



# Understanding and Characterization of Pangenomics

Noah Kaplan, Yufeng Gu, Reetuparna Das

# Linear Alignment

---

Reference: ...G T G T A T G T C C A T G T G T C A T T G T G T C A A T G C T T C T C ...

Read: A T G T G T C A

# Linear Alignment

---

Reference: ...G T G T A T G T C C A T G T G T C A T T G T G T C A A T G C T T C T C ...

Read: A T G T G T C A



# The Variant Problem

---

Reference: ...G T G T A T G T C C A T G T G T C A T T G T G T C A A T G C T T C T C...

Read: C T G G G C T A

# The Variant Problem

---

Reference: ...G T G T A T G T C C A T G T G T C A T T G T G T C A A T G C T T C T C ...

Read: C T G G G C T A

# Pangenomes

---

Reference: ...G T G T A T G T C C A T G T G T C A T T G T G T C A A T G C T T C T C...

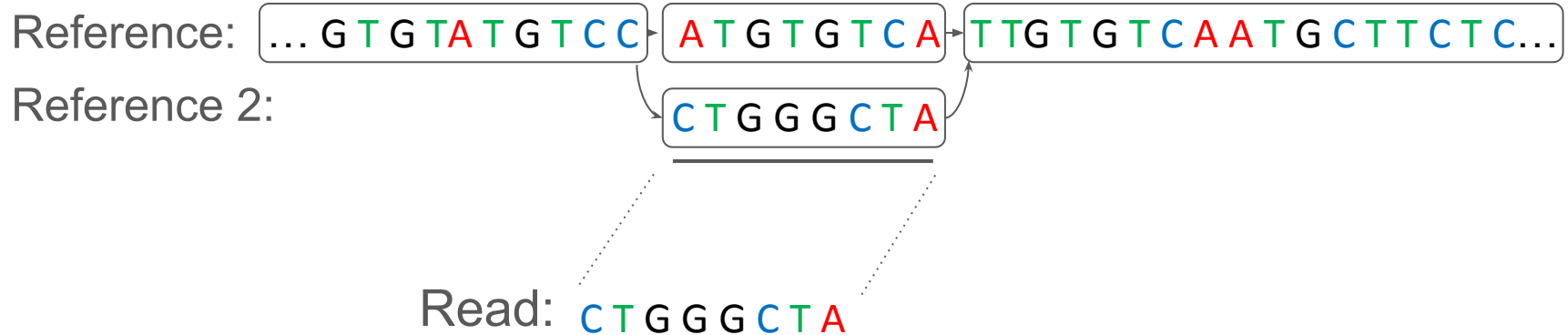
Reference 2: ...G T G T A T G T C C C T G G G C T A T T G T G T C A A T G C T T C T C...

Read: C T G G G C T A



# Pangenomes

---

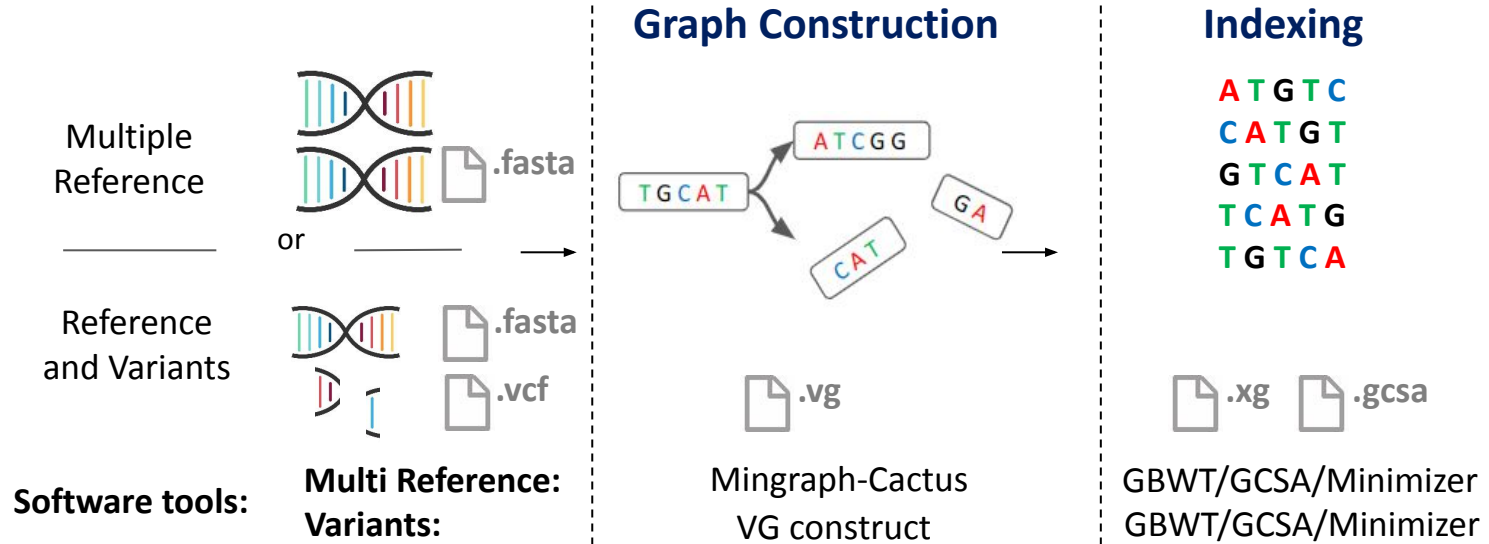


# Outline

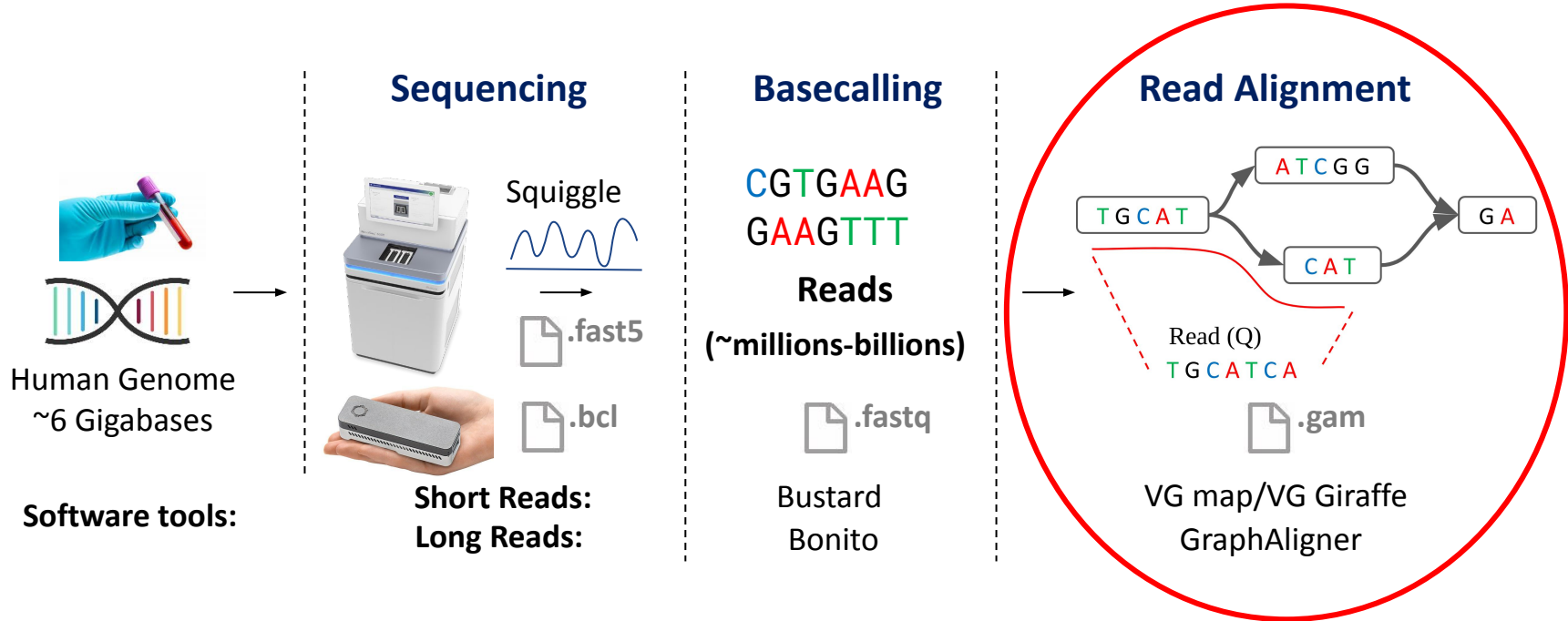
- **Workflows in Pangenomics**
- **Common Tools**
  - Vg Map (Short Reads)
  - GraphAligner (Long Reads)
- **Profiling Results**
- **Conclusion**



# Pangenomics Reference Graph Building



# Pangenomics Alignment Workflow



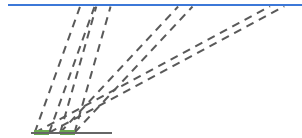
# Outline

- **Workflows in Pangenomics**
- **Common Tools**
  - Vg Map (Short Reads)
  - GraphAligner (Long Reads)
- **Profiling Results**
- **Conclusion**

# Algorithms in Vg Map

```
$ A T G G T T
A T G G T T $
G G T T $ A T
G T T $ A T G
T $ A T G G T
T G G T T $ A
T T $ A T S G
```

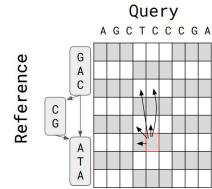
Seeding  
(GCSA)



Clustering  
(Distance Based)



Cluster Filtration  
(Overlap Heuristic)



Extension  
(GSSW)

# Seeding: GCSA

---

Reference: **A**TGGTT\$

\$	A	T	G	G	T	T
<b>A</b>	T	G	G	T	T	\$
G	G	T	T	\$	A	<b>T</b>
G	T	T	\$	A	T	G
<b>T</b>	\$	A	T	G	G	<b>T</b>
<b>T</b>	G	G	T	T	\$	<b>A</b>
<b>T</b>	T	\$	A	T	G	G

# Seeding: GCSA

---

Query: GGT

\$	A	T	G	G	T	T
A	T	G	G	T	T	\$
G	G	T	T	\$	A	T
G	T	T	\$	A	T	G
T	\$	A	T	G	G	T
T	G	G	T	T	\$	A
T	T	\$	A	T	G	G

# Seeding: GCSA

---

Query: GGT

\$	A	T	G	G	T	T
A	T	G	G	T	T	\$
G	G	T	T	\$	A	T
G	T	T	\$	A	T	G
T	\$	A	T	G	G	T
T	G	G	T	T	\$	A
T	T	\$	A	T	G	G

# Seeding: GCSA

---

Query: GGT

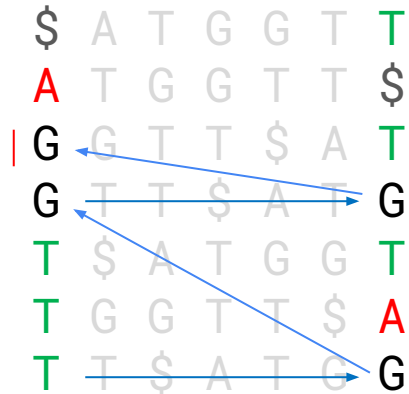
\$	A	T	G	G	T	T
A	T	G	G	T	T	\$
G	G	T	T	\$	A	T
G	T	T	\$	A	T	G
T	\$	A	T	G	G	T
T	G	G	T	T	\$	A
T	T	\$	A	T	G	G



# Seeding: GCSA

---

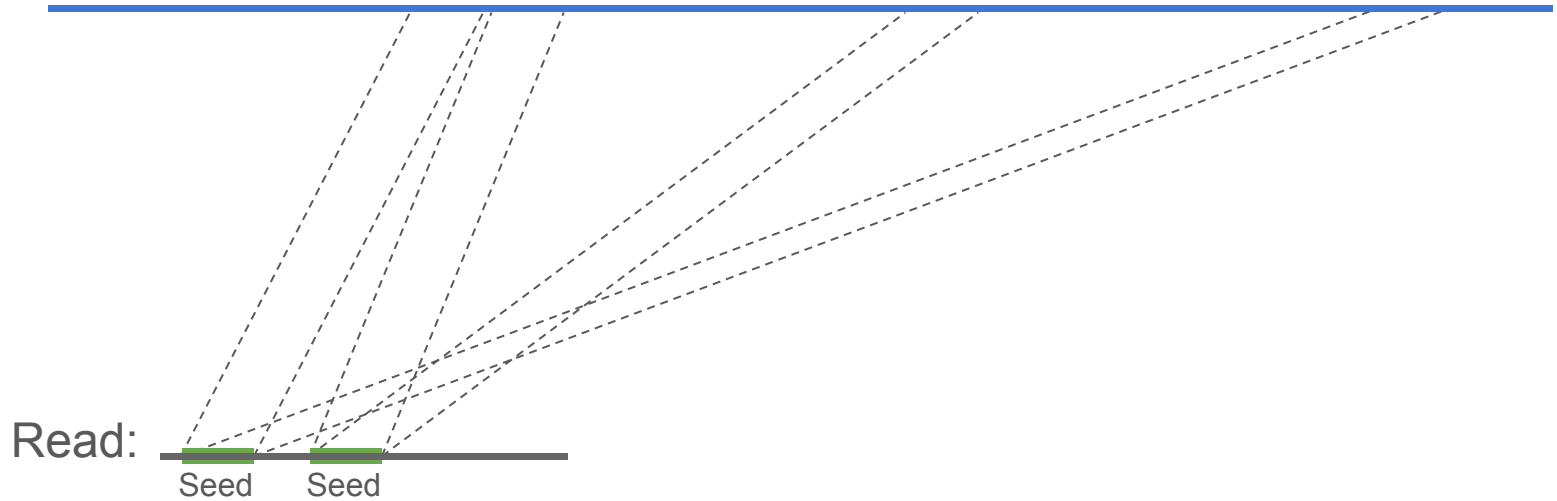
Query: GGT



# Clustering

---

Reference:



# Cluster Filtration

---

- If a cluster is too short, drop it
- If a cluster overlaps with another cluster, drop it



# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

A	G	C	T	C	C	C	G	A	
									G
									A
									G
									T
									...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

Query

	A	G	C	T	C	C	C	G	A	
	0	3	1	0	0	0	0	3	1	G
	3	1	0	0	0	0	0	1	6	A
	1	6	4	3	2	1	0	4	4	G
	0	4	3	7	5	4	3	2	2	T
	...	...	...	...	...	...	...	...	...	...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

	A	G	C	T	C	C	C	G	A	
	0	3	1	0	0	0	0	3	1	G
	3	1	0	0	0	0	0	1	6	A
	1	6	4	3	2	1	0	4	4	G
	0	4	3	7	5	4	3	2	2	T
	...	...	...	...	...	...	...	...	...	...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

Query

	A	G	C	T	C	C	C	G	A	
	0	3	1	0	0	0	0	3	1	G
	3	1	0	0	0	0	0	1	6	A
	1	6	4	3	2	1	0	4	4	G
	0	4	3	7	5	4	3	2	2	T
	...	...	...	...	...	...	...	...	...	...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

A	G	C	T	C	C	C	G	A	
									G
									A
									G
									T
									...

Reference



# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

	A	G	C	T	C	C	C	G	A	
1										G
										A
										G
										T
										...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

	A	G	C	T	C	C	C	G	A	
1	2									G
										A
										G
										T
										...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

	A	G	C	T	C	C	C	G	A	
1	2	3								G
										A
										G
										T
										...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

	A	G	C	T	C	C	C	G	A	
1										G
2										A
3										G
4										T
5										...

# Extension: Graph Simd Smith-Waterman (GSSW)

---

Query

	A	G	C	T	C	C	C	G	A	
1										G
10										A
										G
										T
										...

# Extension: Graph Simd Smith-Waterman (GSSW)

---

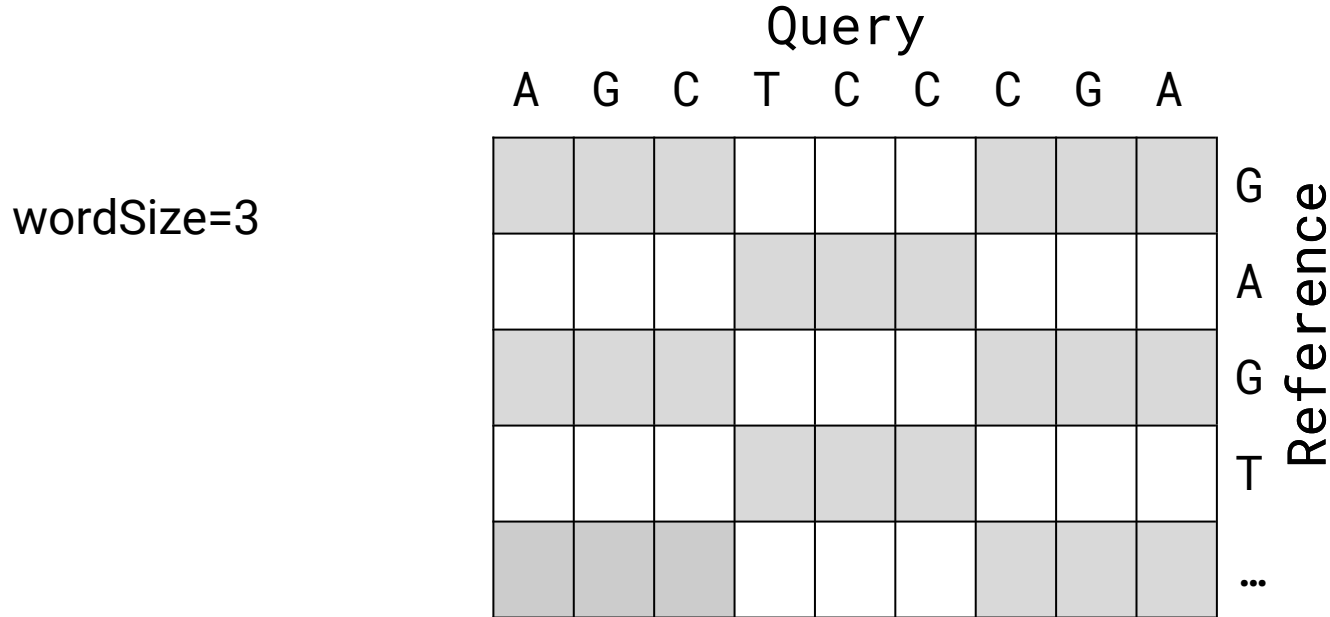
Query

	A	G	C	T	C	C	C	G	A	
1										G
2										A
3										G
4										T
5										...
6										
7										
8										
9										
10										

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

---



# Extension: Graph Simd Smith-Waterman (GSSW)

wordSize=3

		Query									
		A	G	C	T	C	C	C	G	A	
G A G T ...	G	1			1			1			
	A										
	G										
	T										
	...										



# Extension: Graph Simd Smith-Waterman (GSSW)

wordSize=3

Query

	A	G	C	T	C	C	C	G	A	
	1	2		1	2		1	2		G
										A
										G
										T
										...

# Extension: Graph Simd Smith-Waterman (GSSW)

wordSize=3

Query

	A	G	C	T	C	C	C	G	A	
	1	2	3	1	2	3	1	2	3	G
										A
										G
										T
										...

# Extension: Graph Simd Smith-Waterman (GSSW)

wordSize=3

Query

	A	G	C	T	C	C	C	G	A	
	1	2	3	1	2	3	1	2	3	G
	4	5	6	4	5	6	4	5	6	A
	7			7			7			G
										T
										...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

wordSize=3

Query

	A	G	C	T	C	C	C	G	A	
	1	2	3	1	2	3	1	2	3	G
	4	5	6	4	5	6	4	5	6	A
	7			7			7			G
										T
										...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)

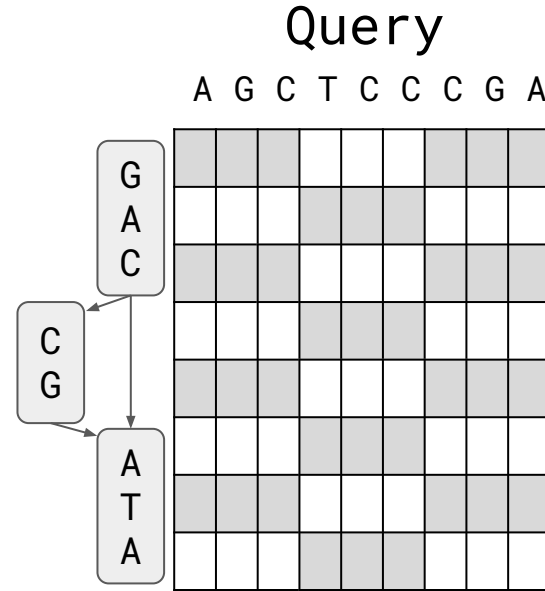
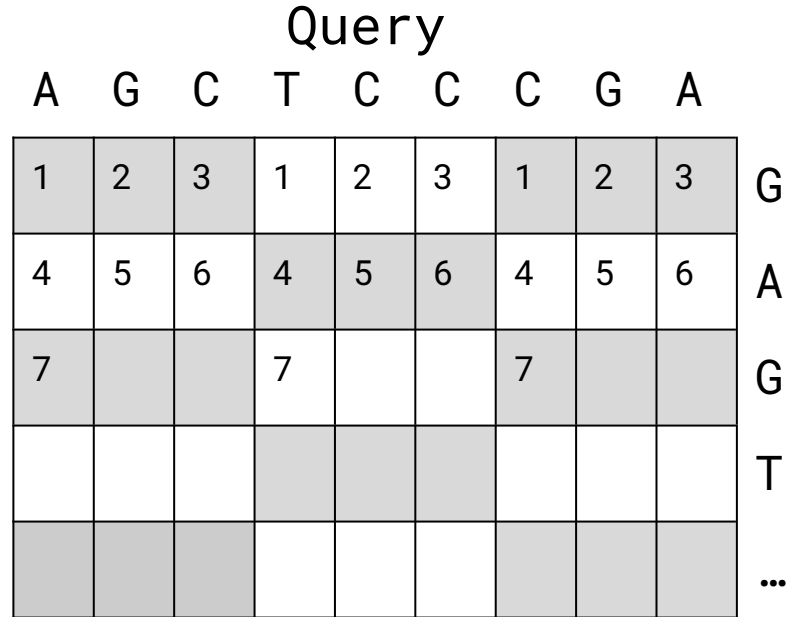
wordSize=3

Query

	A	G	C	T	C	C	C	G	A	
	1	2	3	1	2	3	1	2	3	G
	4	5	6	4	5	6	4	5	6	A
	7			7			7			G
										T
										...

Reference

# Extension: Graph Simd Smith-Waterman (GSSW)



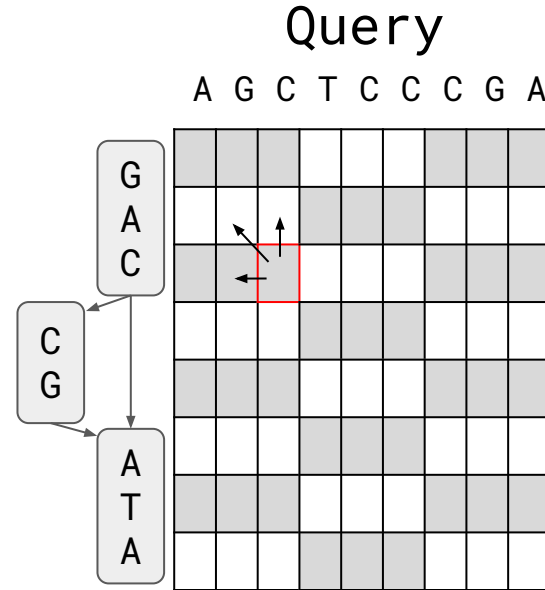
# Extension: Graph Simd Smith-Waterman (GSSW)

Query

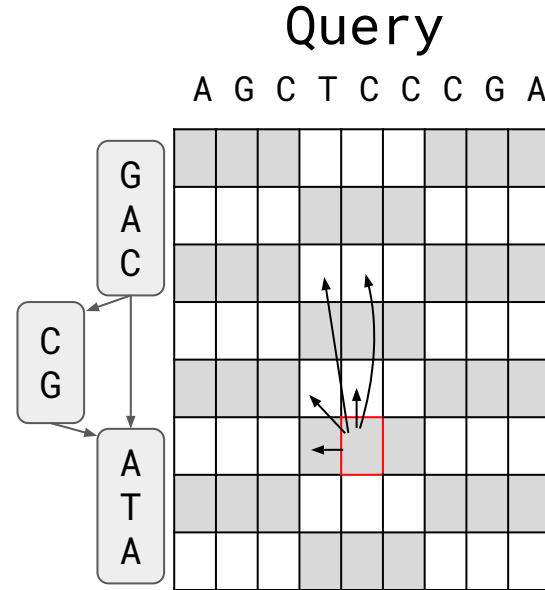
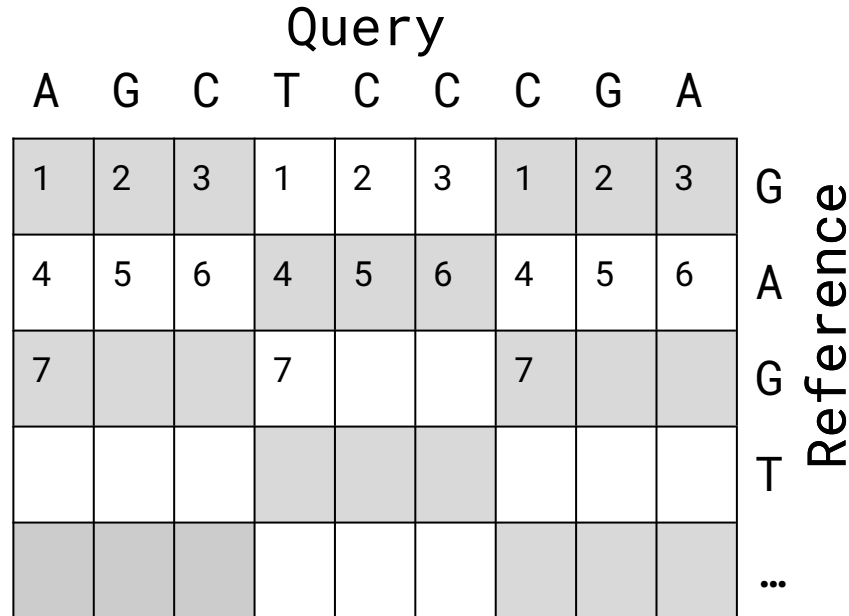
	A	G	C	T	C	C	C	G	A
1	2	3	1	2	3	1	2	3	
4	5	6	4	5	6	4	5	6	
7			7			7			

Reference

G  
A  
G  
T  
...



# Extension: Graph Simd Smith-Waterman (GSSW)





# Outline

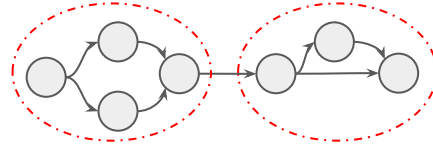
- **Workflows in Pangenomics**
- **Common Tools**
  - Vg Map (Short Reads)
  - **GraphAligner (Long Reads)**
- Profiling Results
- Conclusion

# Algorithms in GraphAligner

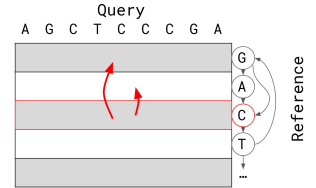
---

CATG**T**GTCATTG**T**GTCA  
↓  
A**T**G**T**

Seeding  
(Minimizer)



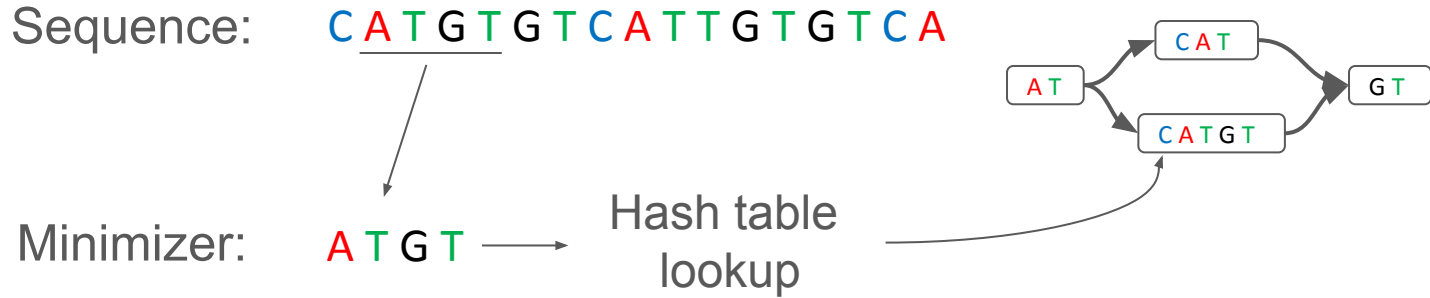
Clustering  
(Distance Based)



Extension  
(Myers Bitvector)

# Seeding: Minimizer

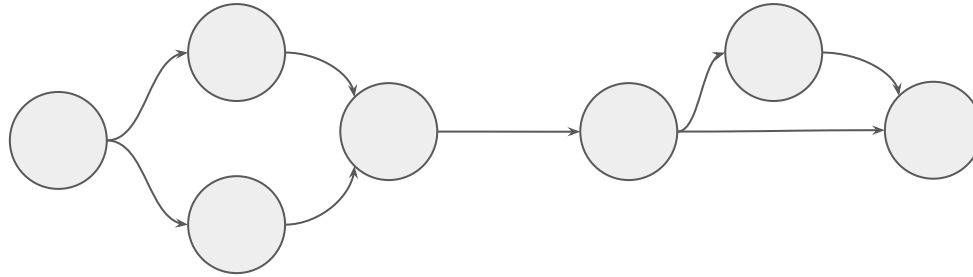
---



# Clustering

---

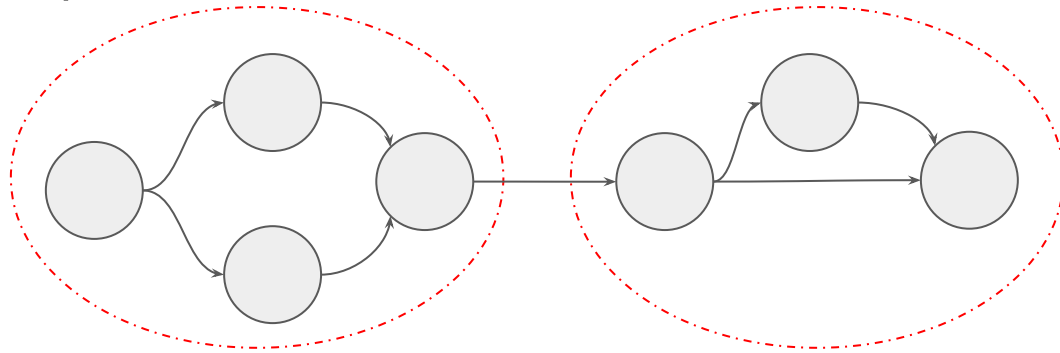
Reference Graph:



# Clustering

---

Reference Graph:



Create linear chains of  
“superbubbles”

# Extension: Myers Bitvector

---

Query

A	G	C	T	C	C	C	G	A	
									G
									A
									C
									T
									⋮

Reference



# Extension: Myers Bitvector

---

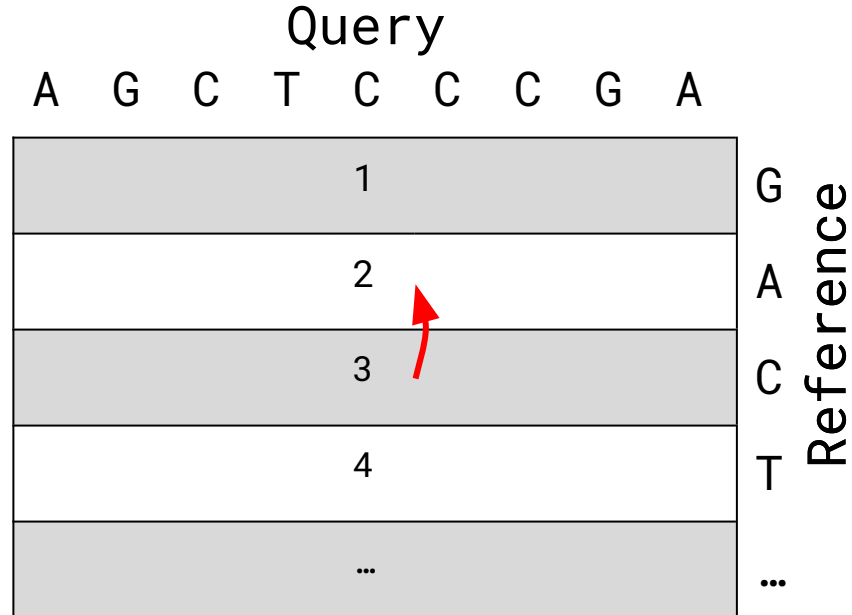
Query										
A	G	C	T	C	C	C	G	A		
				1						G
				2						A
				3						C
				4						T
				...						...

Reference



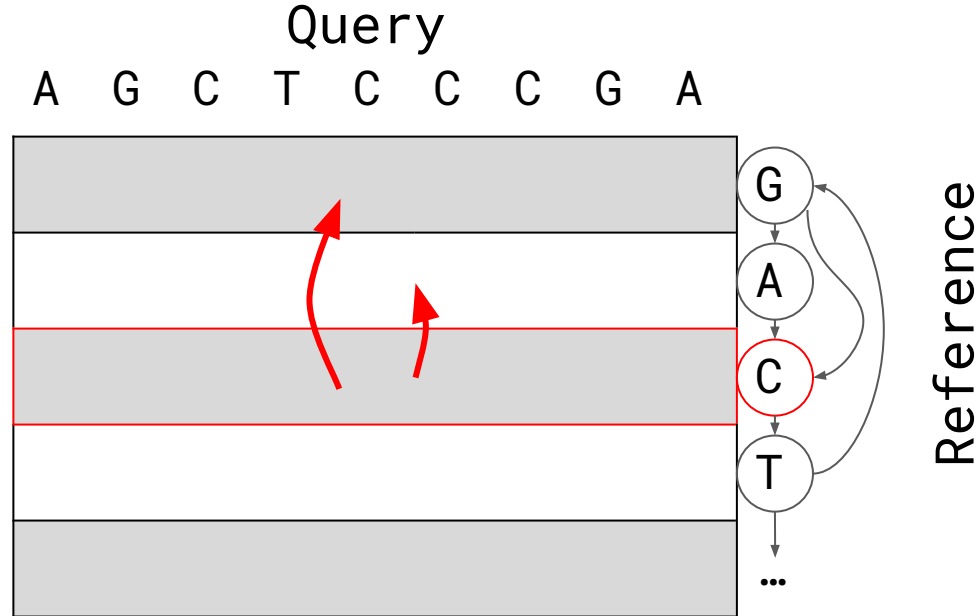
# Extension: Myers Bitvector

---



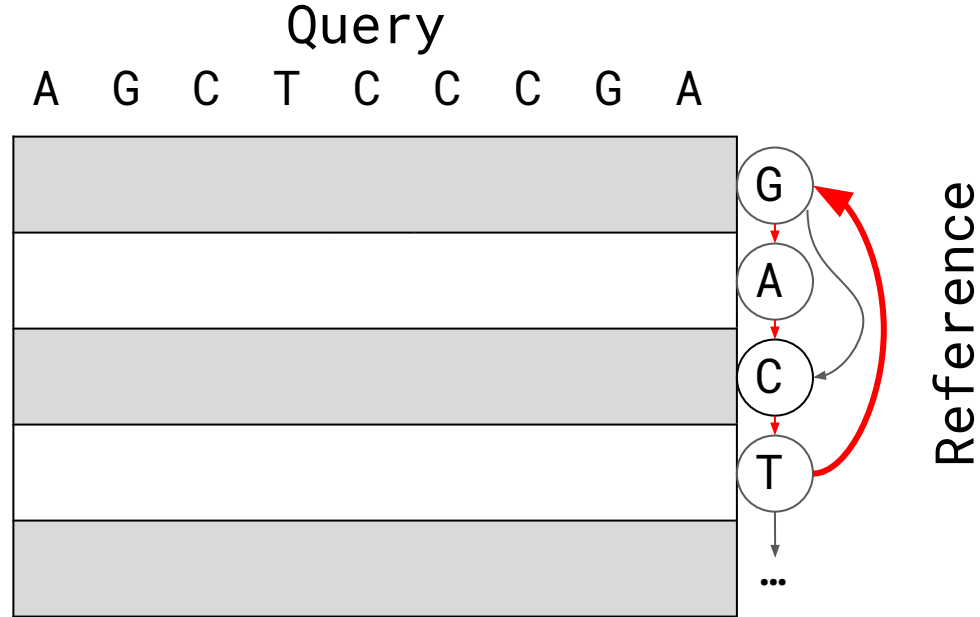
# Extension: Myers Bitvector

---

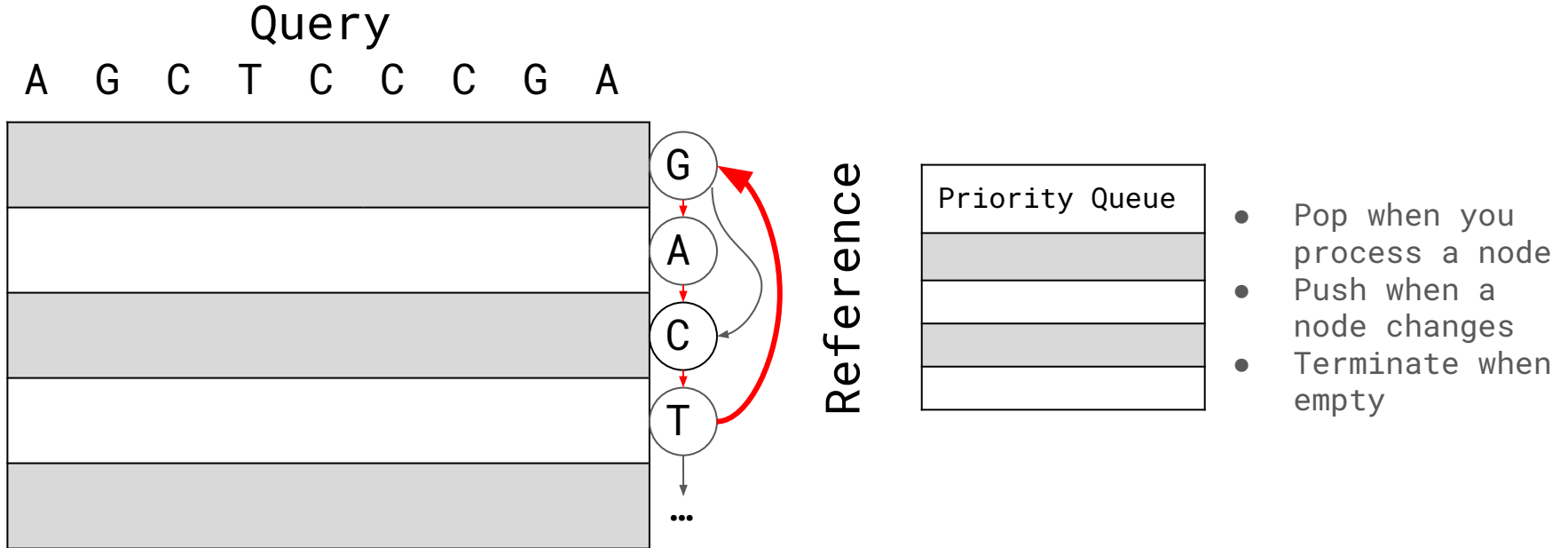


# Extension: Myers Bitvector

---



# Extension: Myers Bitvector

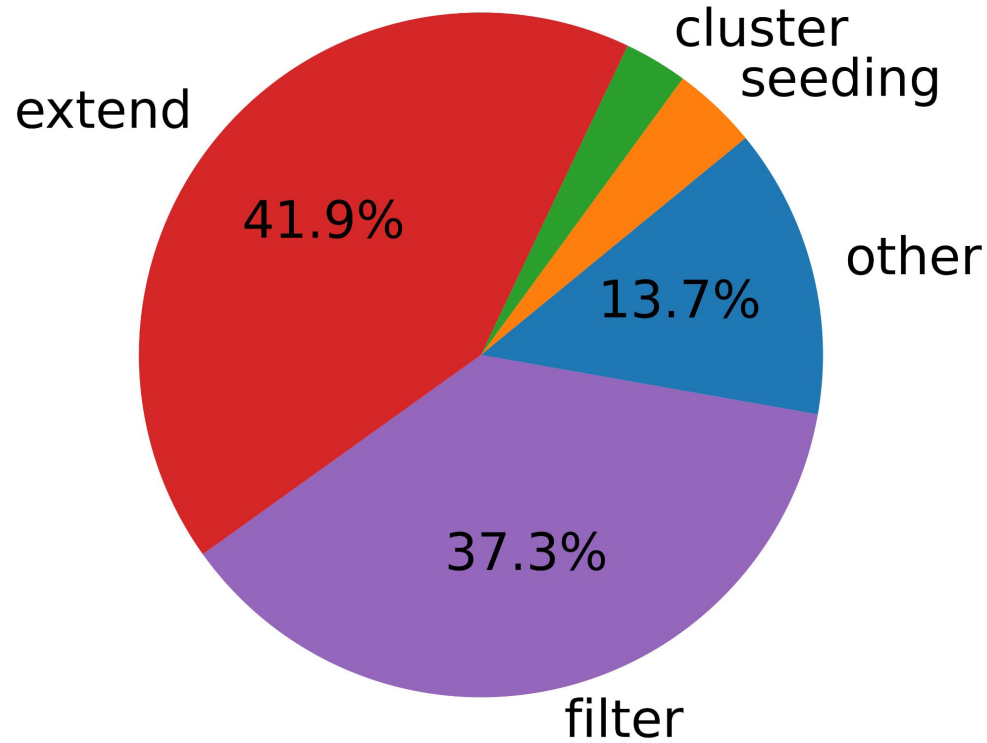


# Outline

- Workflows in Pangenomics
- Common Tools
  - Vg Map (Short Reads)
  - GraphAligner (Long Reads)
- **Profiling Results**
- Conclusion

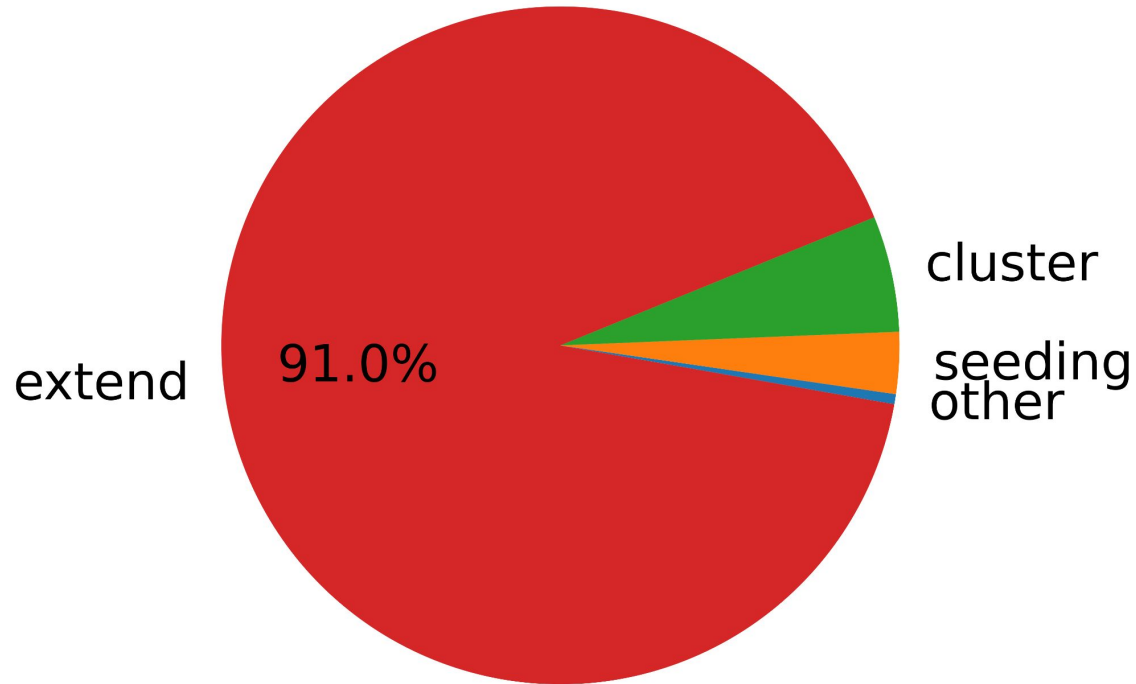
# Vg Map Timing Breakdown

---



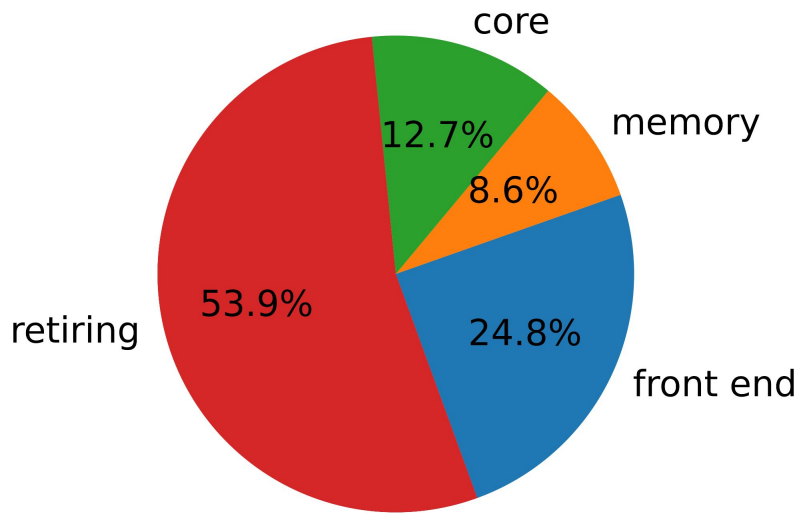
# GraphAligner Timing Breakdown

---

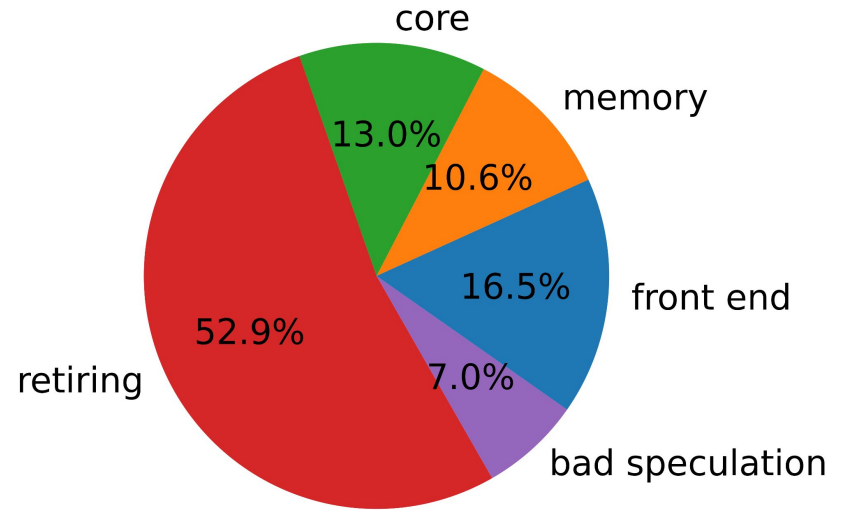


# Microarchitecture Utilization

---



Vg Map



GraphAligner



# Outline

- Workflows in Pangenomics
- Common Tools
  - Vg Map (Short Reads)
  - GraphAligner (Long Reads)
- Profiling Results
- **Conclusion**

# Conclusion

---

- **Seed extension** is the bottleneck for short and long read alignment
- **Cluster filtration** is also a bottleneck for Vg
- Multiple resource constraints restrict application performance without a single bottleneck



# Understanding and Characterization of Pangenomics

Noah Kaplan, Yufeng Gu, Reetuparna Das