

Scalable solutions for genomic data analysis

Tomasz Gambin

Institute of Computer Science, Warsaw University of Technology



<http://biodatageeks.org>



<http://github.com/ZSI-Bio>



<http://twitter.com/biodatageeks>

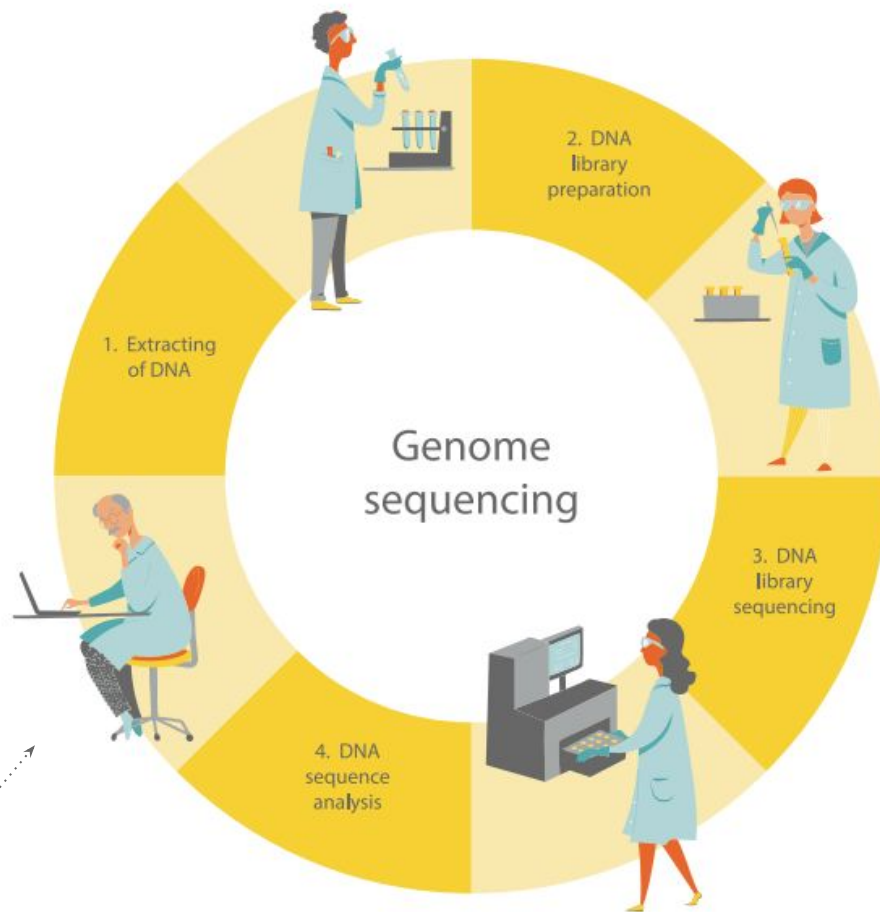
About us

Areas of interest:

- Genomic data analysis
- Distributed solutions for genomics

Applications:

- Distributed pipelines for variant calling
- Depth of coverage analysis, quality control
- Distributed range joins implementation
- Efficient variant data warehousing solutions
- Interpretation tools



That's us

Cumulative Samples: N= 6,599
Sequenced at BCM= ~5,000



- Institutes: 125
- Investigators: 349
- Countries: 26

>300 novel disease genes

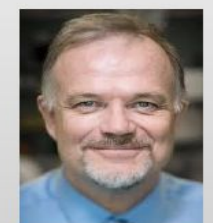
Chong et al. *AJHG* , 2015



James Lupski



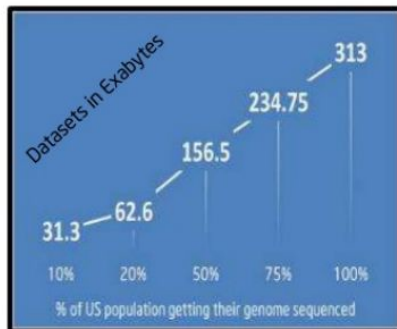
Eric Boerwinkle



Richard Gibbs

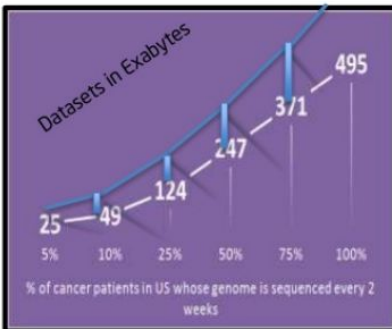
Genomics - Big Data Problem

The day when every newborn gets their DNA sequenced is not far away: <http://www.nih.gov/news/health/sep2013/nhgri-04.htm>.



313 Exabytes

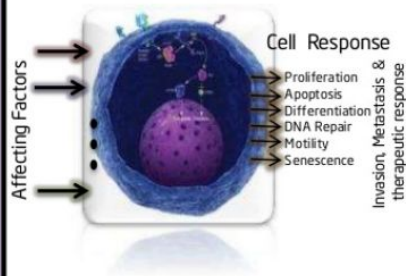
if everyone in the US has their genes sequenced



495 Exabytes

if every cancer patient in the US has their genes sequenced every 2 weeks.

Images, Assays and Drug response data will push it further up as shown in Blue line



Complex interaction of *varied & changing* intrinsic and extrinsic factors determine cell response

With Genomic Data growing rapidly, hospitals and research centers need to access the local data (the ones not shared) and the centralized public/private data for various analysis and analytics for Genomic Research/Development/Medicine.

**Compute has to be done "where data is" and need to be consistent locally and in the cloud.
Energy, Total Cost of Operation are key**

Source: Knights Cancer Institute, Oregon Health Sciences University & Intel

Figure: Genomics as big data problem, Source: Knights Cancer Institute, Oregon Health Sciences University & Intel

Challenges

Performance

Computation intensive steps in pipelines and painful long lasting operations in data analysis.



Data size

Heavy files, inefficient data access, temporary files generation. High storage cost



Data security

'All or nothing' approach is not enough for research projects or clinics

Multi-sample analysis

Data standardization and merging as additional overhead

Our solutions

Distributed calculations

Reimplementation of algorithms using Apache Spark and other Big Data tools



Unified data model

Providing ANSI SQL-compliant interfaces
Table-oriented processing

Optimized algorithms

Ensuring scalability to handle population-scale analysis



Standard technologies

Fine-grained access control

SeQuiLa (range joins)- scalable intersection of interval sets

e.g: What variants (snps) occur WITHIN genes

chr	name	start	end
11	Gene1	100	150
11	Gene2	200	250
11	Gene3	300	350
12	Gene4	100	150

chr	start	end	ref	alt	af
11	101	101	C	T	0.5
11	104	105	CT	C	0.01
11	134	135	AA	TT	0.01
11	201	201	A	G	0.05
12	102	102	T	G	0.05
13	1004	1005	TA	CG	0.04
13	2004	2005	TA	CG	0.04

```
SELECT g.chr, g.name, g.start, g.end, s.start,
s.end, s.af
FROM genes g JOIN snps s ON (
  g.chr = s.chr AND s.start >= g.start AND s.end
<= g.end)
```

- counting overlaps
- additional criterias on overlap (maxGap, minOverlap)

chr	name	start	end	start	end	ref	alt	af
11	Gene1	100	150	101	101	C	T	0.5
11	Gene1	100	150	104	105	CT	C	0.01
11	Gene1	100	150	134	135	AA	TT	0.01
11	Gene2	200	250	201	201	A	G	0.05
12	Gene4	100	150	102	102	T	G	0.05

real genomic example: $(160 * 10^6) \times (200 * 10^3)$ or even:
 $(2,6 * 10^9) \times (200 * 10^3)$

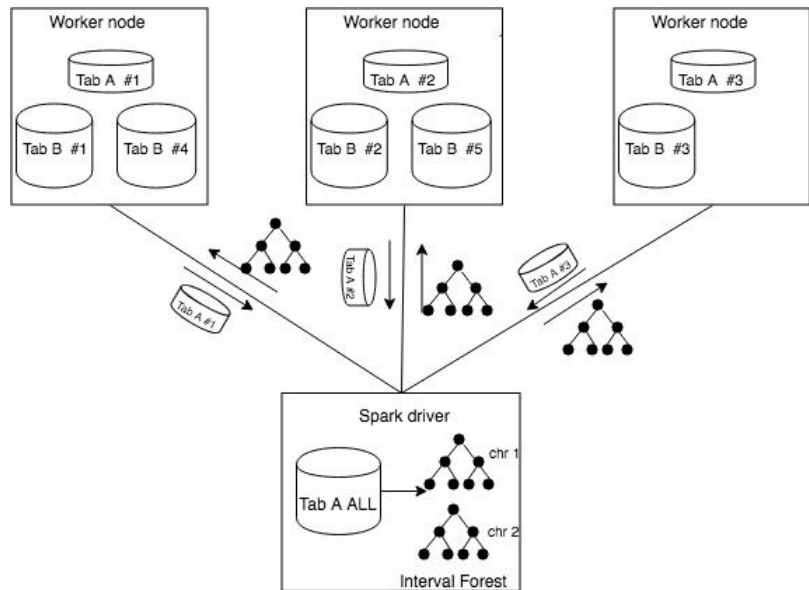
SeQuiLa (range joins): methods

Extension of Catalyst (SparkSQL component)

1. **IntervalTree** structure is used for efficient overlaps search
 - a. Interval Forest (one tree for each chromosome)
2. augmenting IntervalTree with table data if possible

Algorithm for range join table A (small) with table B (big):

1. Send to driver node all table A partitions
2. Build Interval Forest in driver node
3. Broadcast Interval Forest to all worker nodes
4. Perform interval search
5. Join search results with table A if necessary



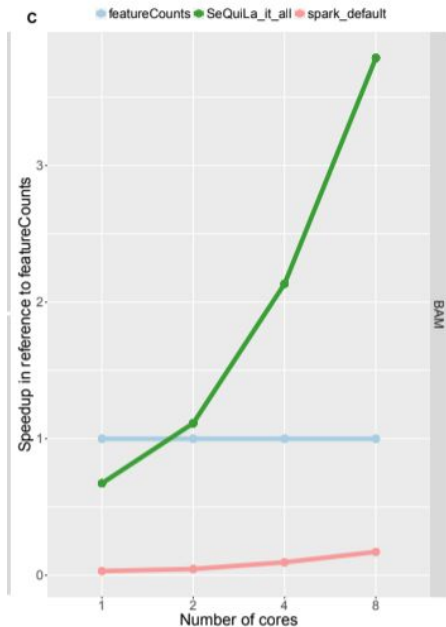
SeQuiLa (range joins)

SeQuiLa: an elastic, fast and scalable SQL-oriented solution for processing and querying genomic intervals

Marek Wiewiórka, Anna Leśniewska, Agnieszka Szmurło, Kacper Stępień, Mateusz Borowiak, Michał Okoniewski, Tomasz Gambin ✉

Bioinformatics, bty940, <https://doi.org/10.1093/bioinformatics/bty940>

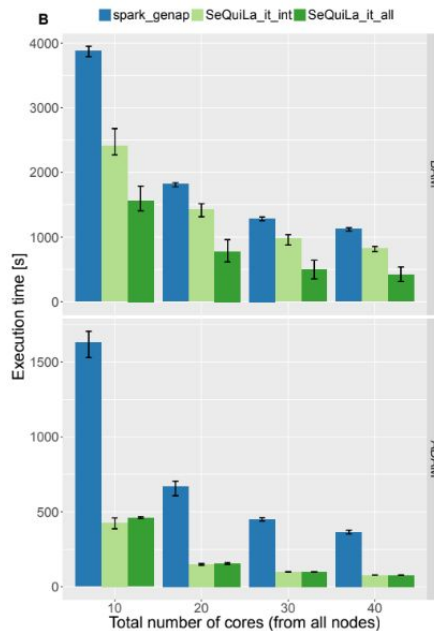
Published: 14 November 2018 **Article history** ▾



- single node
- data: WES (17 GB)
- reads ($160 \cdot 10^6$) x targets ($200 \cdot 10^3$)

Benchmark against:

- featureCounts
- SparkGenap
- spark default



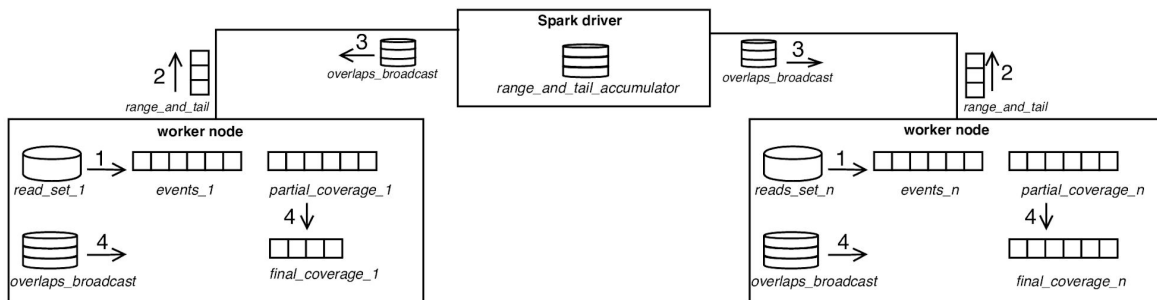
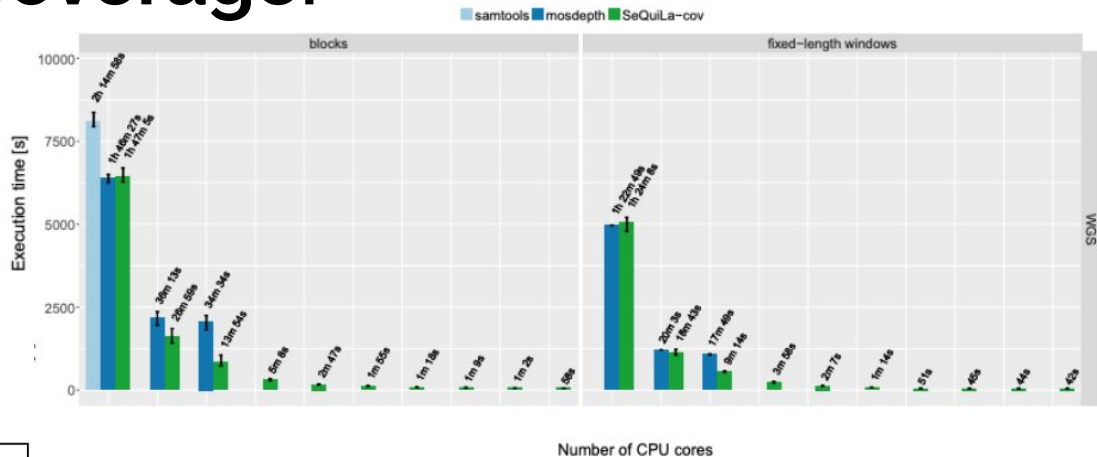
- cluster (4 worker nodes)
- data: WGS (270 GB)
- reads ($2,6 \cdot 10^9$) x targets ($200 \cdot 10^3$)

Benchmark against:

- SparkGenap

Distributed Depth of Coverage:

- ✓ Distributed calculations
- ✓ Simple event-based algorithm
- ✓ Low level optimizations
- ✦ Standalone version



2019

SeQuiLa-cov: A fast and scalable library for depth of coverage calculations

Marek Wiewiórka^{1,*}, Agnieszka Szmurło^{1,*}, Wiktor Kuśmirek¹ and Tomasz Gambin^{1,†}

¹Institute of Computer Science, Warsaw University of Technology, ul. Nowowiejska 15/19, 00-665 Warsaw, Poland,



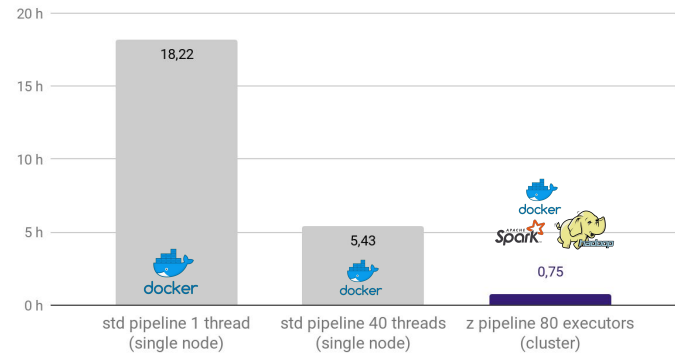
<http://biodatageeks.org/sequila/>

Distributed variant annotation pipeline:

- ✓ Automatic execution
- ✓ Customizations
- ✓ Monitoring of task execution
- ✦ Distributed calculations

Benchmark:

VCF pipeline processing time [hours]

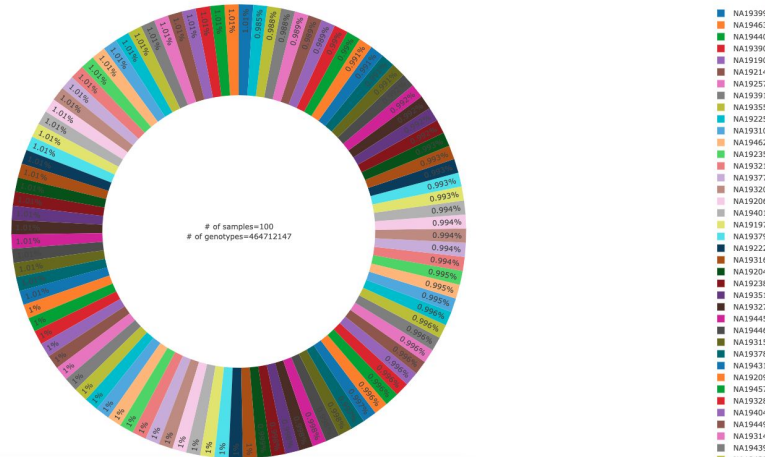


<input type="checkbox"/>	Off	cleanup ⓘ		zettagene				
<input checked="" type="checkbox"/>	On	cov_only	None	zettagene	2	2019-03-18 13:02 ⓘ	1	
<input checked="" type="checkbox"/>	On	coverage_cleanup	None	zettagene	4	2019-03-19 09:37 ⓘ	8 9	
<input checked="" type="checkbox"/>	On	coverage_pipeline	None	zettagene	5	2019-03-20 13:05 ⓘ	3 9	
<input checked="" type="checkbox"/>	On	vcf_cleanup	None	zettagene				
<input checked="" type="checkbox"/>	On	vcf_pipeline	None	zettagene	7	2019-03-08 19:10 ⓘ	50 79	

Interpretation tools:

- ✓ Tabular view on variants and genotypes
- ✓ Charts, statistics, breakdowns
- ✓ Fine grained access control
- ✦ IGV view of variants and aligned reads

Displaying maximum of 1000 rows Statistics of selected variant calls Statistics of all variant calls Population breakdown



Update View

Sample filters

Variant filters

Impact: HIGH MODERATE

Clin. significance: benign

Consequence:

Gene: ADACL2

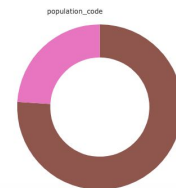
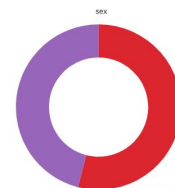
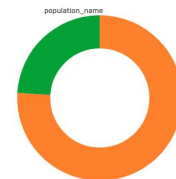
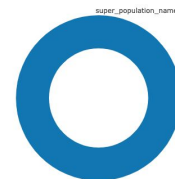
max AF: 0.05 - 0.1

Showing 1 to 15 of 781 entries

sample_id	gt	chr	pos	ref	alt	variant_id	af	gnomad_af	max_af	impact	consequence	gene_symbol	exon	aminoacids	zift	
1	NA19438	A/G	18	2890770	G	A	rs16943977	0.0219649	0.004643	0.0809	MODERATE	missense_variant	EMILIN2	4/8	A/T	tolerated(1)
2	NA19438	A/C	18	2890855	C	A	rs30646935	0.0151757	0.003314	0.056	MODERATE	missense_variant	EMILIN2	4/8	T/K	tolerated(0.05)
3	NA19438	A/C	18	5891054	C	A	rs577863972	0.00239585	0.009921	0.06699	MODERATE	missense_variant	TMEM200C	3/3	R/L	tolerated(0.05)
4	NA19438	A/G	18	5891055	G	A	rs559244640	0.00239585	0.009825	0.06634	MODERATE	missense_variant	TMEM200C	3/3	R/W	deleterious(0.01)
5	NA19438	G/A	18	6908977	G	A	rs116232392	0.0249601	0.00633	0.0877	MODERATE	missense_variant	ARHGAP28	16/17	C/Y	tolerated(0.19)
6	NA19438	J/A	18	8784372	A	T	rs918272	0.0231629	0.007106	0.0779	MODERATE	missense_variant	MTCL1	6/17	E/D	tolerated(0.75)
7	NA19438	G/C	18	9221885	C	G	rs2298548	0.0463259	0.06377	0.08982	MODERATE	missense_variant	ANKRD12	8/13	P/A	tolerated(0.31)
8	NA19438	T/C	18	9254787	C	T	rs17489752	0.0461262	0.06357	0.0878	MODERATE	missense_variant	ANKRD12	9/13	T/I	deleterious_low_conf
9	NA19438	A/C	18	9255786	C	A	rs72939232	0.0595048	0.06773	0.09	MODERATE	missense_variant	ANKRD12	9/13	T/N	tolerated_low_conf
10	NA19438	C/T	18	9258539	T	C	rs3744822	0.0621006	0.067	0.09986	MODERATE	missense_variant	ANKRD12	9/13	S/P	tolerated(0.32)
11	NA19438	A/G	18	23421436	G	A	rs76709352	0.0239585	0.00561	0.0946	MODERATE	missense_variant	TMEM241	3/15	T/M	deleterious(0.03)
12	NA19438	G/A	6	158066494	A	G	rs61601143	0.0241613	0.004605	0.0885	MODERATE	missense_variant	SYNJ2	12/27	N/D	tolerated(0.19)
13	NA19438	G/A	6	168030801	G	A	rs34049091	0.0215655	0.00498	0.0802	MODERATE	missense_variant	KIF25	7/13	A/T	tolerated(0.36)
14	NA19438	A/G	6	168062944	G	A	rs902393	0.0696885	0.071	0.092	MODERATE	missense_variant	FRMD1	7/11	R/C	tolerated(0.14)
15	NA19438	C/G	6	170584376	G	C	rs74482927	0.0211661	0.004123	0.0779	MODERATE	missense_variant	PDCD2	1/6	P/R	tolerated(0.13)

Showing 1 to 15 of 781 entries

Previous 1 2 3 4 5 ... 53 Next



African
Luhya
Yoruba
female
male
LWK
YRI

Bioinformatics pipelines



NGS variant calling



Coverage QC



RNA-Seq



Data analysis



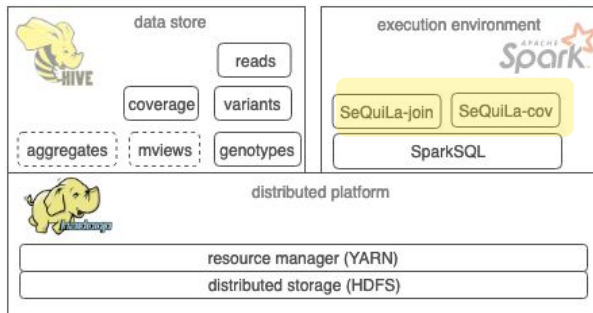
Interpretation tools



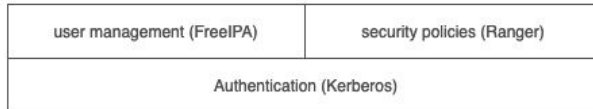
Custom analyses



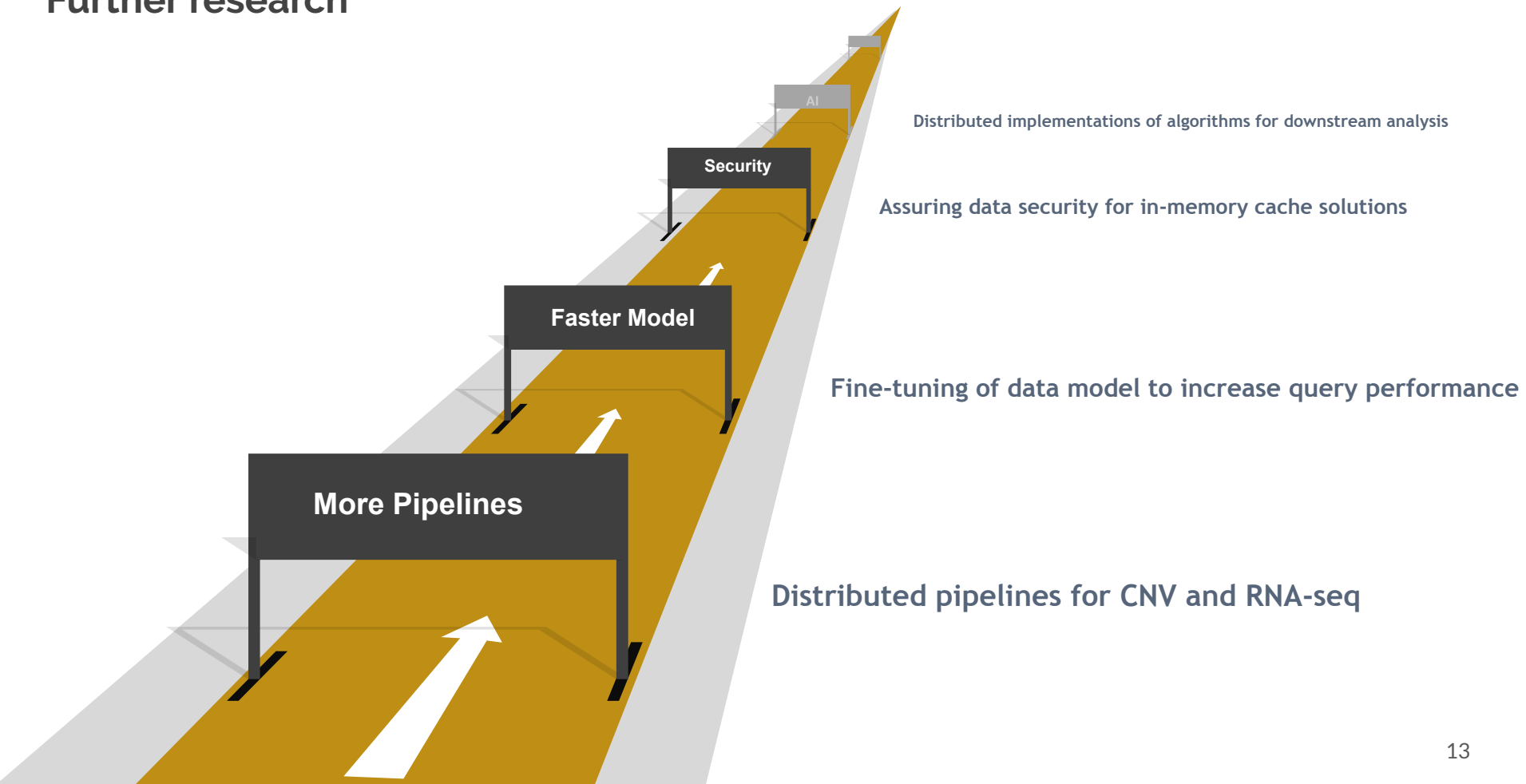
Data



Security



Further research





Thanks!

Any questions?

You can find us at blodatageeks.org