

Genomic Scale Data

Midwest BioDefense Summit

Lynn Langit

Big Data Considerations

Applied to Bioinformatics









- = Thymine
 - = Cytosine
 - 🔲 = Guanine

99.9% alike

VOLUME

3 Billion *.1 ≈ 300 Million

= Phosphate backbone

How big is 40 exabytes?

Genomics projects will generate 40 exabytes of data in the next decade.

Each shark = 100,000,000 GB of data

VELOCITY

<hr/>



Whole genome sequencing

Whole exome sequencing

Targeted sequencing











source



CCACCA---TGTGTGTGTGTGTGTGTGTGTGT > Contig 1096708082468

```
CCACCA------GTGTGTGT < Read 1094846374870
-----GTGTGTGT < Read 1089083146323
CCACCA-
                   -----
| | | | | | <sub>|</sub>_
CCACCAGTGTGTGTGTGTGTGTGTGTGTGTGT < Read 1089057878307
CCACCAGTGTGTGTGTGTGTGTGTGTGTGTGT < Read 1089053585558
CCACCAGTGTGTGTGTGTGTGTGTGTGTGTGT < Read 1094851158703
. . . . . .
           -----GTGTGTGTGT < Read 1095846055349
CCACCA-
```



Research Challenges

Analysis takes too long!
 Emergency Scalability!
 Out of storage space!*

*and more...



Analysis Phases



A high-level overview of NGS data processing



Use the Cloud

- Google Cloud: GCP
- Amazon: AWS
- Microsoft Azure





Practical Considerations for Developers



Dev Env



Sample Data



Pattern 1: Cloud-Native Data Lake







Traditionally the columns (features) of large "raw" datasets needed to be scaled down to allow advanced analytics to take place. This process of eliminating information potentially biases the result and impacts on the sensitivity of the analysis. VariantSpark is different, it can deal with the raw data directly resulting in a more accurate output.

Variant Spark was tested on datasets with 3000 samples, each containing 80 Million features, in either unsupervised clustering approaches (e.g. k-means) or supervised applications (e.g. Random Forests) with target/truth values that are categorical (classification) or continuous (regression). It can cluster data according to required profiles or identify predictive markers of events in just 30 minutes. By building a model on the full dataset, this allows, for the first time, to identify sets of markers that together have a stronger predictive power than the previously identified independent markers. Source



Results

1 Speed (30,000 -> 30 minutes)

(2)

Predictable Costs

Cloud Data Lake Reproducibility / Collaboration

Population-based study reveals strong genetic links with multiple cardiometabolic diseases and traits

By Nora Bradford

Broad scientists analyzed sequencing data from more than 200,000 people and found rare genetic variants associated with diseases like diabetes and heart failure.





Pattern 2: Serverless Data Analysis Pipeline



Serverless Pipeline



MVP

	1W	1M	ЗМ	6M	All	
Cumulative Volume =						
40M	 Cumulative Tests Completed Cumulative Positives 			Tuesday, Jan 24, 00:00 Cumulative Tests Complet Cumulative Positives:	ted: 37,408,482 906,858	
30M						
10M						
0						
0	Apr-1	Sep_1	May-1	pec-1 han-1	Nov-1 Sep-1 Sep-1	



Results

Scalability (0-> 37 M Tests)!

2

Zero Down-time

Predictable Costs

Serverless Pipeline

Analyze Everything!







T1037 / 6vr4 90.7 GDT (RNA polymerase domain) **T1049 / 6y4f** 93.3 GDT (adhesin tip)

Experimental result

Computational prediction

Spatial Multiomics









Pipelines



Data Mesh





Results - Space + Speed

File Storage

Data Lakes Grouping Lakes into DataMesh

2

Experiments / Analytics Serverless Pipelines

2 weeks from idea to deployment

Learn on GitHub

Notes, Code and Videos

- TeamTeri Repo
- GCP-for-Bioinformatics Repo
- AWS-for-Bioinformatics Repo
- Learn Quantum Programming Repo
- Learn Apache Spark Repo

📮 lynnlangit / TeamTeri Public	nlangit / TeamTeri Public					
<> Code 🕑 Issues 🕅 Pull reque	ests 🕞 Actions [] Projects 1	🕮 Wiki 🕕 Security 🗠 Insights 🕸				
위 master → 위 1 branch 💿 0 tag	S	Go to file Add file - Code -				
Iynnlangit Update README.md		5329a15 4 days ago 🕚 397 commits				
1_Concepts_Terms	Update README.md	5 months ago				
2_Lab_Testing	Update README.md	12 months ago				
3_Genomic_Tools	Update README.md	2 months ago				
4_Genomics_Platforms	Update README.md	4 days ago				
5_Public_Cloud_Genomics	Update README.md	2 months ago				
Images	added nf-on-aws concepts	4 days ago				
🗅 .gitignore	Cleanup ignore	5 years ago				
🗅 README.md	Update README.md	2 months ago				

