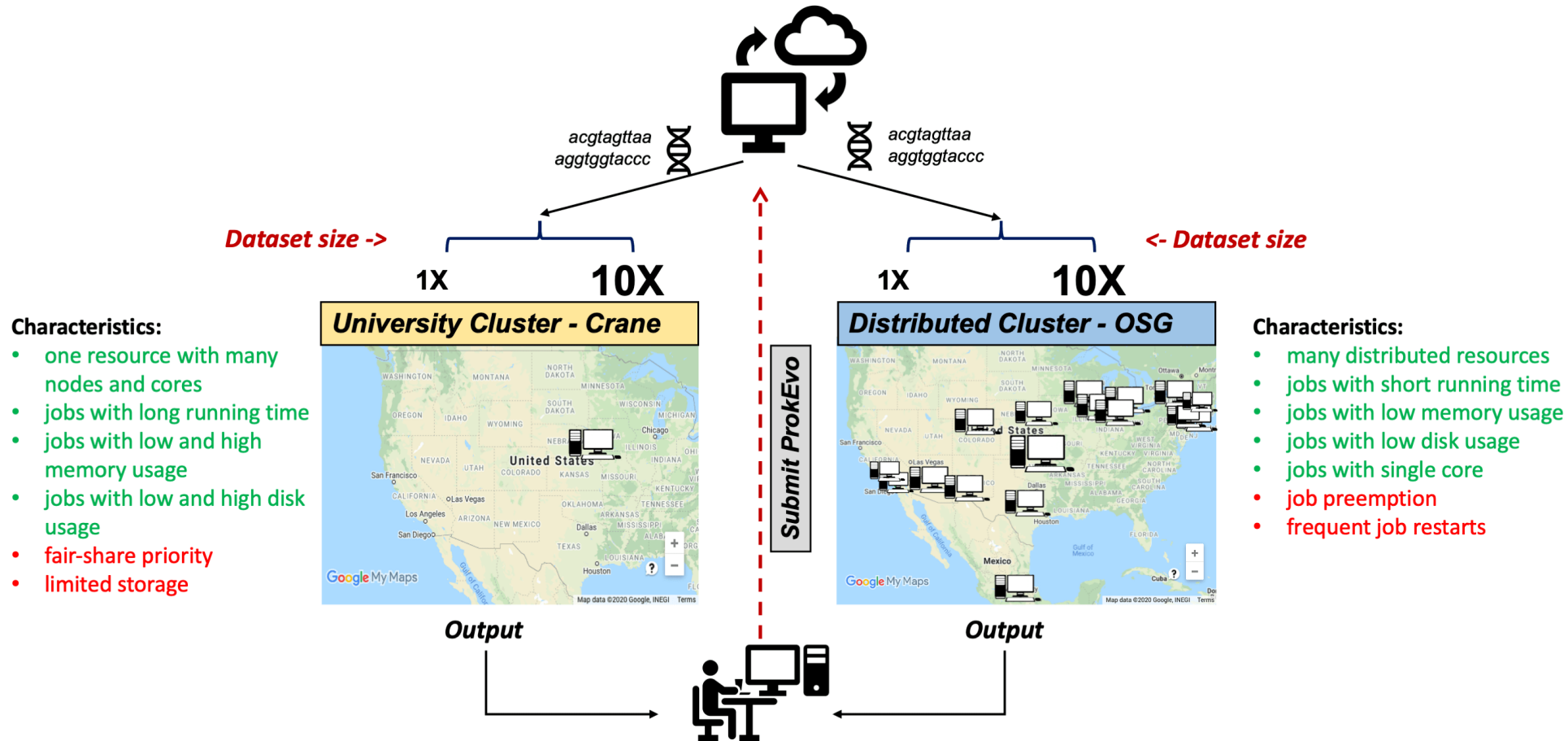# ProkEvo: Pegasus Workflow Management System

Part 1

# ProkEvo: Pegasus Workflow Management System

# Computational execution platforms



acgtagttaa
aggtggtaccc

acgtagttaa
aggtggtaccc

**Dataset size ->**

**<- Dataset size**

1X  **10X**

1X  **10X**

*University Cluster - Crane*

*Distributed Cluster - OSG*

**Submit ProkEvo**

**Characteristics:**
- one resource with many nodes and cores
- jobs with long running time
- jobs with low and high memory usage
- jobs with low and high disk usage
- fair-share priority
- limited storage

**Characteristics:**
- many distributed resources
- jobs with short running time
- jobs with low memory usage
- jobs with low disk usage
- jobs with single core
- job preemption
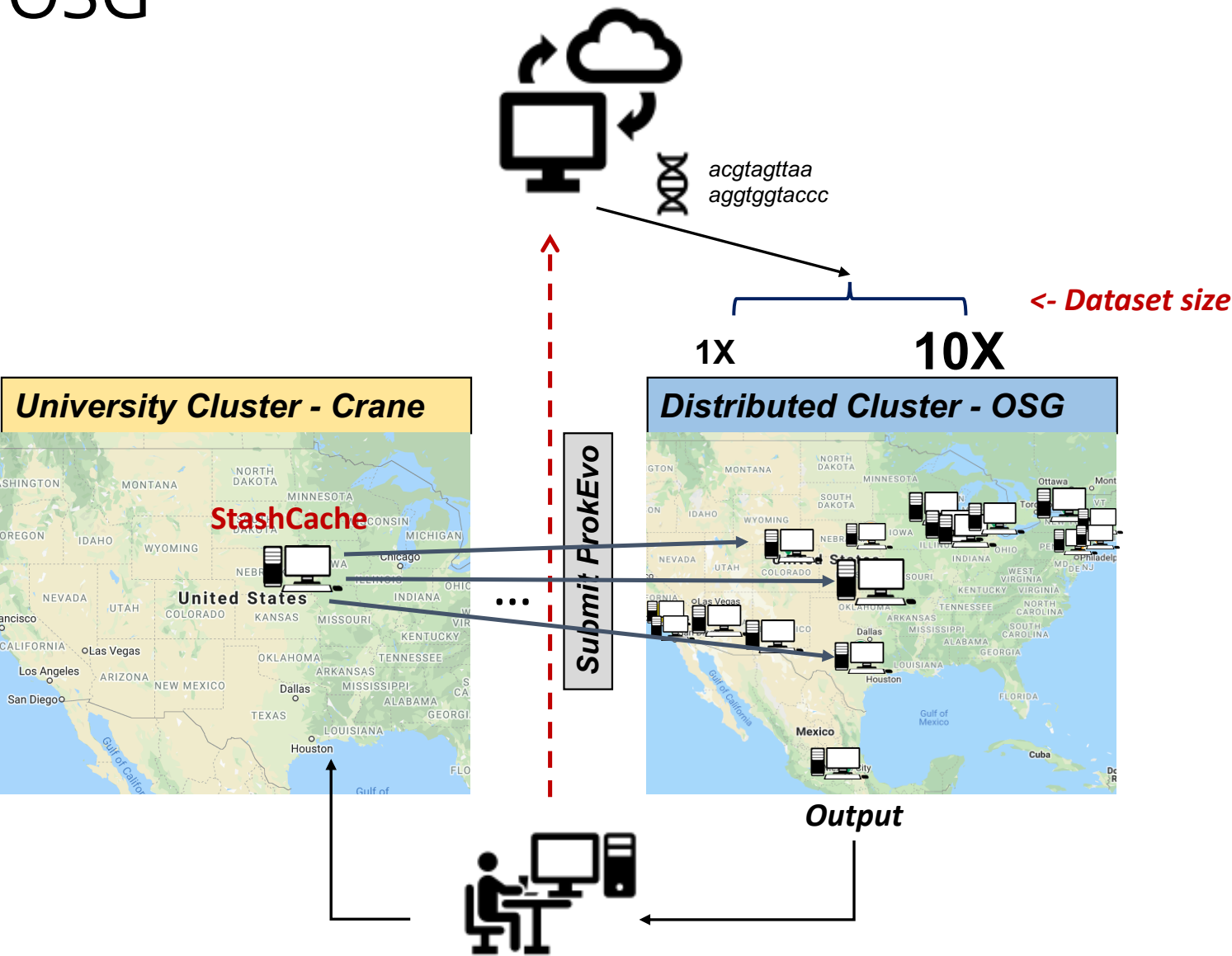- frequent job restarts

*Output*

*Output*

# ProkEvo on OSG

- OSG Connect
- Copy data to Crane
- Non-shared file-system
- Singularity
- Downloading data from NCBI
  - Intermittent network errors
  - Download limitations

# ProkEvo on OSG

# ProkEvo on HCC

- Store data to Crane

- Shared file-system

- Anaconda package manager

- Downloading data from NCBI
  - Intermittent network errors
  - Download limitations

# Results

- Performance evaluation

# Performance evaluation

| | Crane | OSG | Crane | OSG |
|---|---|---|---|---|
| **Number of genomes** | 2,392 | | 23,045 | |
| **Total <u>distributed</u> running time\*** | 3 days 15 hours | 7 days 4 hours | 15 days 22 hours | 26 days 6 hours |
| **Total <u>estimated sequential</u> running time\*\*** | 115 days 18 hours | 1 year 69 days | 2 years 268 days | 13 years 5 days |
| **Maximum jobs ran in a day\*\*\*** | 2,377 | 8,608 | 12,382 | 25,540 |
| **Total number of jobs ran** | 9,281 | 16,624 | 217,942 | 232,422 |
| **Output data size** | 131 GB | | 1.2 TB | |

\*Total distributed running time is calculated when many independent tasks are executed simultaneously utilizing single cores. This is the default behavior of ProkEvo.

\*\*Total estimated sequential running time is calculated when all steps from the pipeline are assumed to be run sequentially, on one single core.

\*\*\*The number of maximum jobs ran in a day depends on the type and length of the job, and is not linear.

# Conclusion

- ProkEvo
  - Automated, open-source pipeline for population genomics analyses that uses Pegasus Workflow Management System
  - Utilizing high-performance and high-throughput computational resources

- Applications

# Acknowledgement

# Future plans

- Submit ProkEvo from GUI (Science Gateway or OpenOnDemand) and run it on different cyberinfrastructures

# Our Pegasus Feedback

- Pegasus is awesome ☺

Notes:
- Modify executables after failure and re-run
- Use only Slurm instead of HTCondor?
- Singularity support for "sharedfs" data configurations
- If a job fails, don't mark the pipeline as failure, but only skip the tasks that follow after the failed job
- Command line options/wrappers (the tasks defined in the config files can be overridden from the command line)

# Thank You!