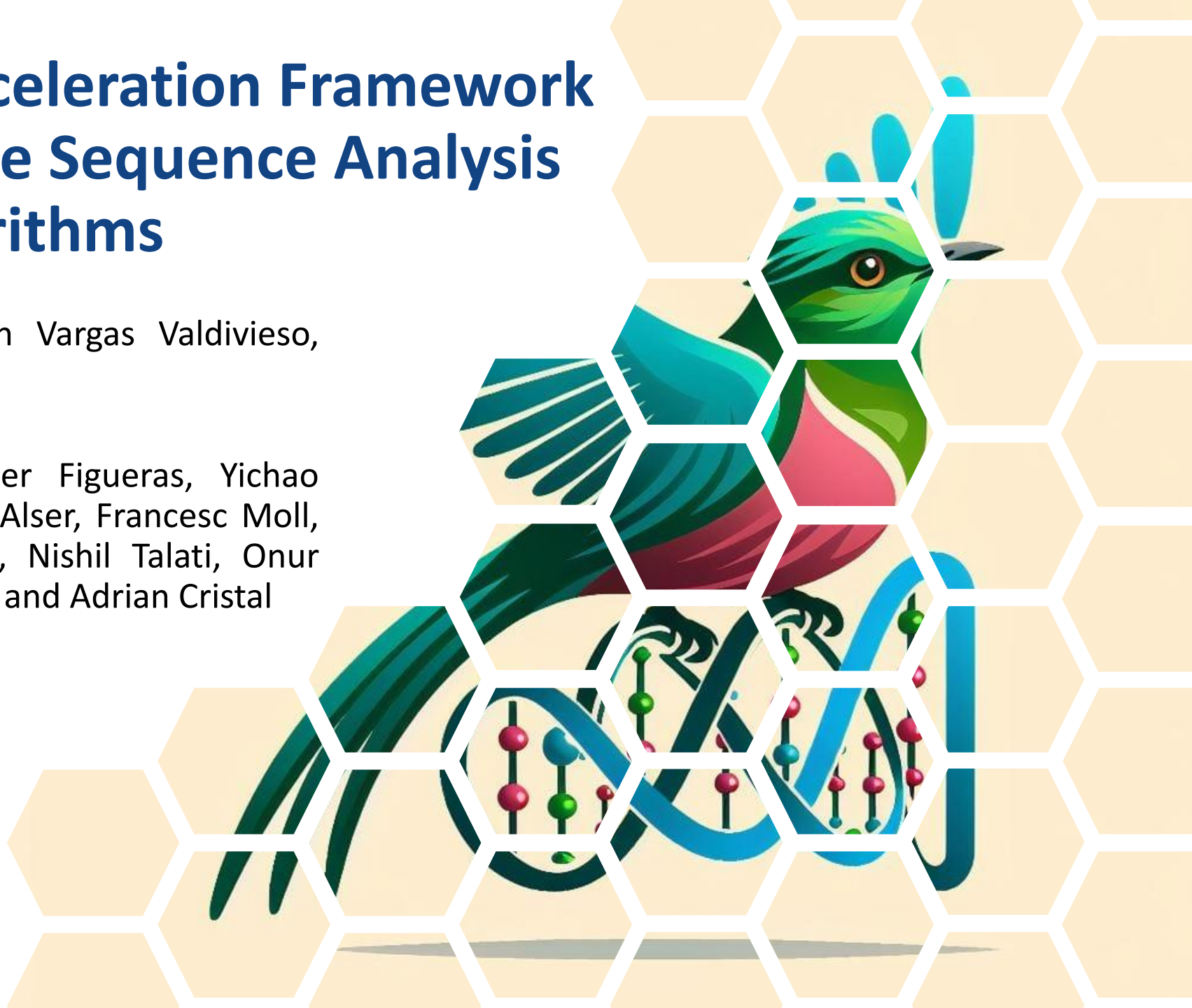


QUETZAL: Vector Acceleration Framework for Modern Genome Sequence Analysis Algorithms

Lead authors: Julian Pavon, Ivan Vargas Valdivieso,
Carlos Rojas, Cesar Hernandez

Co-authors: Mehmet Aslan, Roger Figueras, Yichao
Yuan, Joel Lindegger, Mohammed Alser, Francesc Moll,
Santiago Marco Sola, Oguz Ergin, Nishil Talati, Onur
Mutlu, Osman Unsal, Mateo Valero and Adrian Cristal

ETH zürich



Fun facts

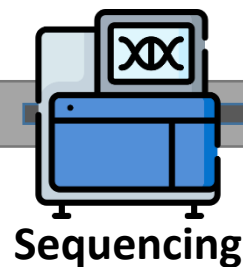
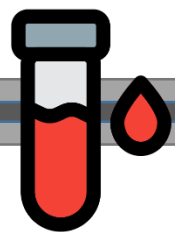
QUETZAL Team:



In the initial QUETZAL team we were all Mexicans.

In ancient cultures in Mexico, QUETZAL was a holy bird.

Genome Analysis Pipeline



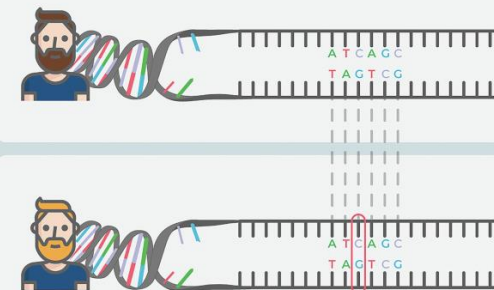
Sequencing

Sequence
Q4LMM4_9BURK
Q8XA99_ECO57
Q7CPQ5_SALTY
Q9I6D8_PSEAE
Q88JC4_PSEPK
Q8P751_XANCP
Q73CX9_BACCI

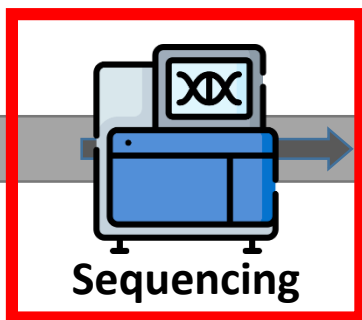
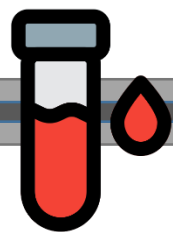
-----MQQRGGSMKKIAVLAVDEFEDSE
-----MSGKLDHCKVAILAVDGFEEAE
MGNSPFPQQRGGSMKKIAVLITDEFEDSE
-----MSKKIAVLITDEFEDSE
-----MTQSLHGKVVAAVLVDGFEQVE
-----MMSAQLNGKRVAPLVTDGFEQVE
-----MTHSLSGKTVAVLATSGFEQSE
-----MSKKIATLITDYFEDTE

Read Mapping

Genomic Variant



Genome Analysis Pipeline

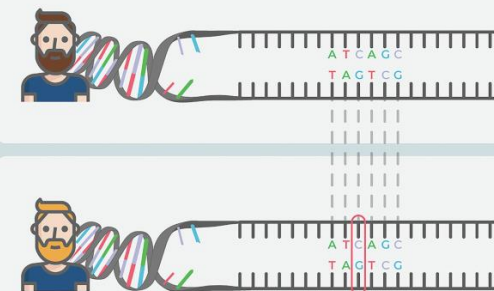


Sequence
Q4LMM4_9BURK
Q8XA99_ECO57
Q7CPQ5_SALTY
Q9I6D8_PSEAE
Q88JC4_PSEPK
Q8P751_XANCP
Q73CX9_BACCI

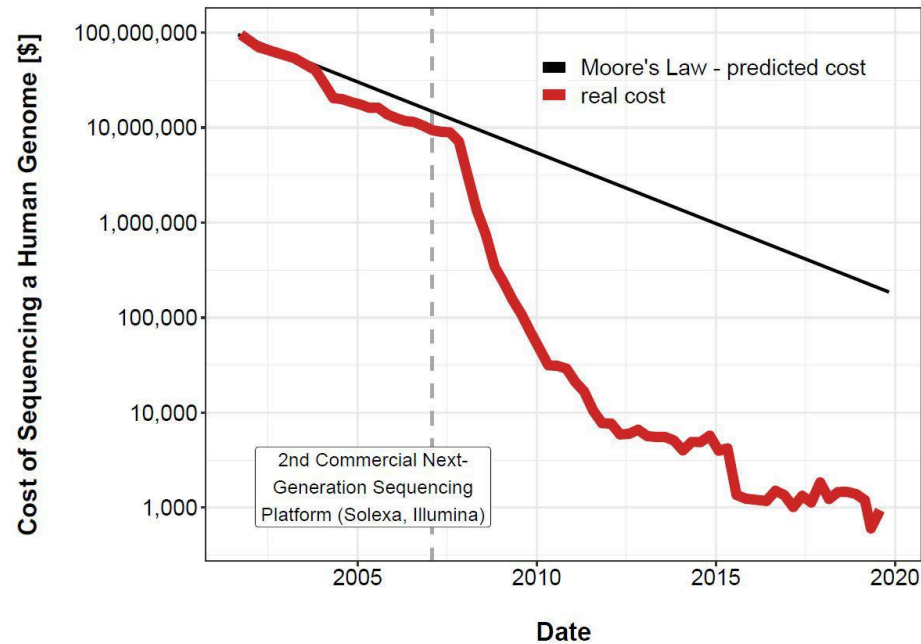
-----MQQGGGMSKKI AVLAVDEFEDSE
-----MSGKLDHCKVAILAVDGFEEAE
MGNSPHMQQGGGMSKKI AVLITDEFEDSE
-----MSKKI AVLITDEFEDSE
-----MTQSLHGKVVAAVLVDGFEQVE
-----MMSAQLNGKRVAPLVTDGFEQVE
-----MTHSLSGKTVAVLATSGFEQSE
-----MSKKIATLITDYFEDTE

Read Mapping

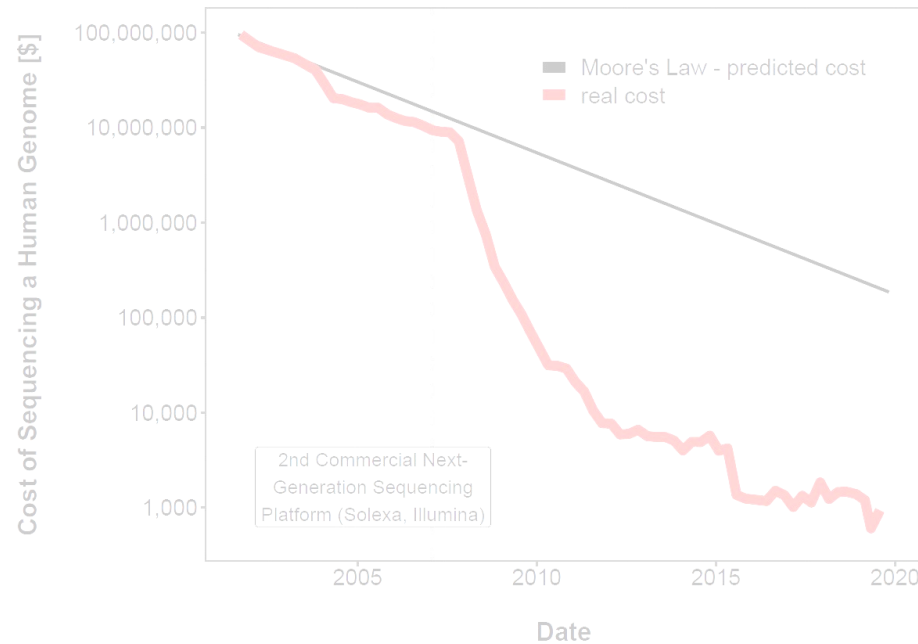
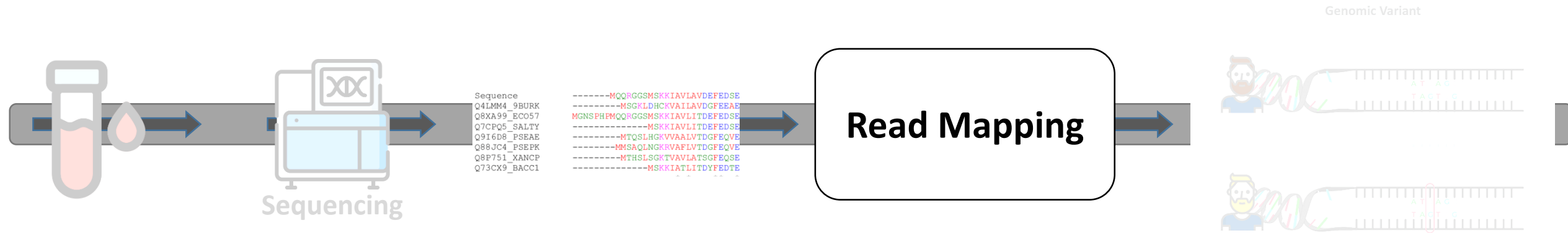
Genomic Variant



Genome Analysis Pipeline

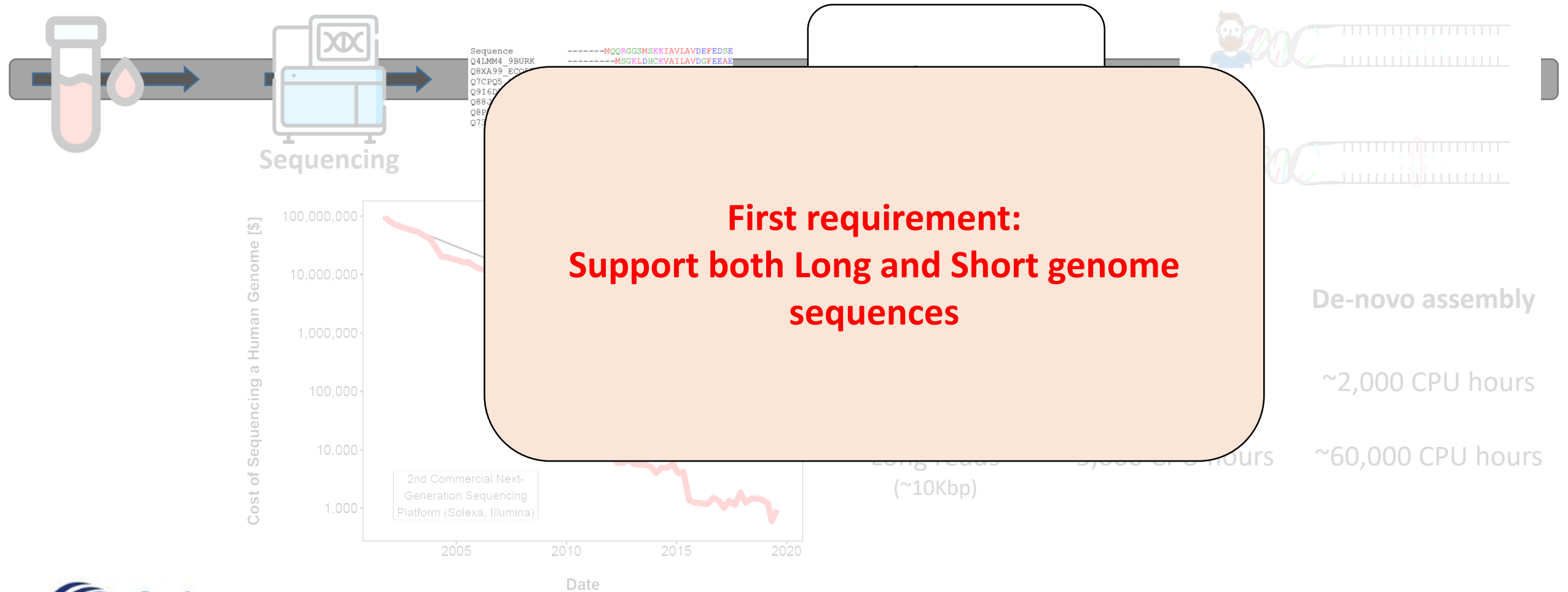


Genome Analysis Pipeline

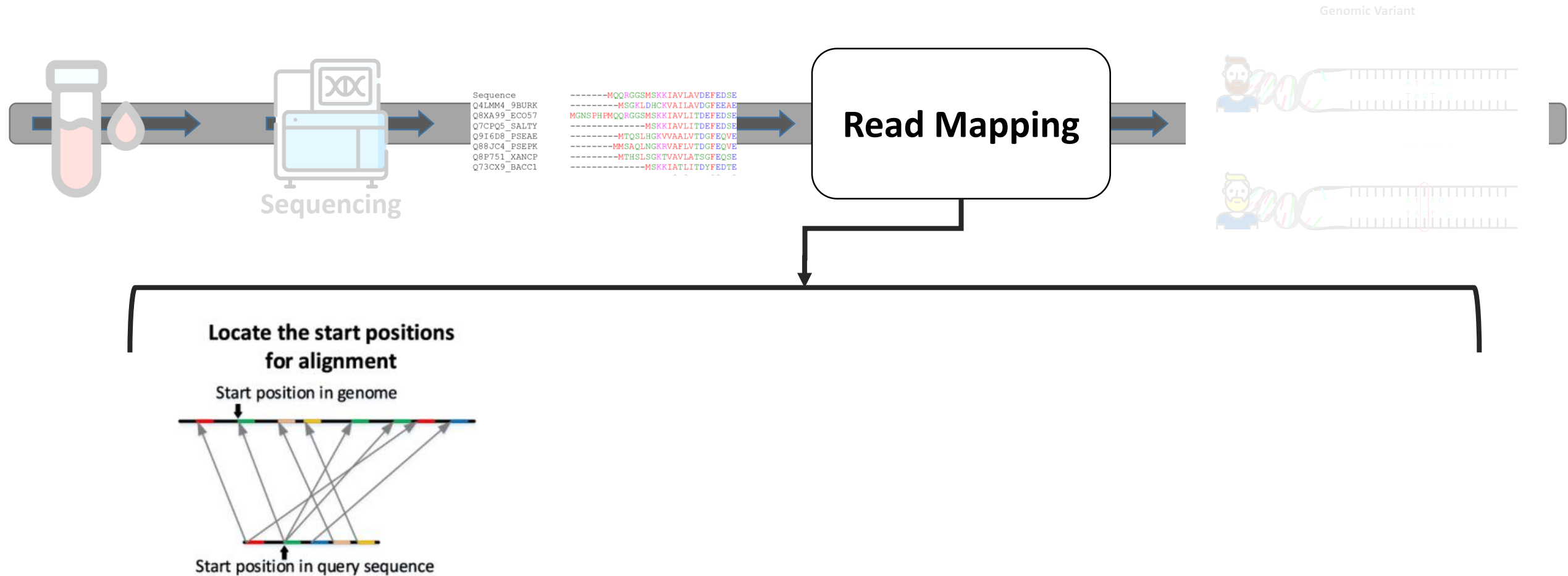


	Reference-guided assembly	De-novo assembly
Short reads (~100bp)	~200 CPU hours	~2,000 CPU hours
Long reads (~10Kbp)	~5,000 CPU hours	~60,000 CPU hours

Genome Analysis Pipeline

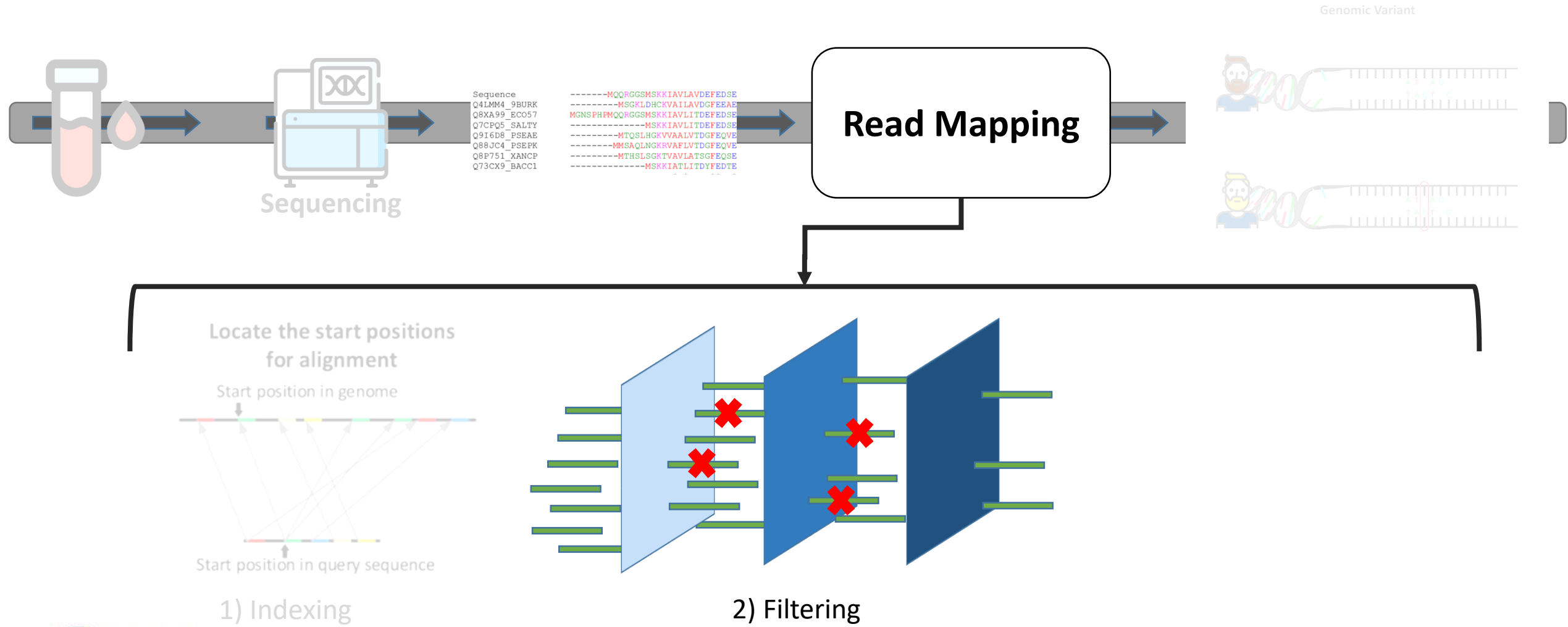


Genome Analysis Pipeline

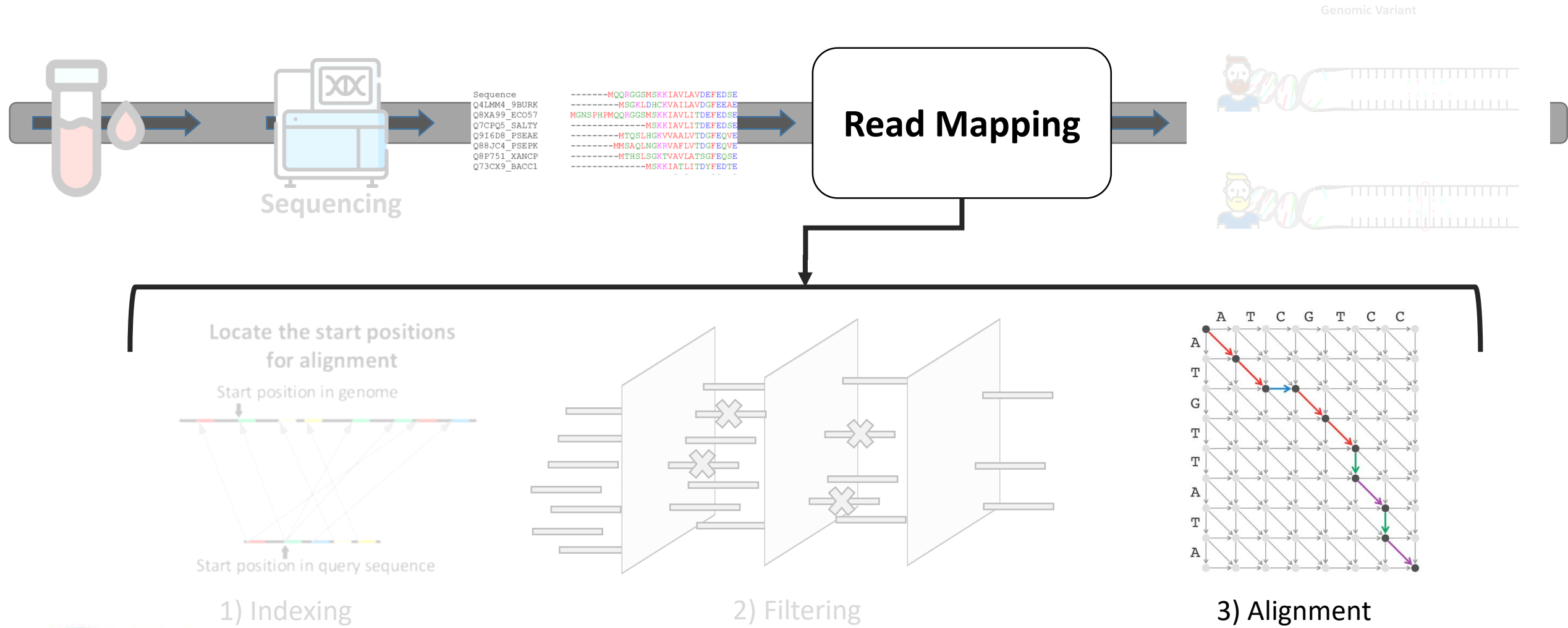


1) Indexing

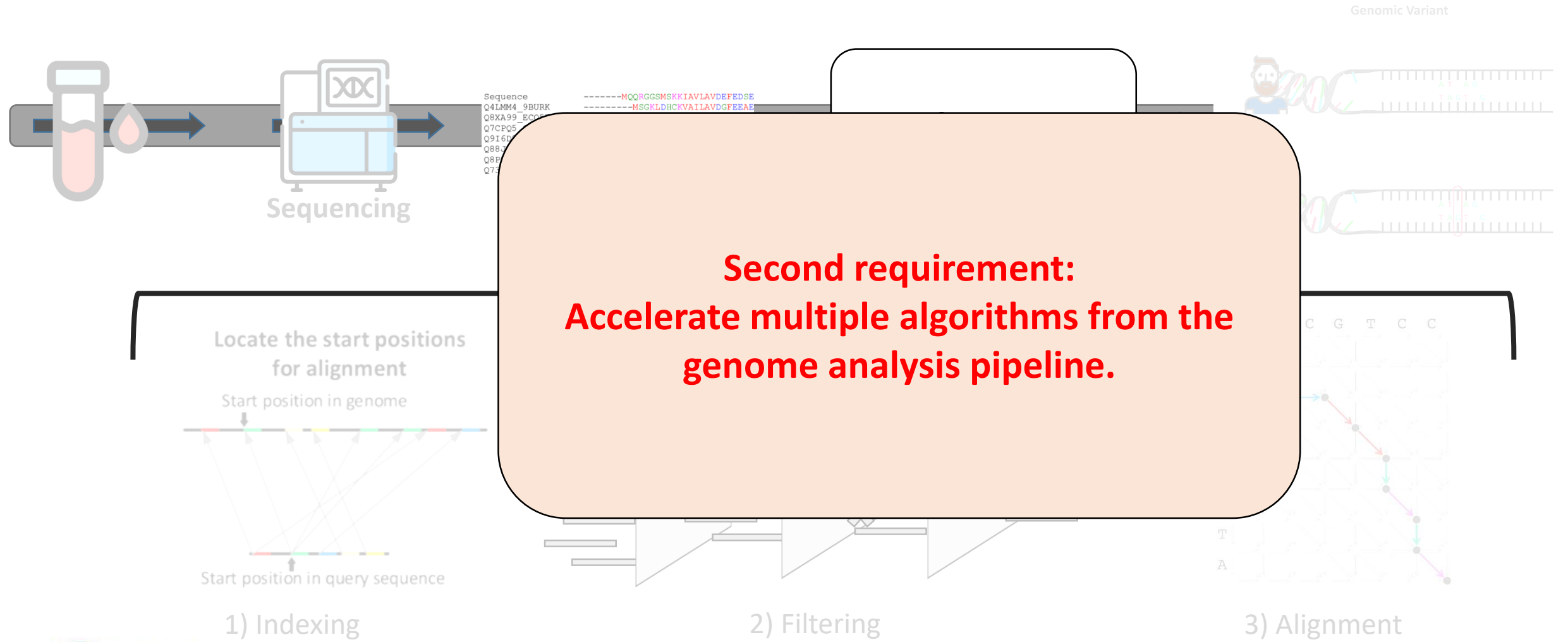
Genome Analysis Pipeline



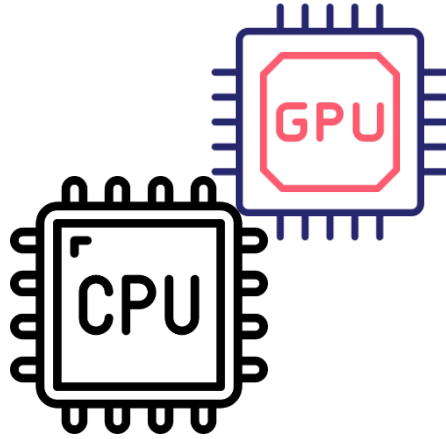
Genome Analysis Pipeline



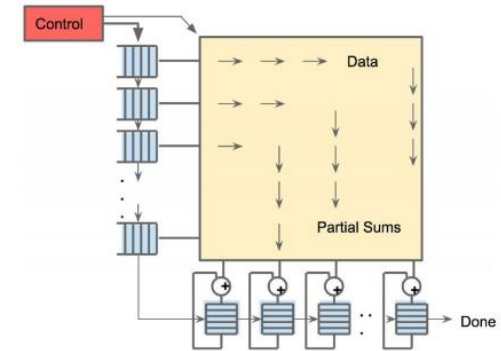
Genome Analysis Pipeline



Accelerating the Genome Analysis Pipeline

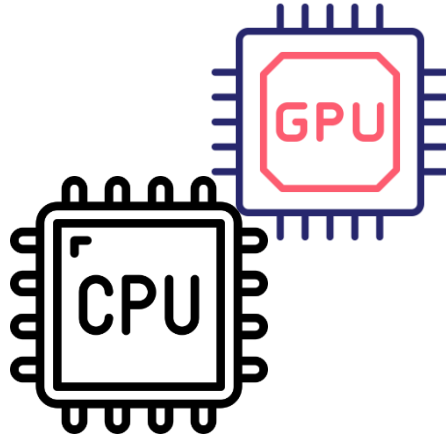


CPU/GPU



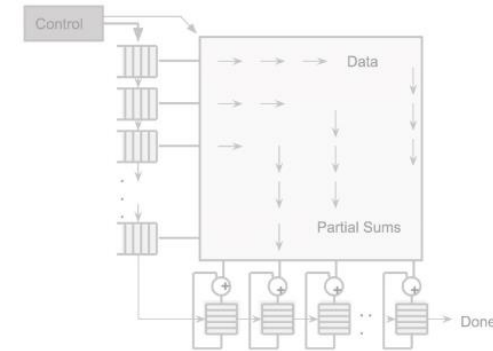
Custom Accelerator

Accelerating the Genome Analysis Pipeline



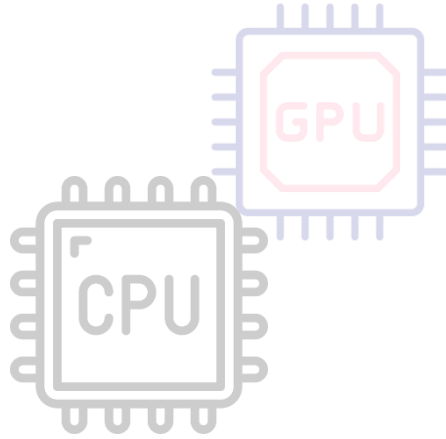
CPU/GPU

- Flexible.
- Low entry-cost.
- Generality and flexibility limits their performance.



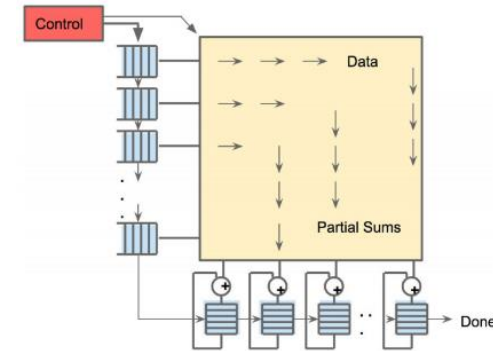
Custom Accelerator

Accelerating the Genome Analysis Pipeline



CPU/GPU

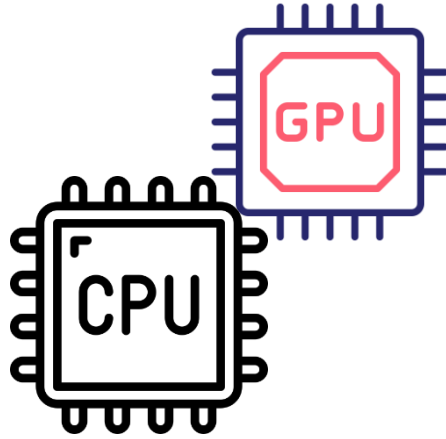
- Flexible.
- Low entry-cost.
- Generality and flexibility limits their performance.



Custom Accelerator

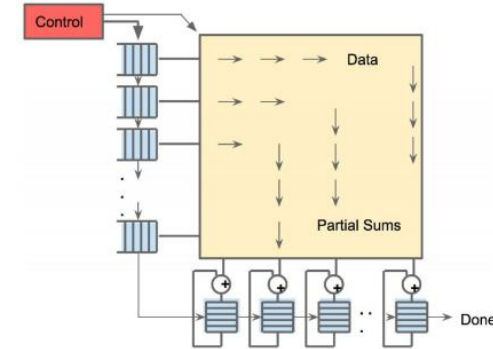
- High performance and Energy efficiency.
- They are tight to a single algorithm or a single sequence length.
- High design and entry-cost.

Accelerating the Genome Analysis Pipeline



CPU/GPU

- Flexible.
- Low entry-cost.
- Generality and flexibility limit their performance.

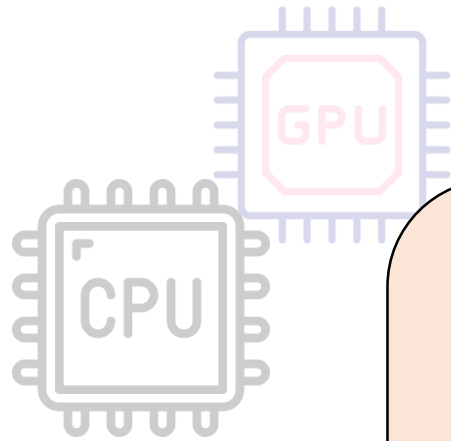


Custom Accelerator

- High performance and Energy efficiency.
- They are tight to a single algorithm or a single sequence length.
- High design and entry-cost.

??

Accelerating the Genome Analysis Pipeline



CPU/GPU

- Flexible.
- Low entry-cost.
- Generality and flexibility in performance.

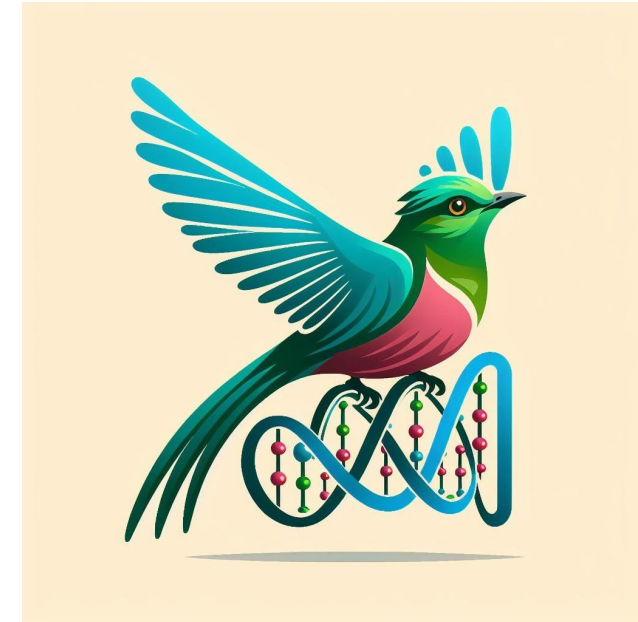
Answer:
**We proposed QUETZAL, a hardware-
software co-designed vector acceleration
framework**



Key Contributions

Key Contributions:

- QUETZAL addresses the bottlenecks from vector architectures when processing genome sequence analysis workloads.



Key Contributions

Key Contributions:

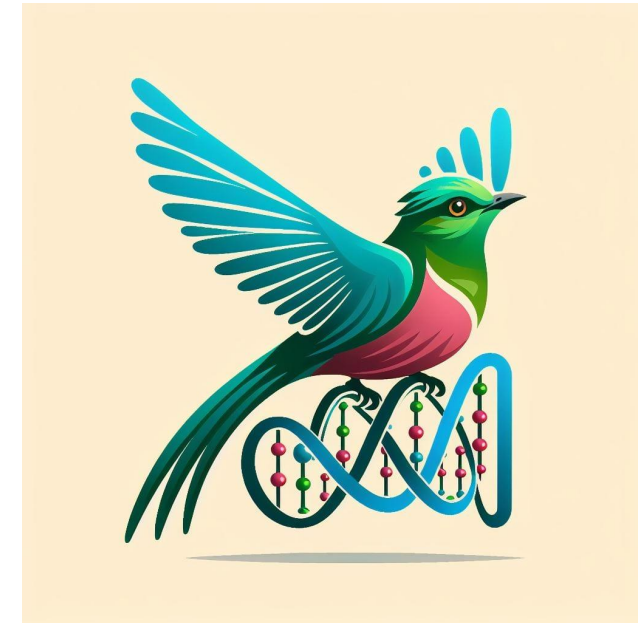
- QUETZAL addresses the bottlenecks from vector architectures when processing genome sequence analysis workloads.
- A CPU featuring QUETZAL, achieves 5.7x better performance compared to the baseline CPU architecture.



Key Contributions

Key Contributions:

- QUETZAL addresses the bottlenecks from vector architectures when processing genome sequence analysis workloads.
- A CPU featuring QUETZAL, achieves 5.7x better performance compared to the baseline CPU architecture.
- QUETZAL outperforms GPUs when processing long genome sequences.



Outline

1. Introduction and motivation
2. Background
 - a) Modern Genome Analysis Algorithms
 - b) Bottlenecks
3. QUETZAL
 - a) Insights and Functionality
4. Evaluation
5. Conclusion

Outline

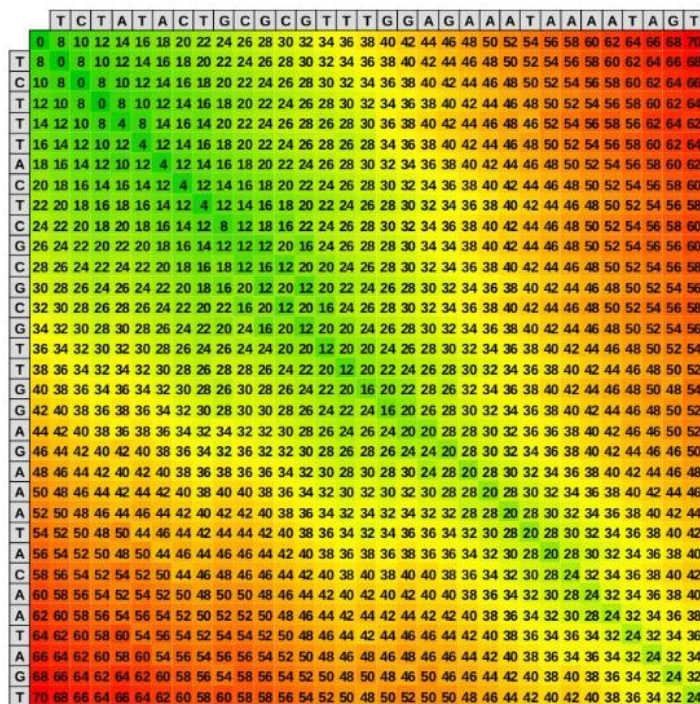
1. Introduction and motivation
- 2. Background**
 - a) Modern Genome Analysis Algorithms
 - b) Bottlenecks
3. QUETZAL
 - a) Insights and Functionality
4. Evaluation
5. Conclusion

Modern Genome Analysis Algorithms

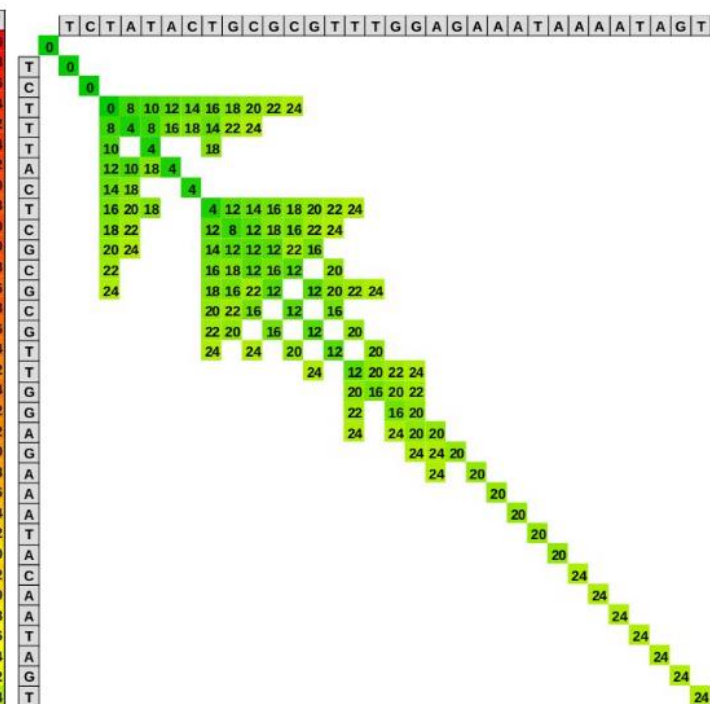
	T	C	T	A	T	A	C	T	G	C	G	C	G	T	T	G	G	A	G	A	A	A	T	A	A	A	T	A	G	T			
	0	8	10	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62	64	66	68	70
T	8	0	8	10	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62	64	66	68
C	10	8	0	8	10	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62	64	66
T	12	10	8	0	8	10	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62	64
T	14	12	10	8	0	8	10	12	14	16	18	20	22	24	26	28	30	36	38	40	42	44	46	48	50	52	54	56	58	62	64	66	68
T	16	14	12	10	12	4	12	14	16	18	20	22	24	26	28	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62	64	66	
A	18	16	14	12	10	12	4	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62
C	20	18	16	14	16	14	12	4	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60
T	22	20	18	16	18	16	14	12	4	12	14	16	18	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58
C	24	22	20	18	16	18	14	12	8	12	18	16	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60	62
G	26	24	22	20	22	20	18	16	14	12	12	18	20	16	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	60
C	28	26	24	22	24	22	20	18	16	18	16	12	16	12	20	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58
G	30	28	26	24	26	24	22	20	18	16	12	10	12	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58
C	32	30	28	26	28	26	24	22	20	22	16	12	20	16	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	
G	34	32	30	28	30	28	26	24	22	20	24	16	12	20	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	
T	36	34	32	30	32	30	28	26	24	26	24	20	12	20	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	
T	38	36	34	32	34	32	30	28	26	28	26	24	22	12	20	22	24	26	28	30	32	34	36	38	40	42	44	46	48	50	52	54	
G	40	38	36	34	36	34	32	30	28	30	28	26	24	22	16	20	22	26	32	34	36	38	40	42	44	46	48	50	52	54	56	58	
G	42	40	38	36	38	36	34	32	30	28	30	28	26	24	22	16	20	26	30	32	34	36	38	40	42	44	46	48	50	52	54	56	
A	44	42	40	38	36	38	36	34	32	34	32	30	28	26	24	20	28	30	32	36	38	40	42	44	46	48	50	52	54	56	58	60	
G	46	44	42	40	42	40	38	36	34	32	32	30	28	26	24	20	28	30	32	34	36	38	40	42	44	46	48	50	52	54	56	58	
A	48	46	44	42	40	42	40	38	36	38	36	34	32	30	28	30	32	34	28	20	30	32	34	36	38	40	42	44	46	48	50	52	
A	50	48	46	44	42	44	42	40	38	40	38	36	34	32	30	32	30	28	28	20	28	30	32	34	36	38	40	42	44	46	48	50	
A	52	50	48	46	44	46	44	42	40	42	40	38	36	34	32	34	32	32	28	28	20	28	30	32	34	36	38	40	42	44	46	48	
T	54	52	50	48	50	44	46	44	42	44	44	42	40	38	36	34	32	34	36	36	34	32	30	28	20	28	30	32	34	36	38	40	
A	56	54	52	50	50	44	46	44	46	46	44	42	40	38	36	38	36	36	34	32	30	28	20	28	30	32	34	36	38	40	42	44	
C	58	56	54	52	54	50	44	46	46	46	44	42	40	38	40	38	40	40	38	36	34	32	30	28	24	32	34	36	38	40	42	44	
A	60	58	56	54	52	54	50	48	50	46	44	42	40	42	40	40	40	38	36	34	32	30	28	24	32	34	36	38	40	42	44	46	
A	62	60	58	56	54	56	52	50	52	50	48	46	44	42	44	42	44	42	40	38	36	34	32	30	28	24	32	34	36	38	40	42	
T	64	62	60	58	56	54	56	54	52	54	52	50	48	46	44	42	44	46	44	42	40	38	36	34	36	34	32	34	36	38	40	42	
A	66	64	62	60	58	60	54	56	54	56	56	54	52	50	48	46	46	48	46	44	42	40	38	36	34	36	34	32	34	36	38	40	
T	68	66	64	62	62	60	58	56	54	58	56	54	52	50	48	50	48	46	50	46	44	42	40	38	40	38	36	34	32	34	36	38	
T	70	68	66	64	66	64	62	60	58	60	58	56	54	52	50	48	50	52	50	48	46	44	42	40	42	40	38	36	34	32	34	36	

Classical approach

Modern Genome Analysis Algorithms



Classical approach



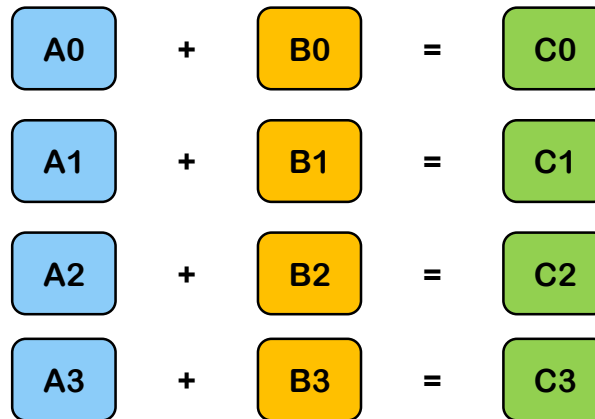
Modern approach (WFA[1])

Vector Architectures

- Vector architectures follow the **SIMD** (Single Instruction Multiple Data) taxonomy.
- Really efficient to exploit **Data-level Parallelism**.

Vector Architectures

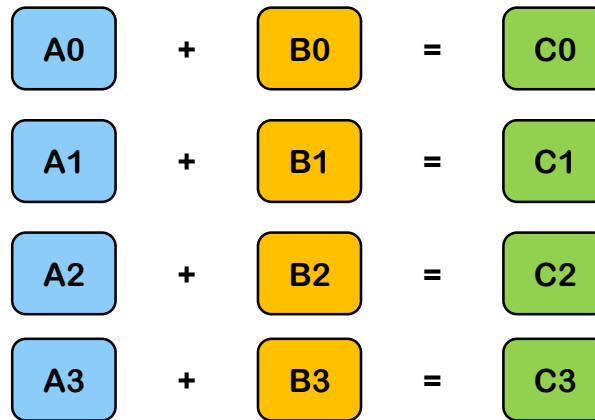
- Vector architectures follow the **SIMD** (Single Instruction Multiple Data) taxonomy.
- Really efficient to exploit **Data-level Parallelism**.



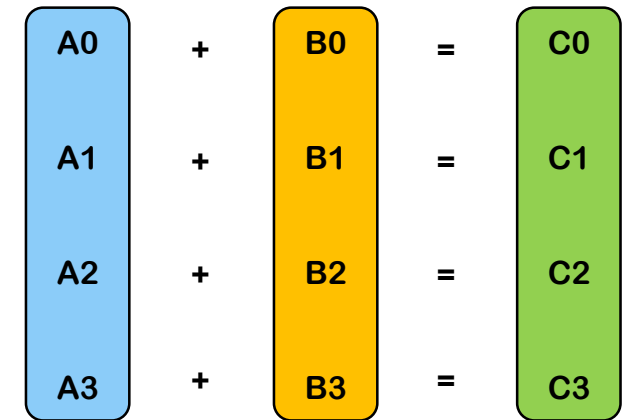
Scalar Operation

Vector Architectures

- Vector architectures follow the **SIMD** (Single Instruction Multiple Data) taxonomy.
- Really efficient to exploit **Data-level Parallelism**.

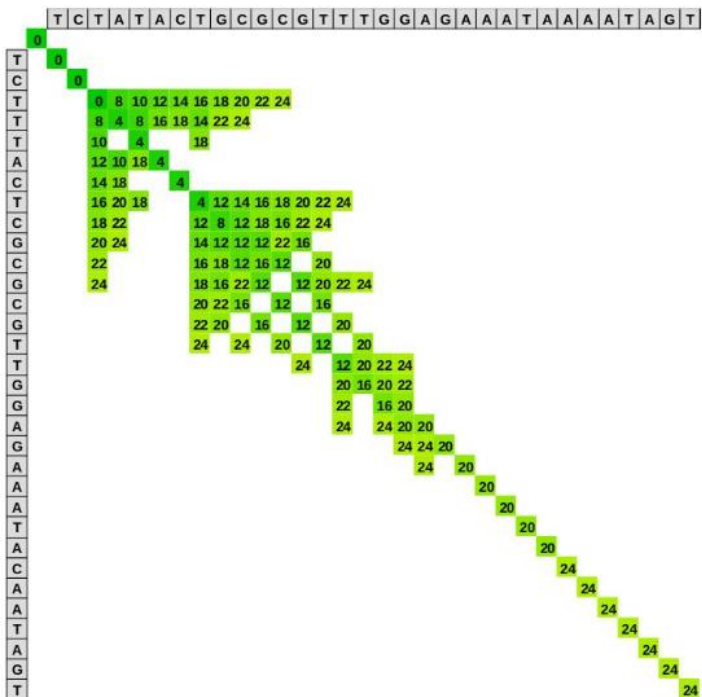


Scalar Operation



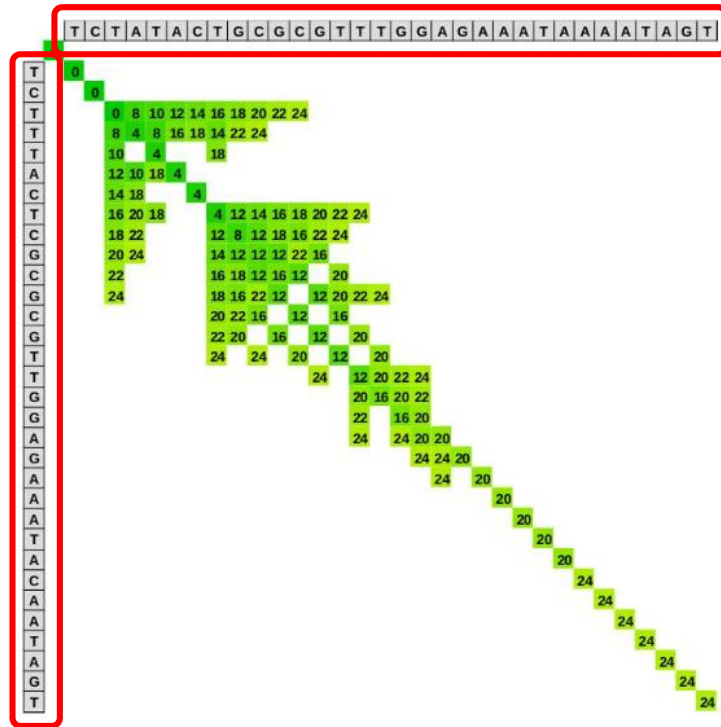
Vector Operation

Bottlenecks for Vector Architectures



Modern approach (WFA[1])

Bottlenecks for Vector Architectures



Modern approach (WFA[1])

Bottlenecks for Vector Architectures

Memory

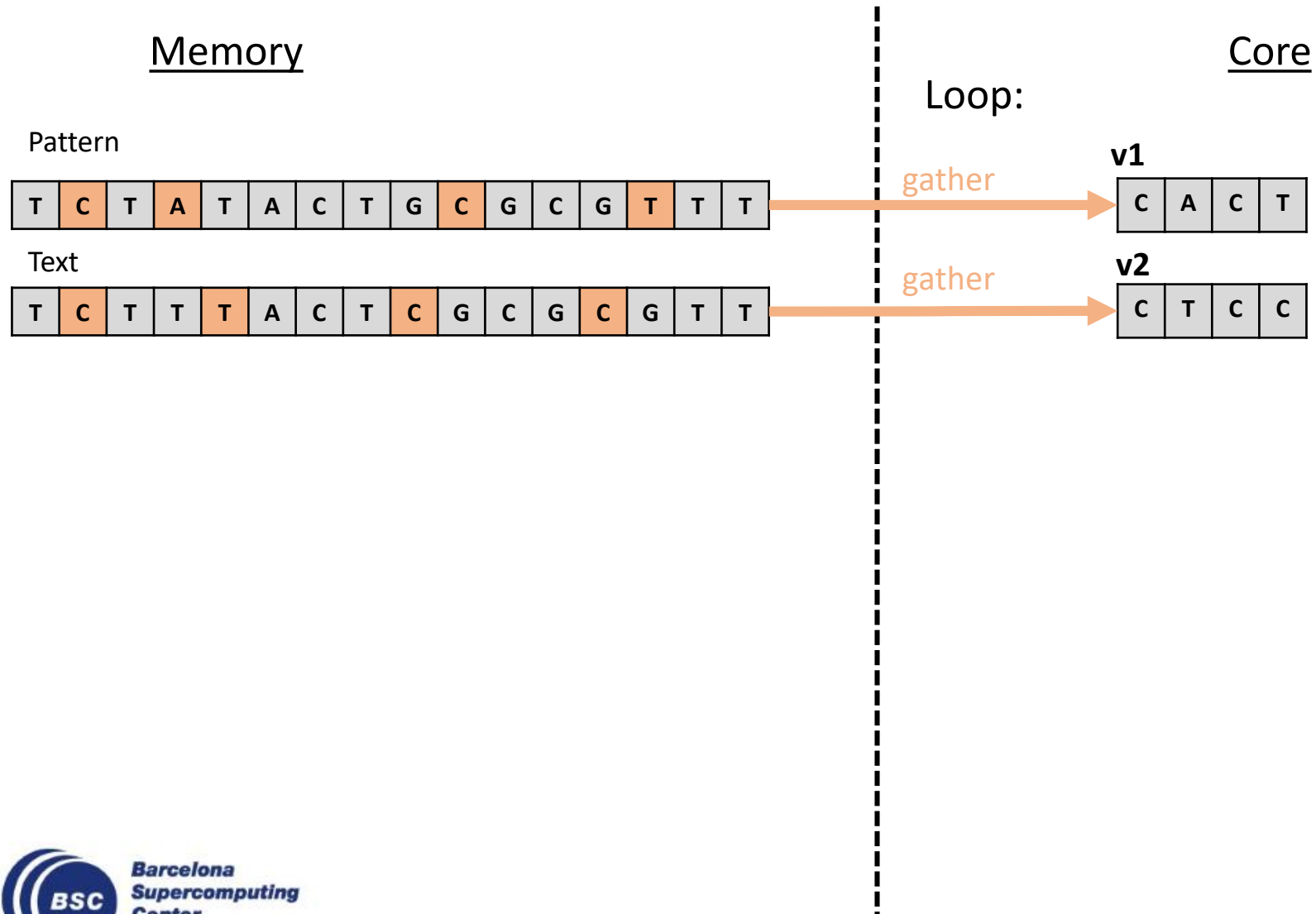
Pattern

T	C	T	A	T	A	C	T	G	C	G	C	G	T	T	T
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

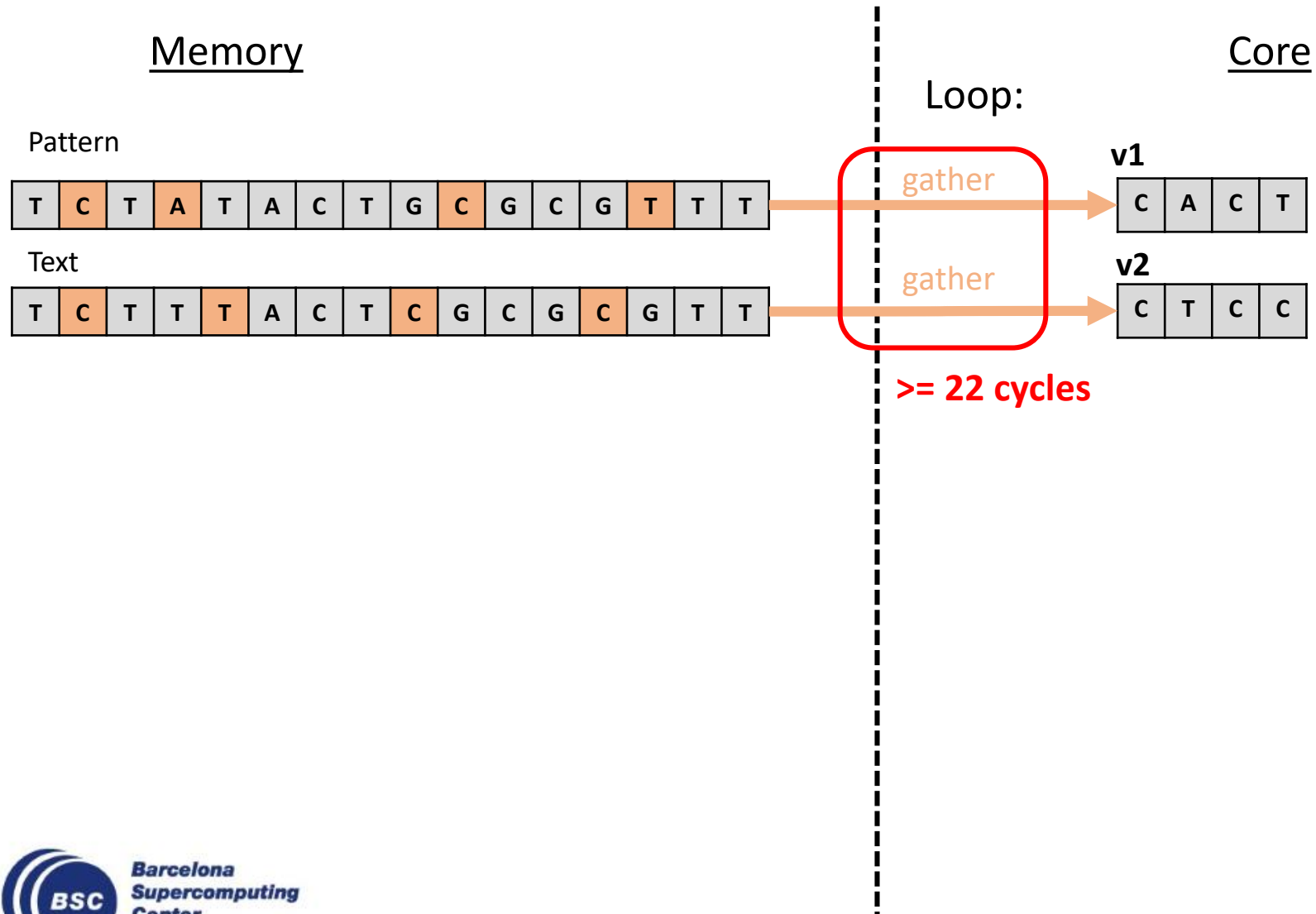
Text

T	C	T	T	T	A	C	T	C	G	C	G	C	G	T	T
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

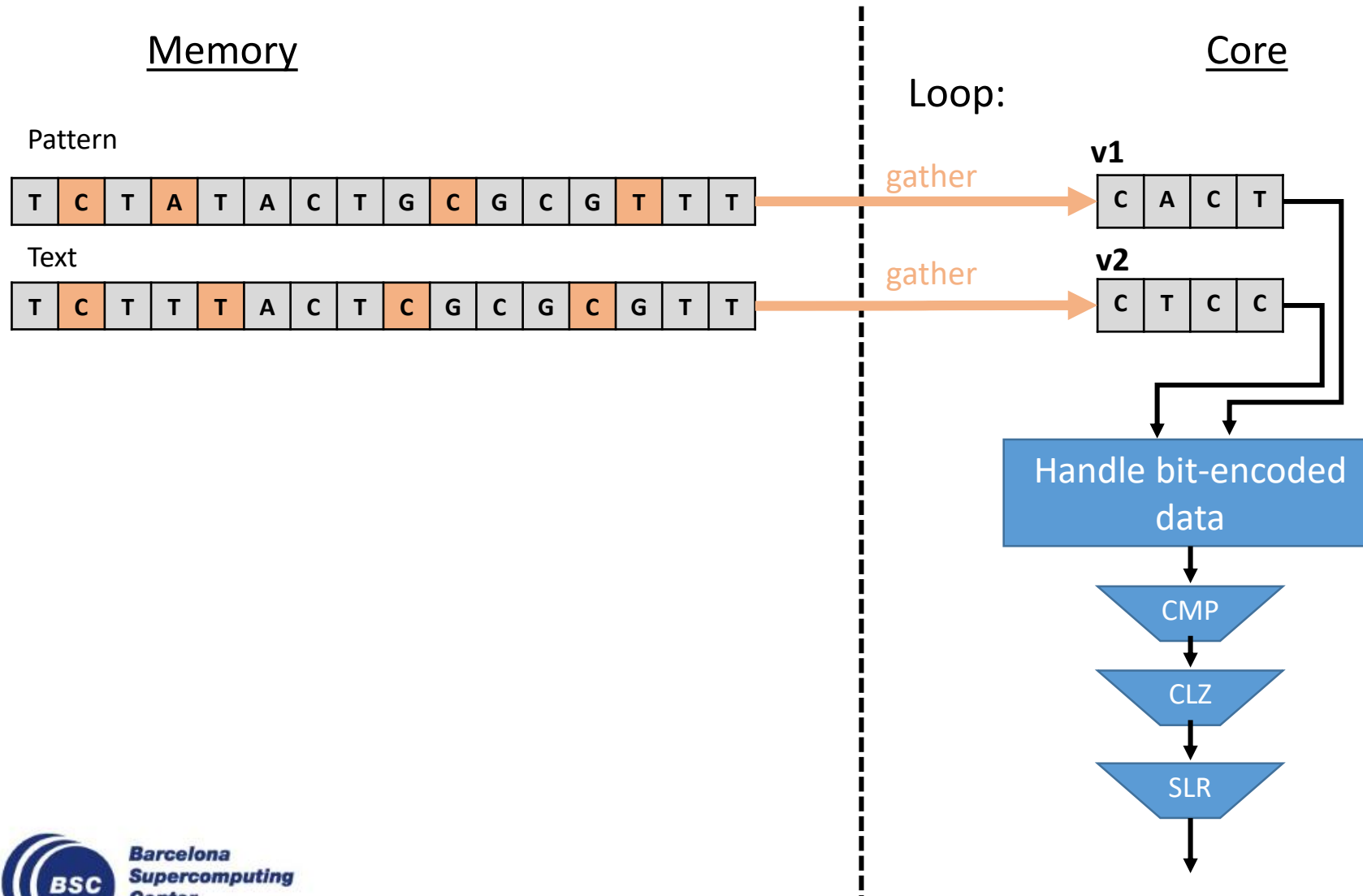
Bottlenecks for Vector Architectures



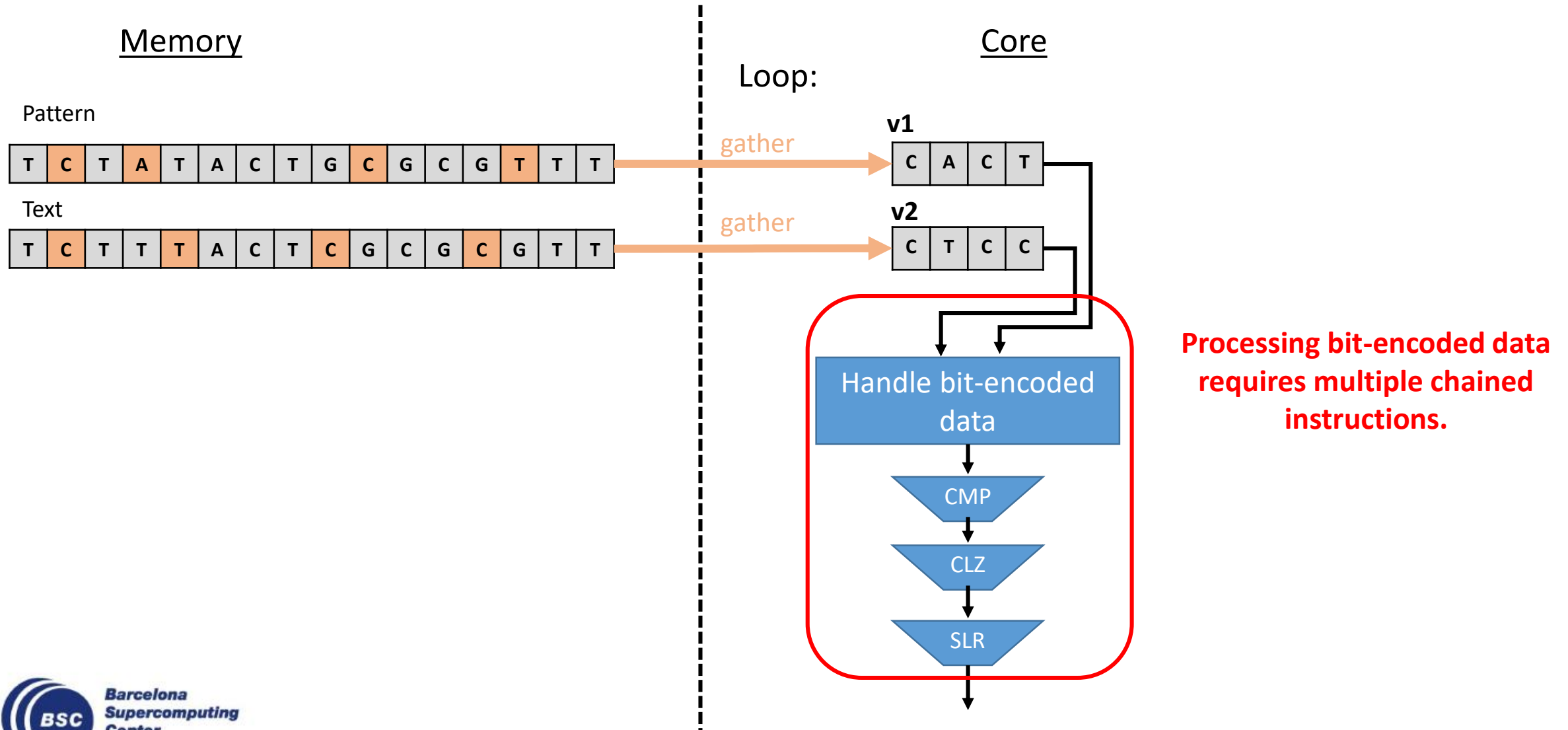
Bottlenecks for Vector Architectures



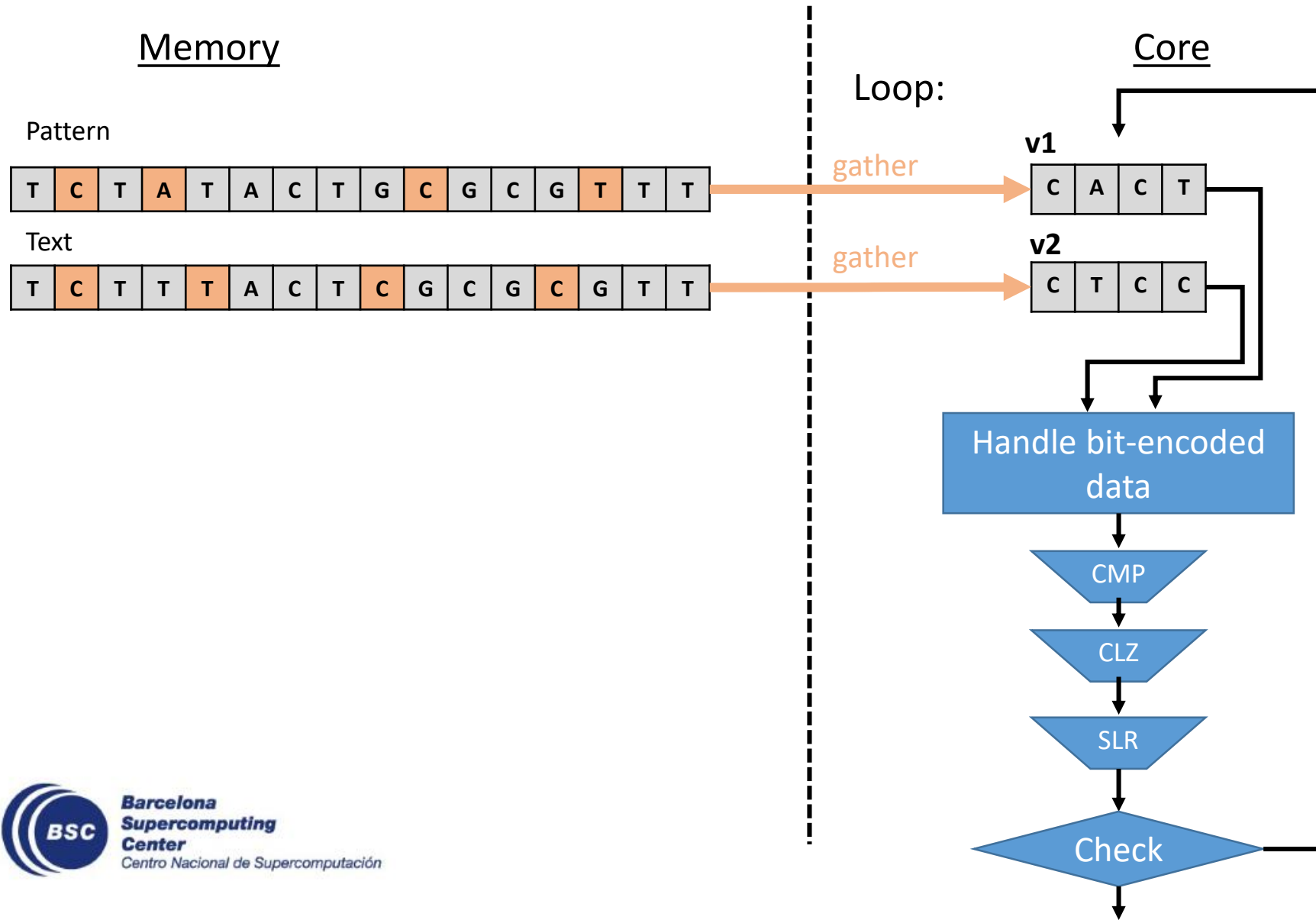
Bottlenecks for Vector Architectures



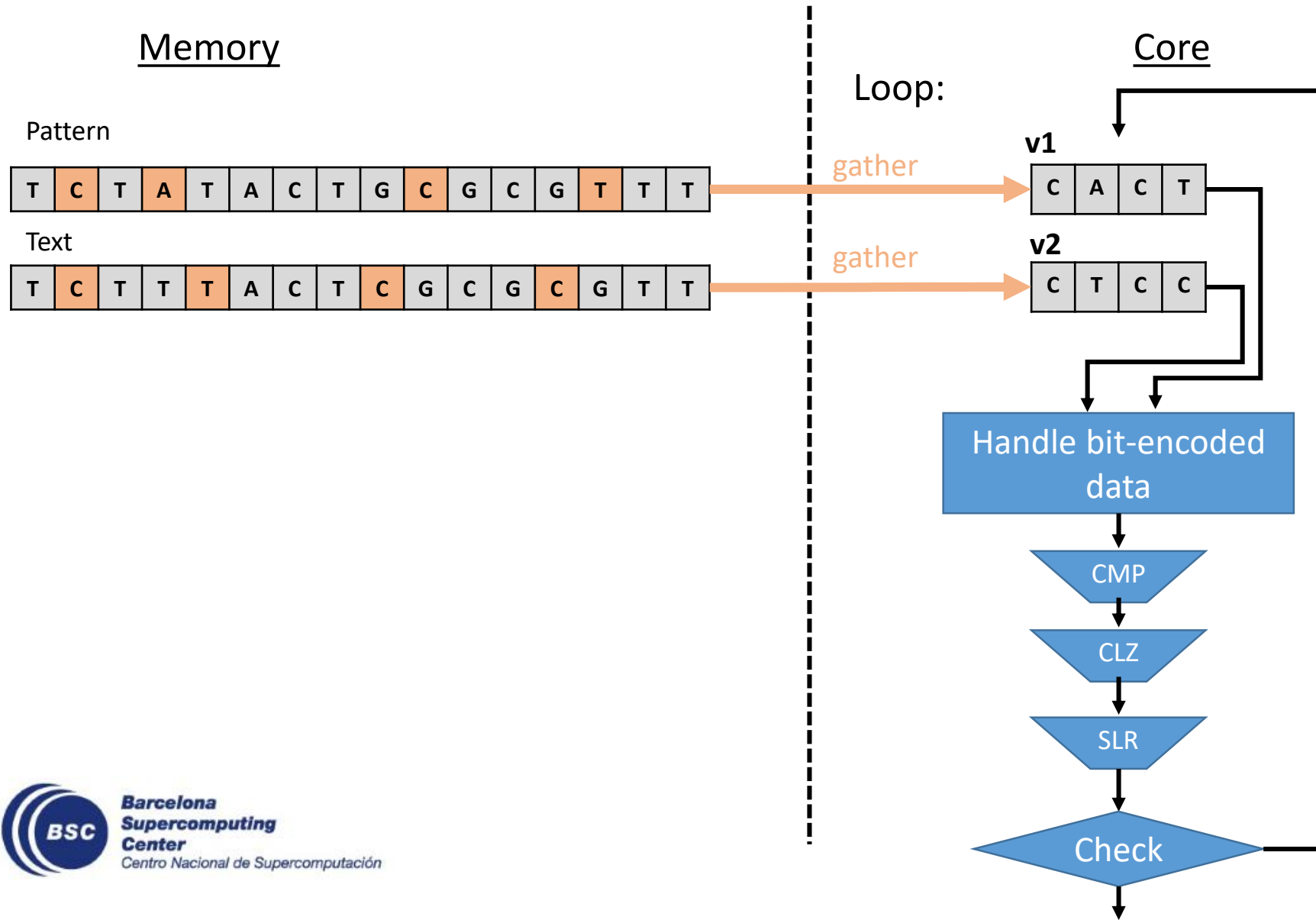
Bottlenecks for Vector Architectures



Bottlenecks for Vector Architectures

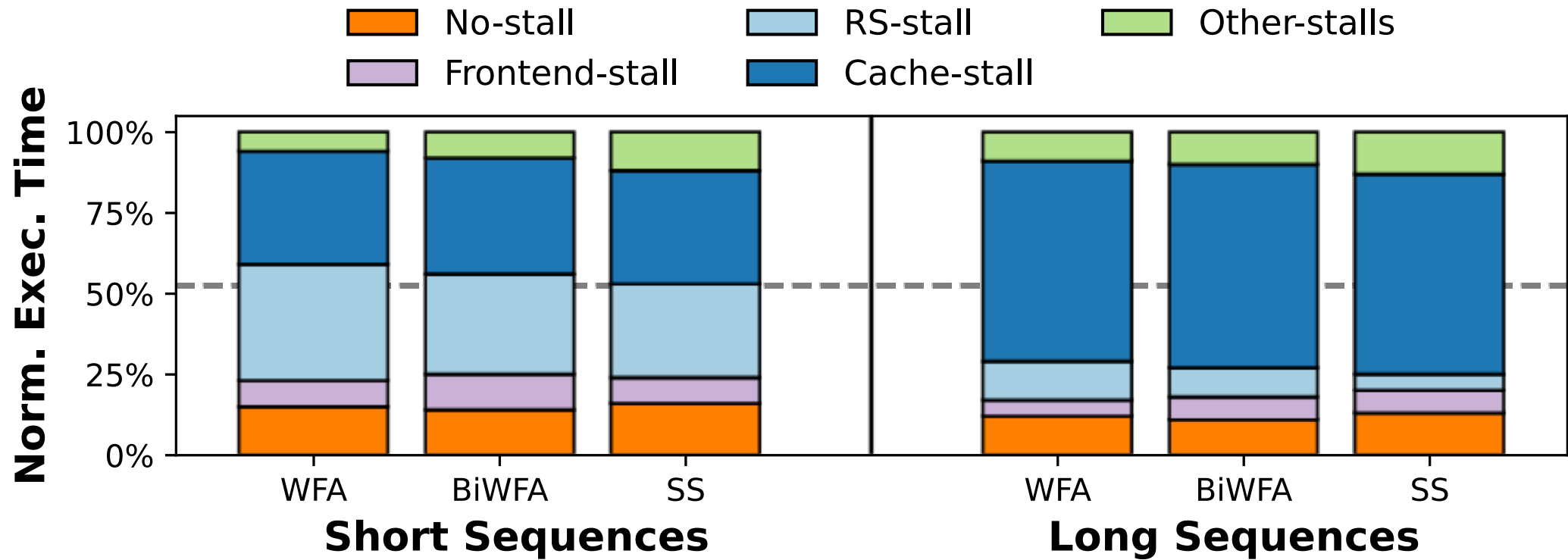


Bottlenecks for Vector Architectures

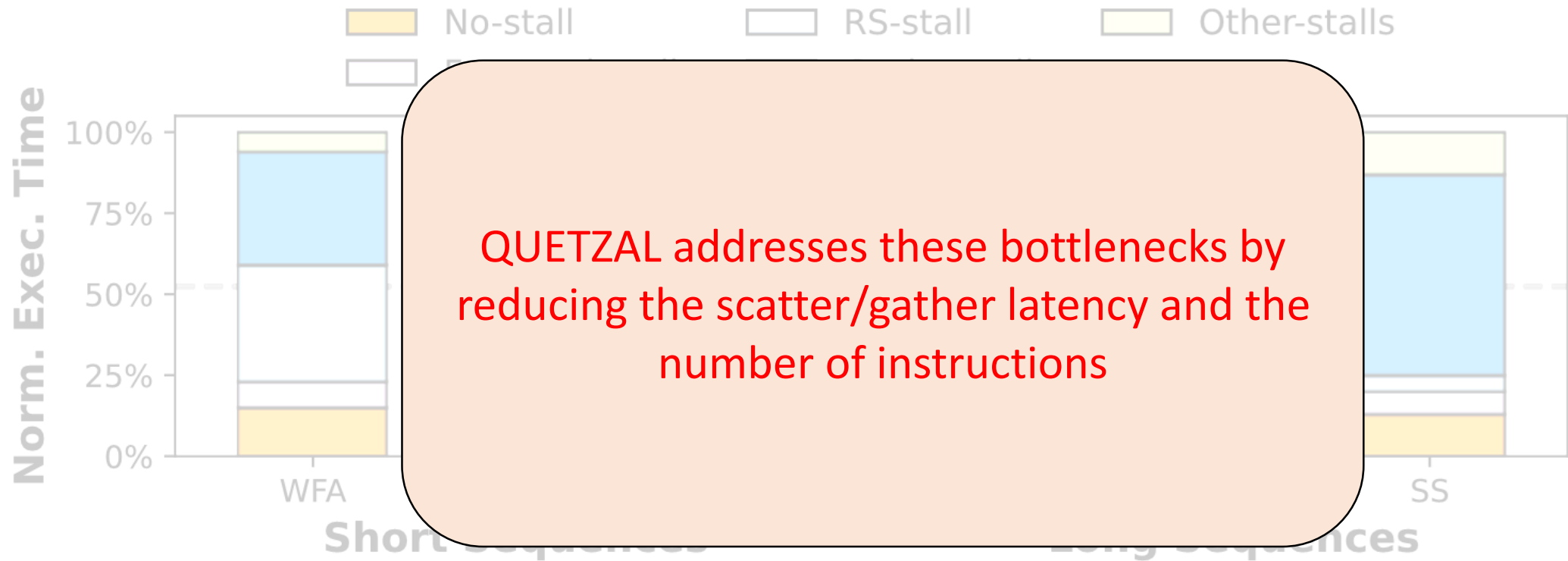


Two major bottlenecks:
1) Scatter/gather instructions
2) The large number of instructions to handle bit-encoded data.

Bottlenecks for Vector Architectures



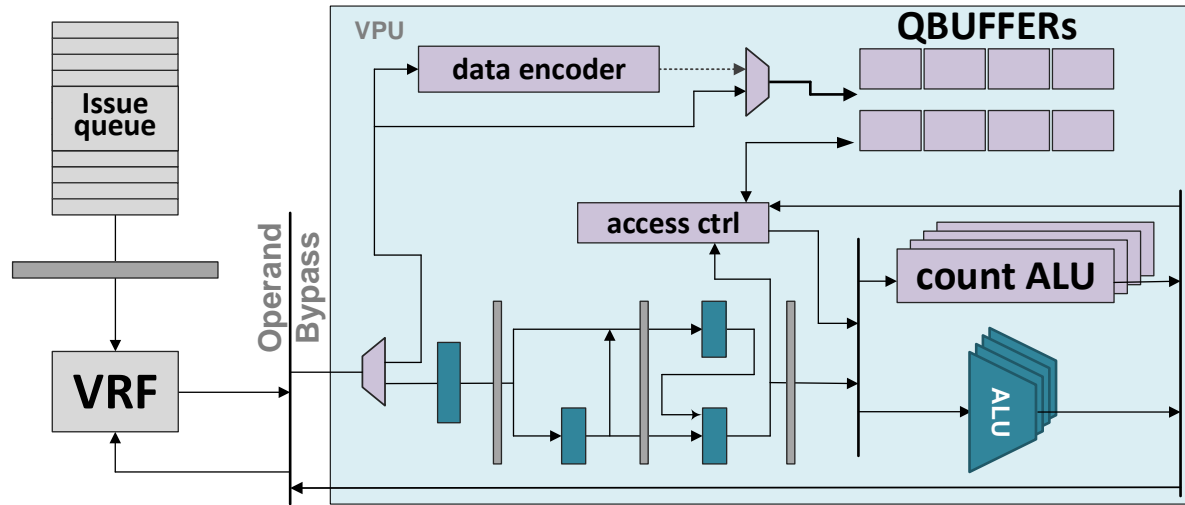
Bottlenecks for Vector Architectures



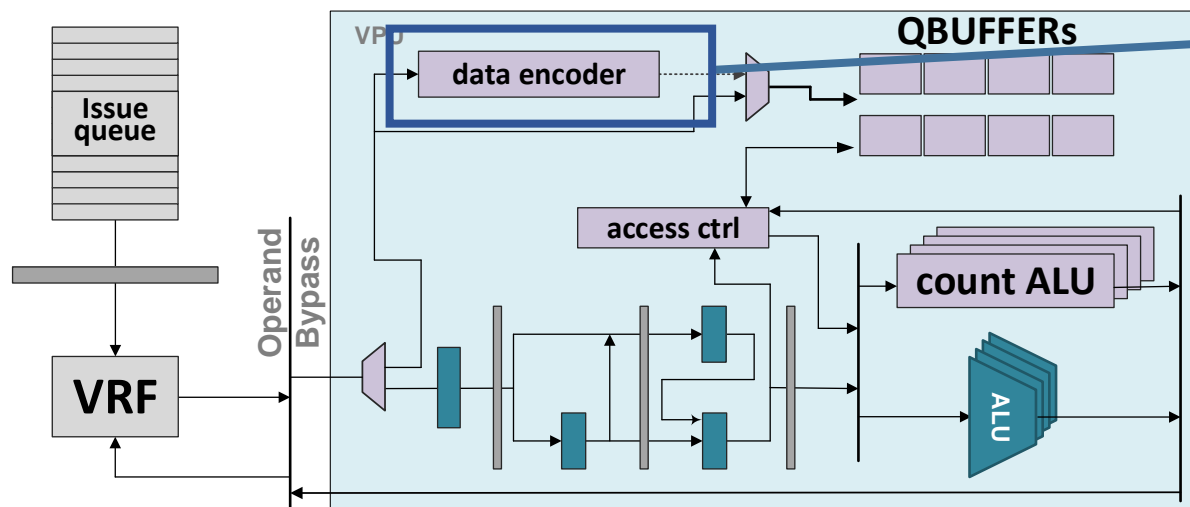
Outline

1. Introduction and motivation
2. Background
 - a) Modern Genome Analysis Algorithms
 - b) Bottlenecks
- 3. QUETZAL**
 - a) Insights and Functionality**
4. Evaluation
5. Summary

QUETZAL

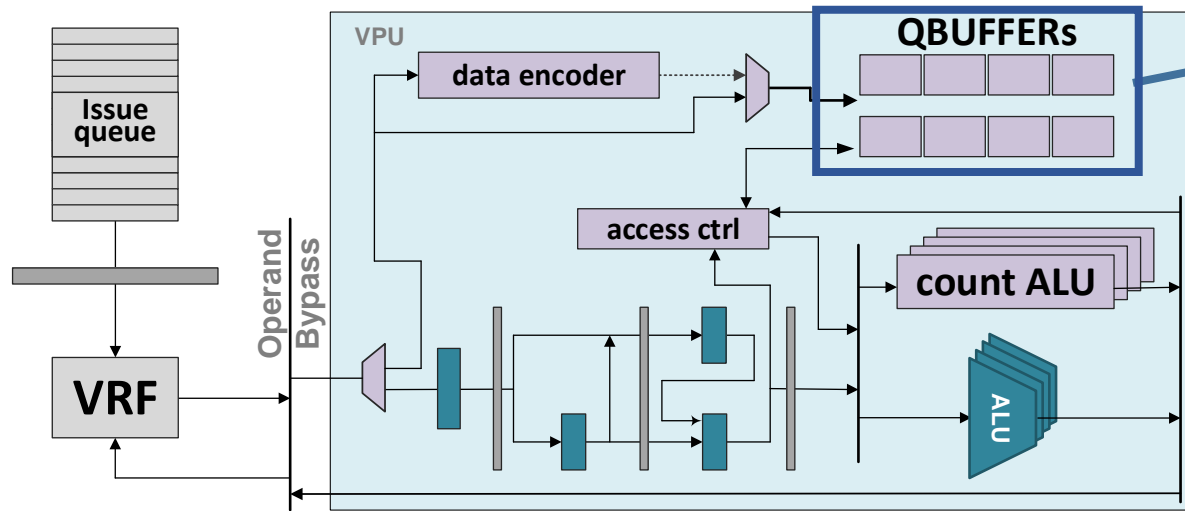


QUETZAL



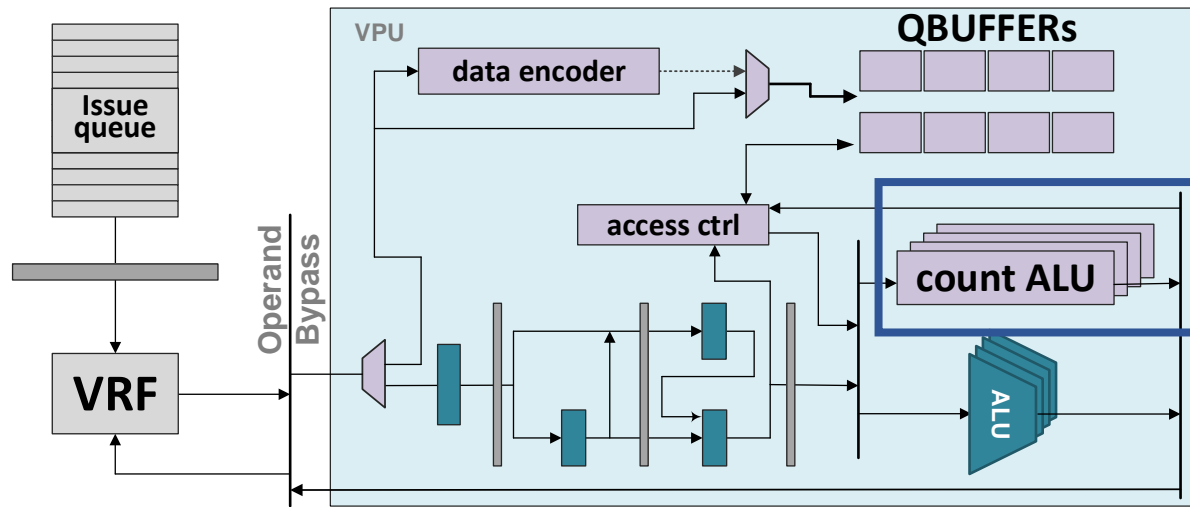
- QUETZAL supports bit-encoded operations without extra instructions.

QUETZAL



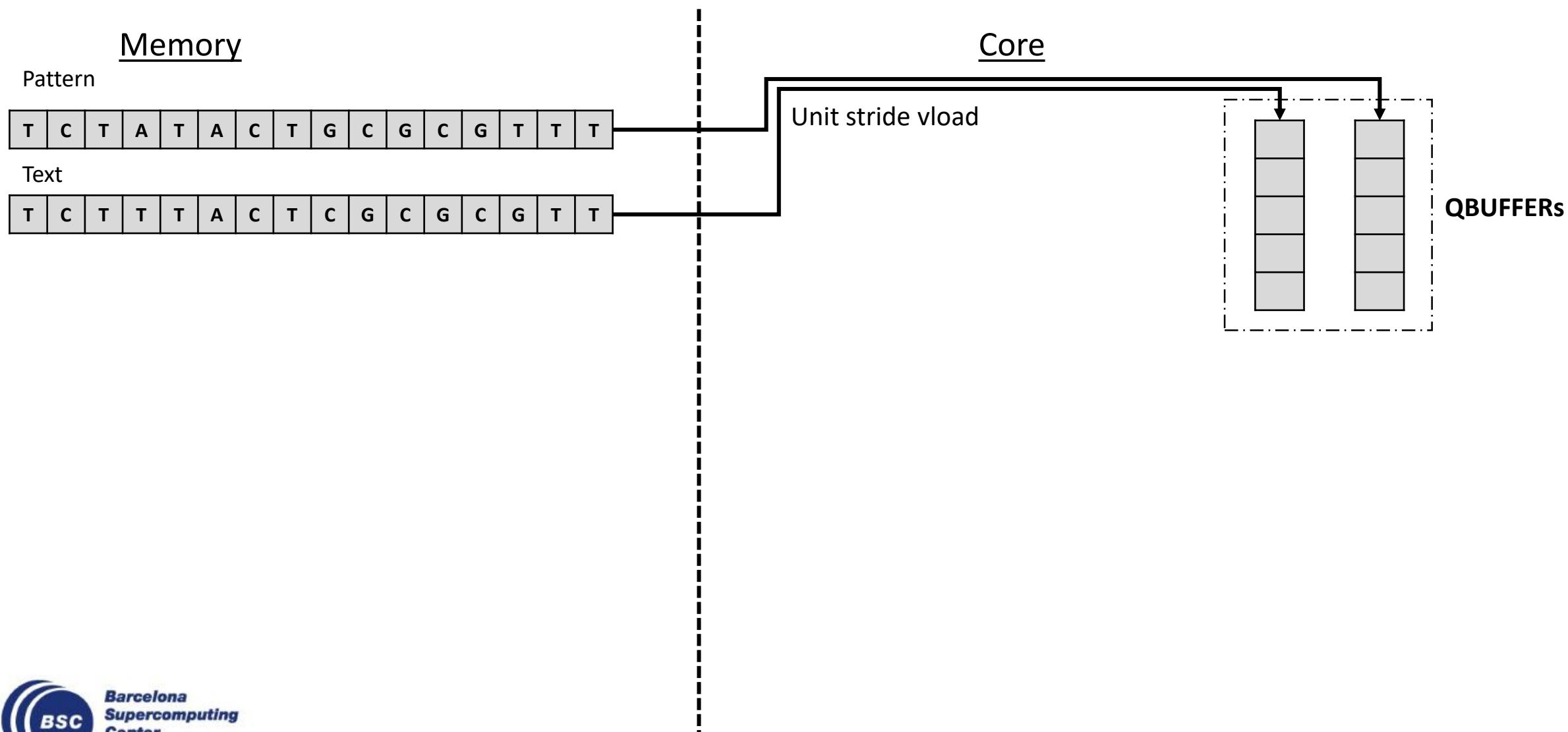
- QUETZAL supports bit-encoded operations without extra instructions.
- QBUFFERS reduce the access latency from 22 cycles to only 2 cycles.

QUETZAL

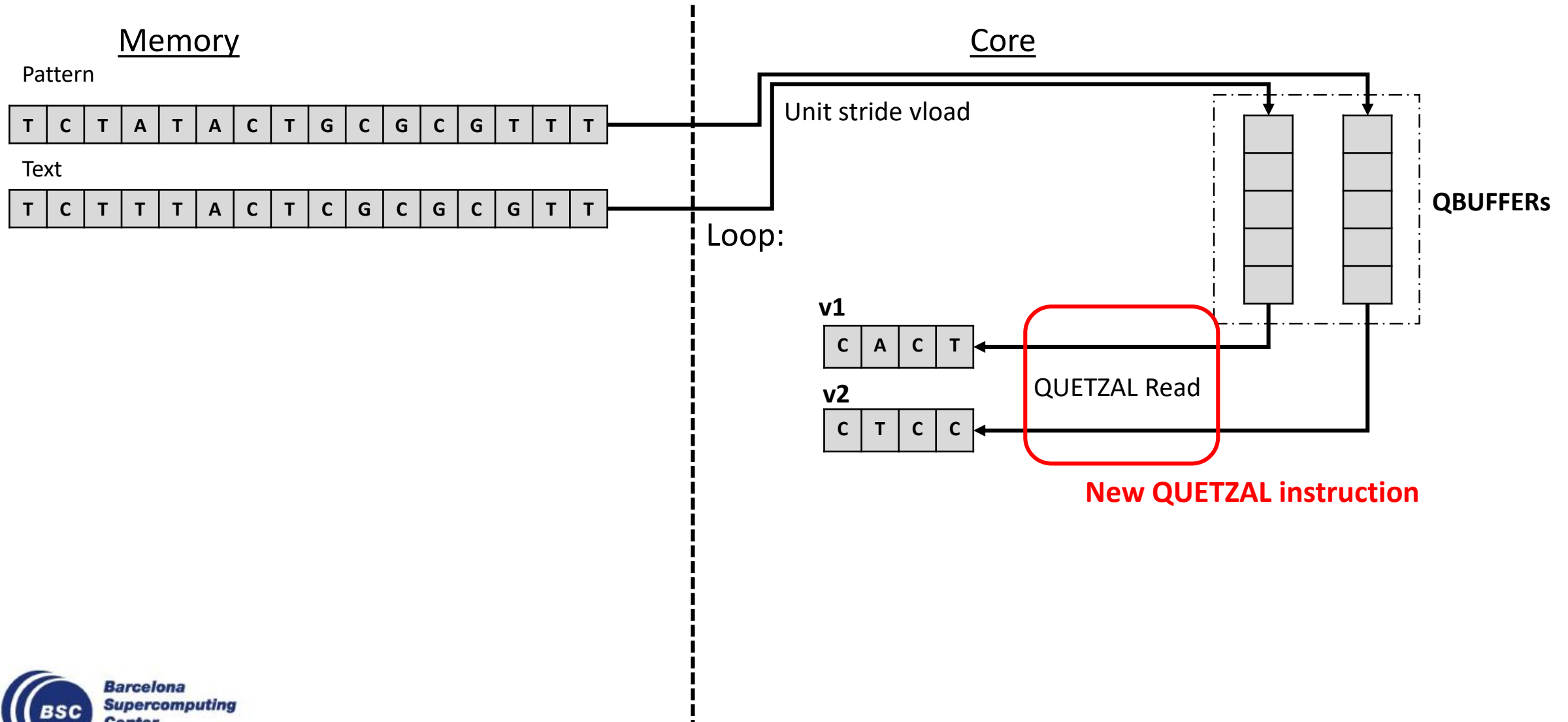


- QUETZAL supports bit-encoded operations without extra instructions.
- QBUFFERS reduce the access latency from 22 cycles to only 2 cycles.
- QUETZAL features custom hardware to calculate the maximum number of exact matches.

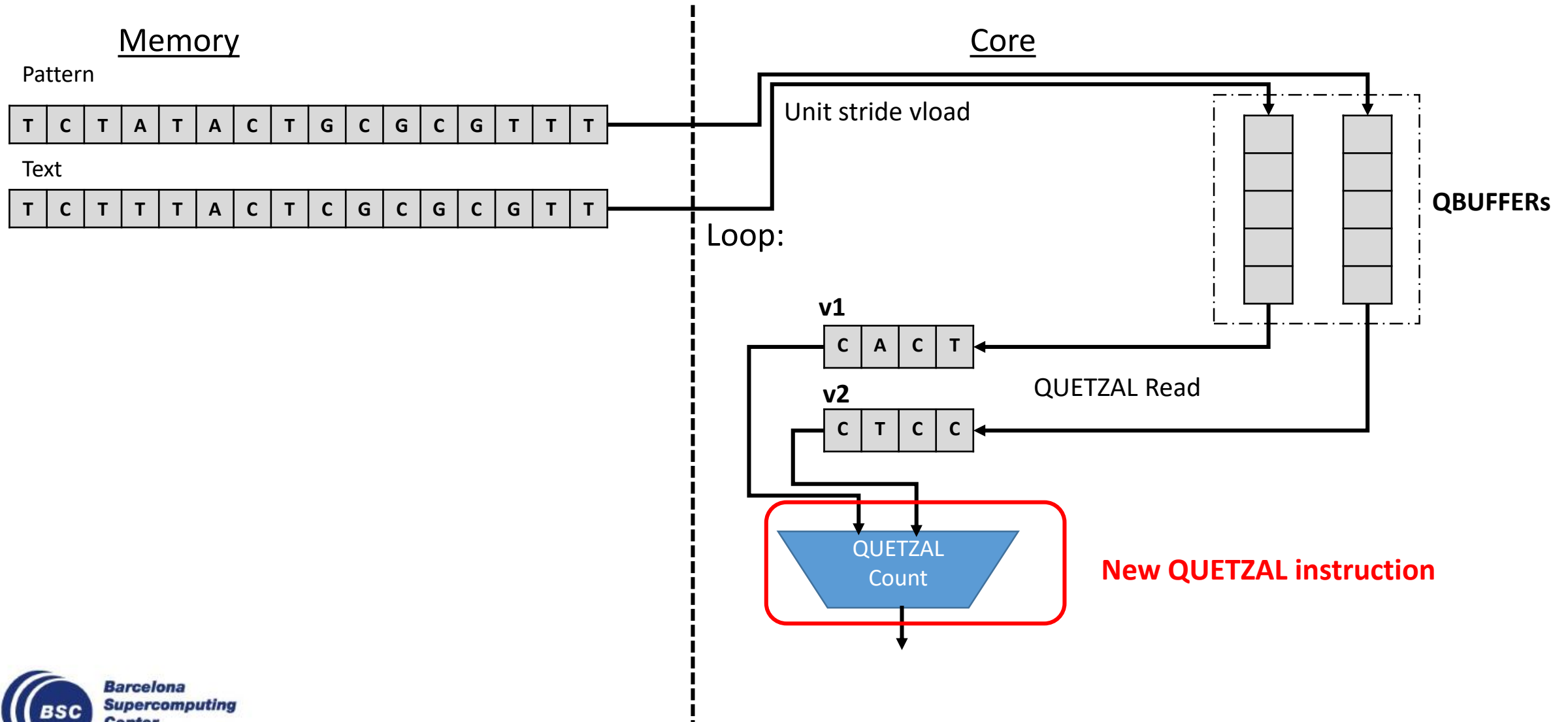
Bottlenecks for Vector architectures



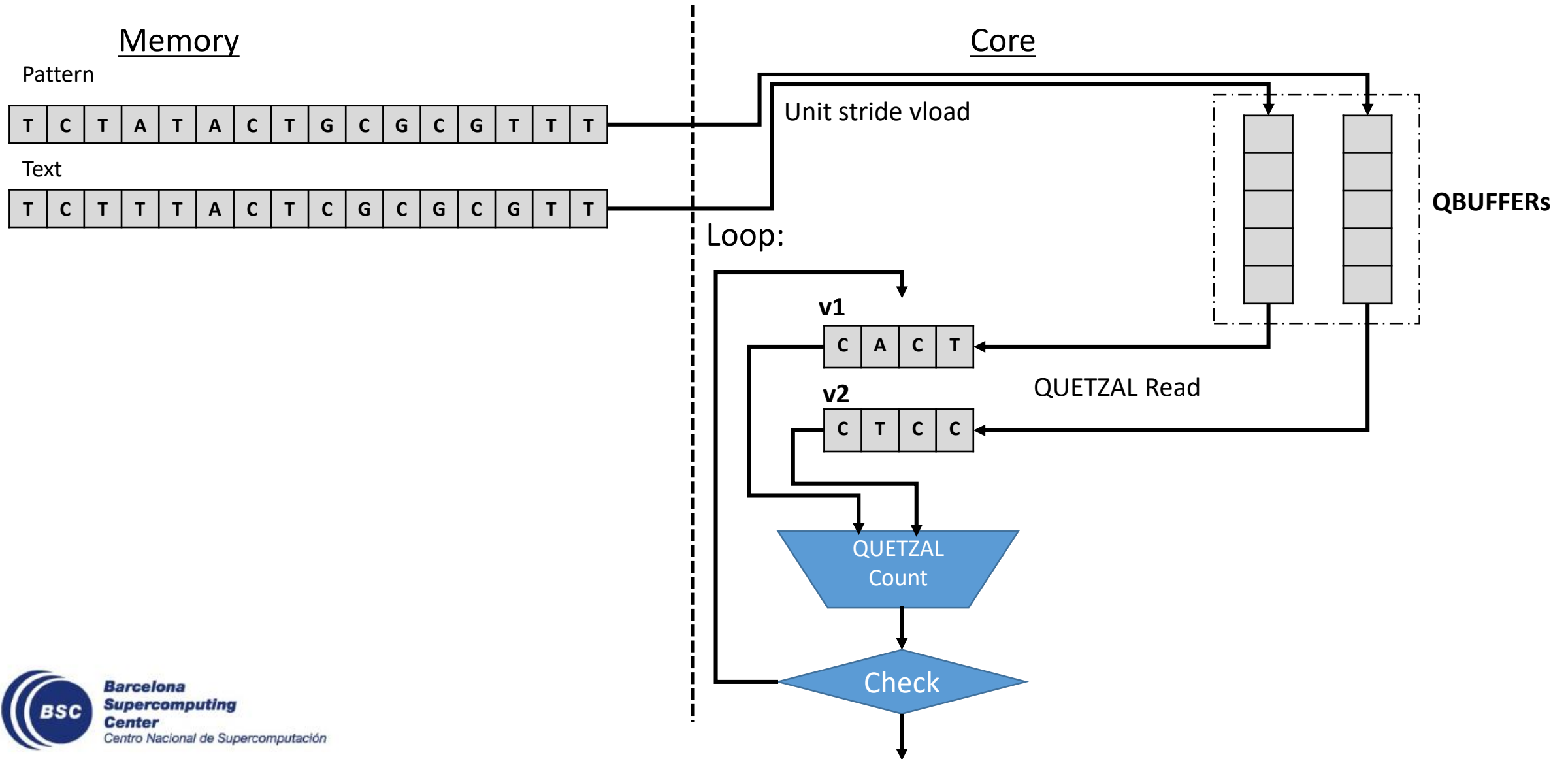
Bottlenecks for Vector architectures



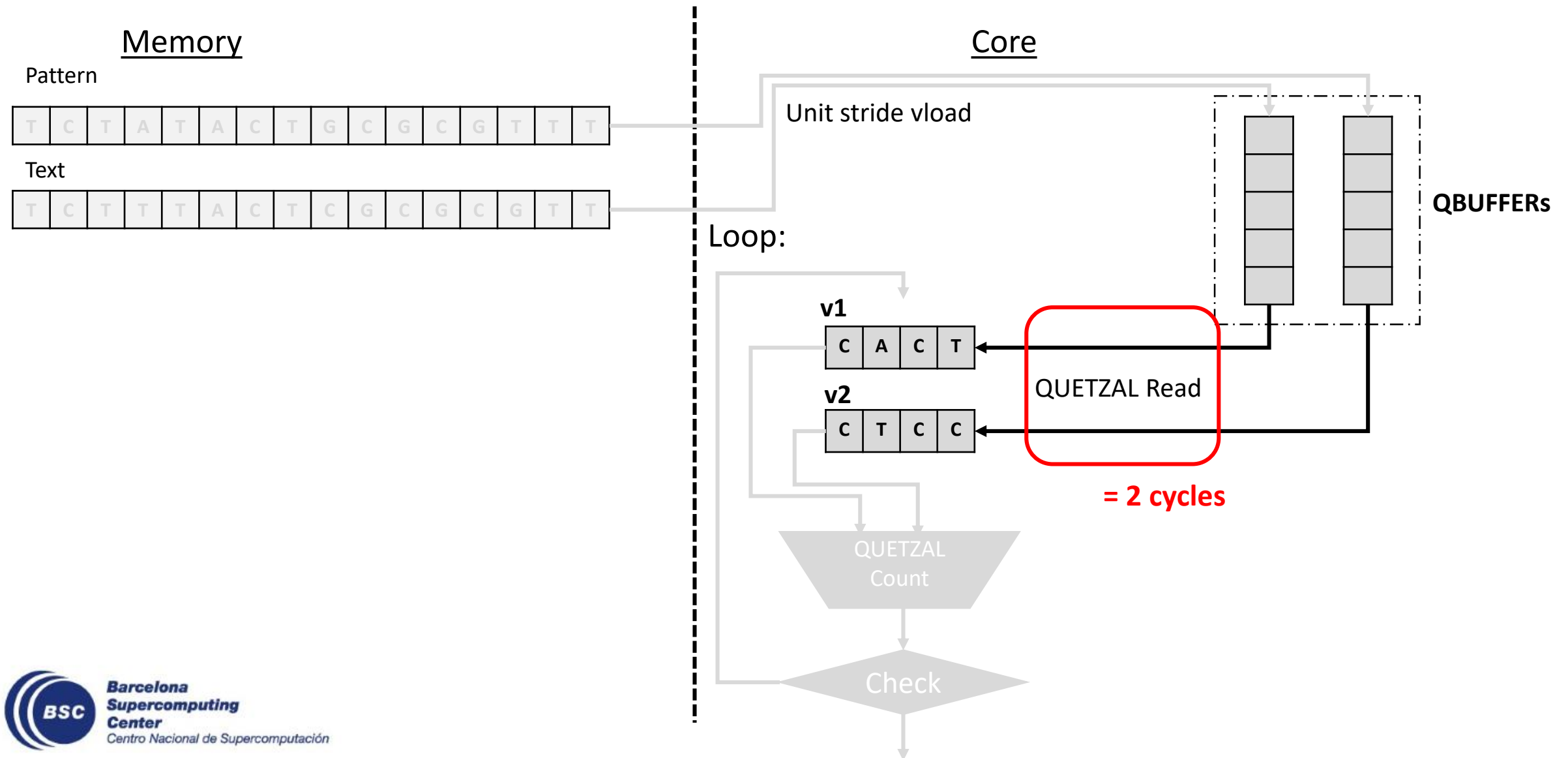
Bottlenecks for Vector architectures



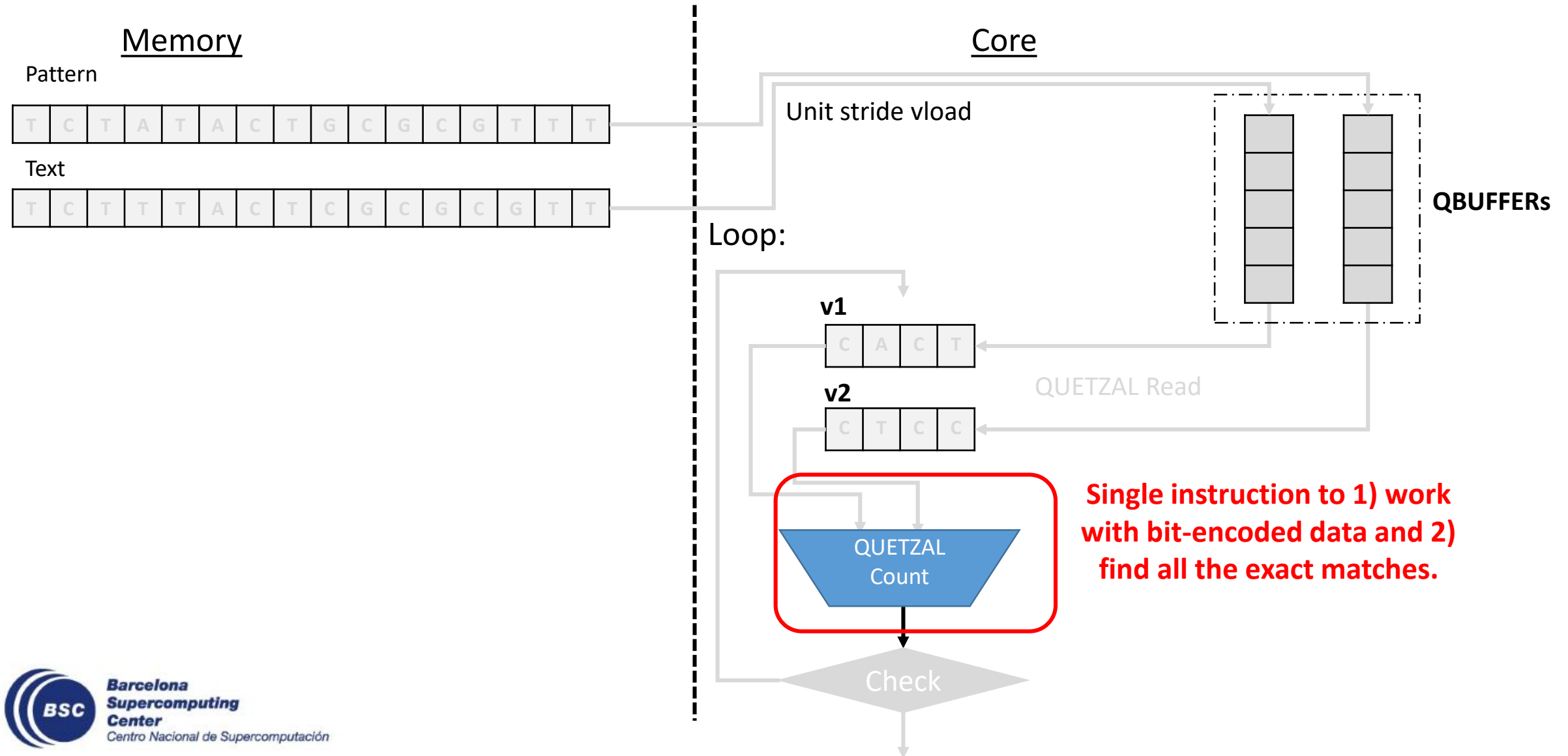
Bottlenecks for Vector architectures



Bottlenecks for Vector architectures



Bottlenecks for Vector architectures

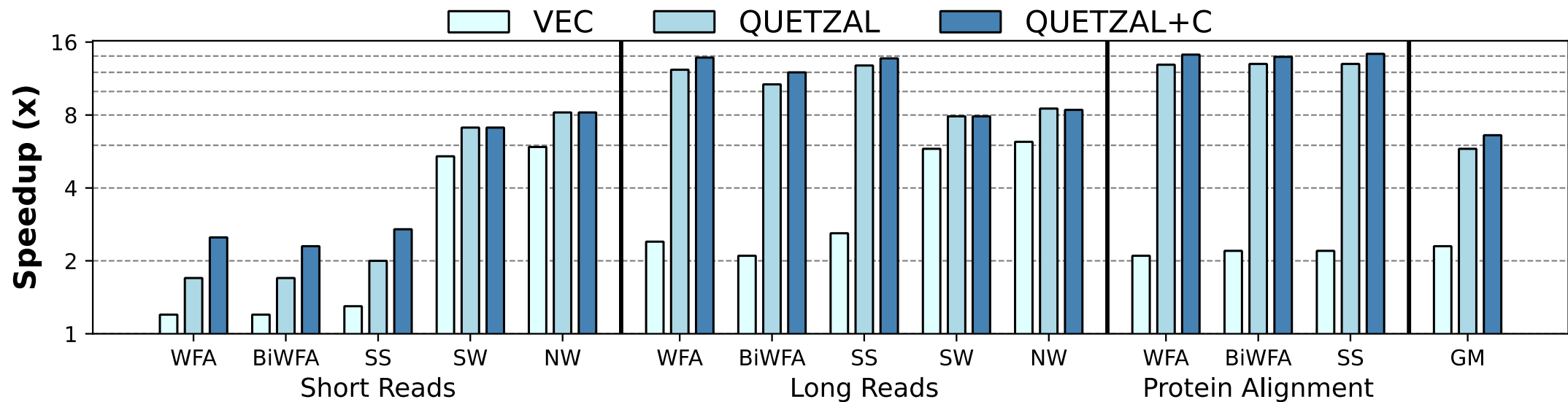


Outline

1. Introduction and motivation
2. Background
 - a) Modern Genome Analysis Algorithms
 - b) Bottlenecks
3. QUETZAL
 - a) Insights and Functionality
- 4. Evaluation**
5. Conclusion

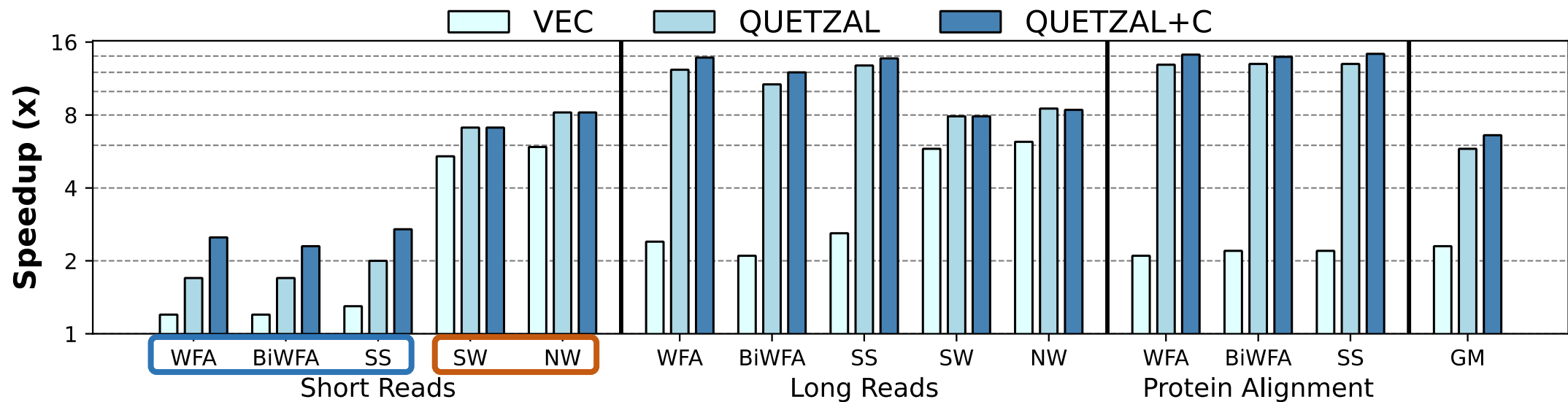
Evaluation

- We use the **gem5 simulator** to model a **16-core ARM 64-bit full-system** running an **Ubuntu 20.04** with a 4.18.0+ Linux Kernel.
- All results are normalized to a CPU baseline.



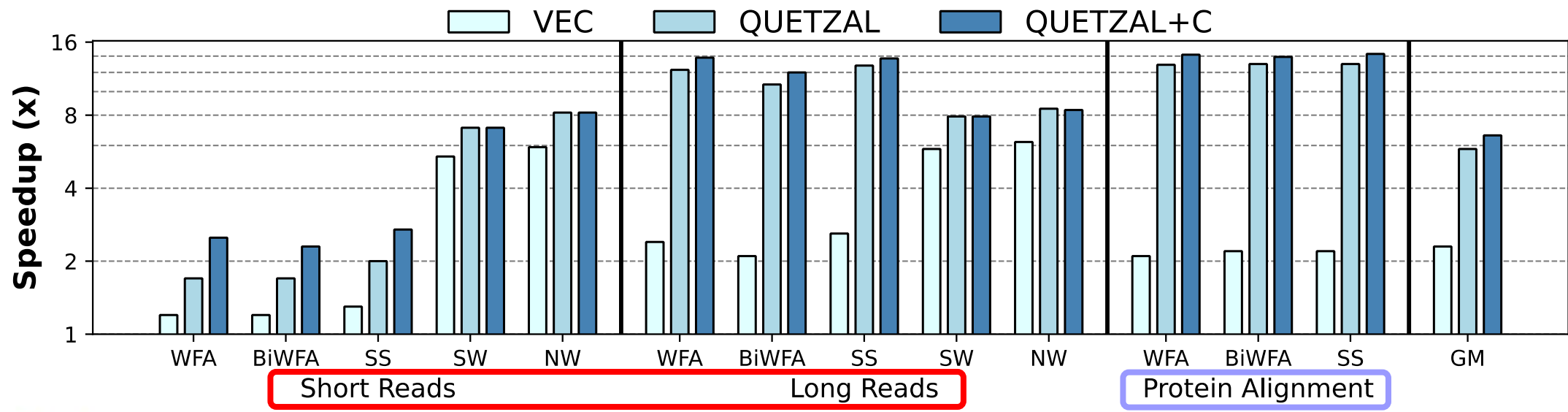
Evaluation

- Five algorithms: **three modern** and **two classical** algorithms.



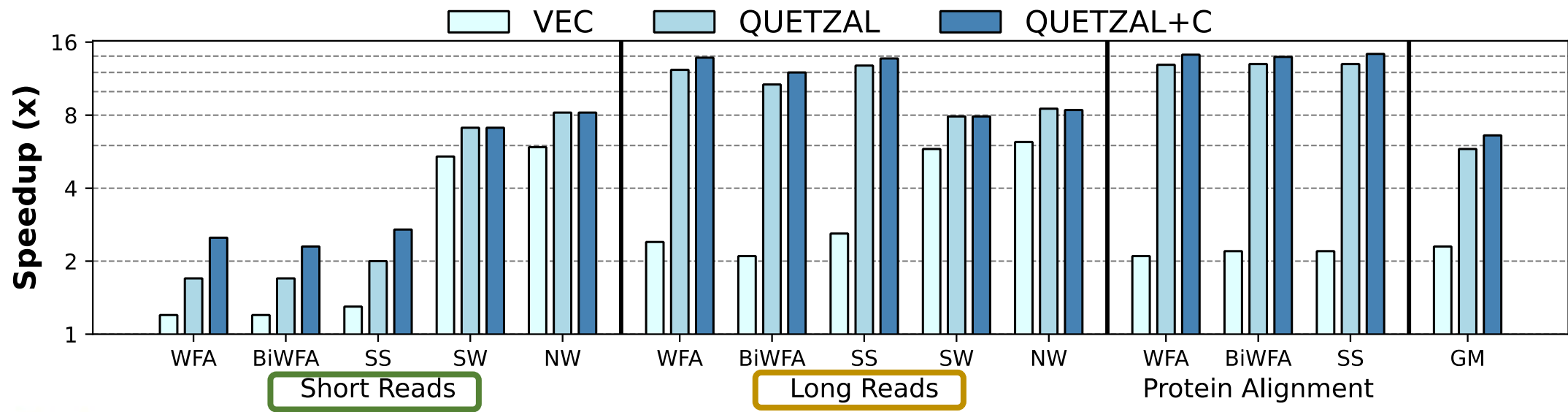
Evaluation

- Five algorithms: **three modern** and **two classical** algorithms.
- Two alphabets: **DNA/RNA** and **Proteins**.



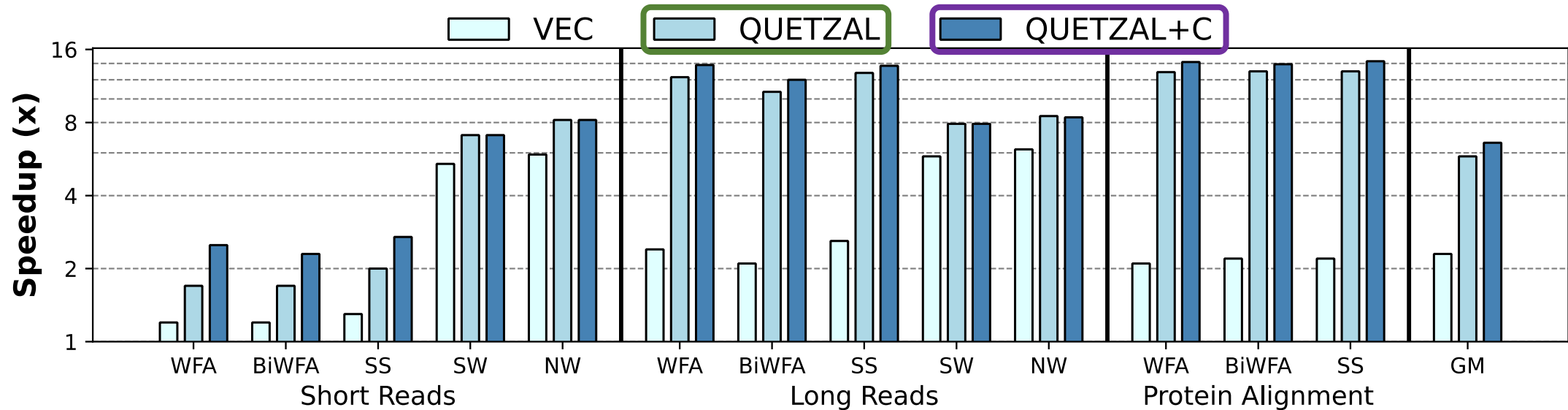
Evaluation

- Five algorithms: **three modern** and **two classical** algorithms.
- Two alphabets: DNA/RNA and Proteins.
- Both **short** and **long** DNA/RNA sequences.



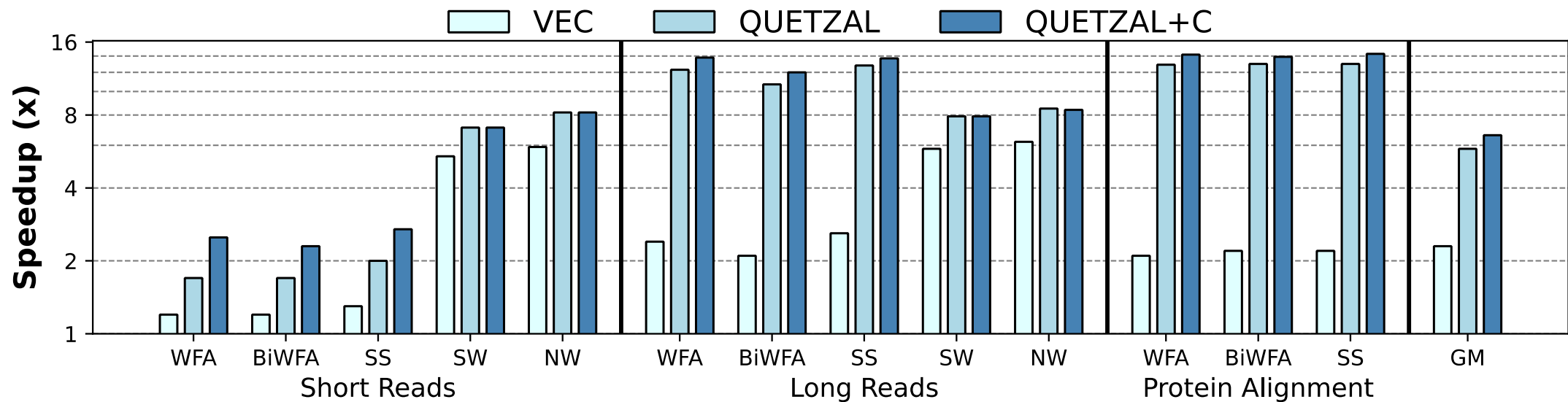
Evaluation

- Five algorithms: **three modern** and **two classical** algorithms.
- Two alphabets: DNA/RNA and Proteins.
- Both short and long DNA/RNA sequences.
- Two QUETZAL implementations



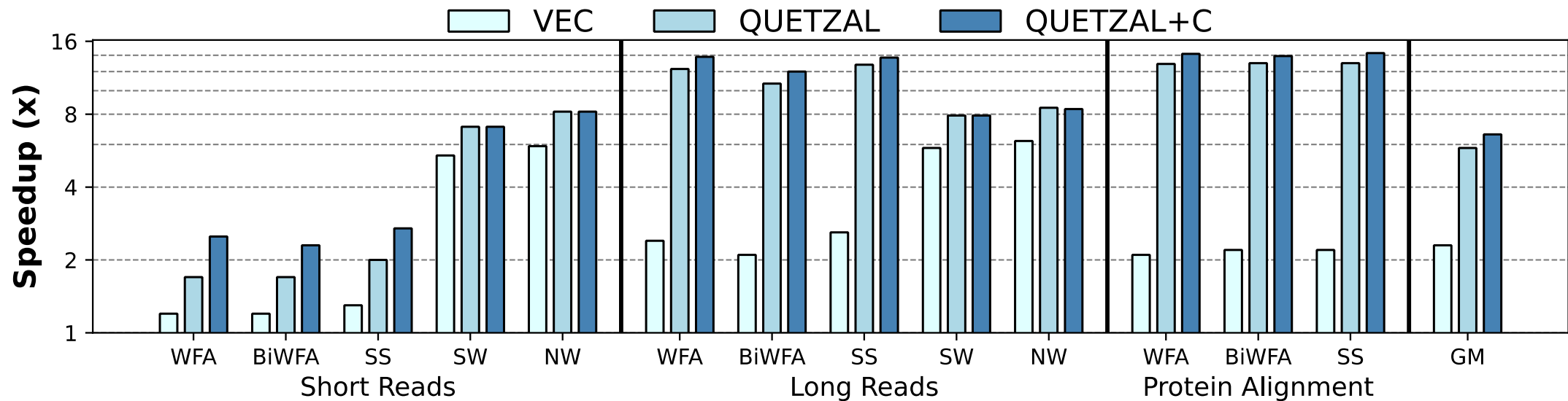
Evaluation

- QUETZAL significantly outperforms all the evaluated algorithms.



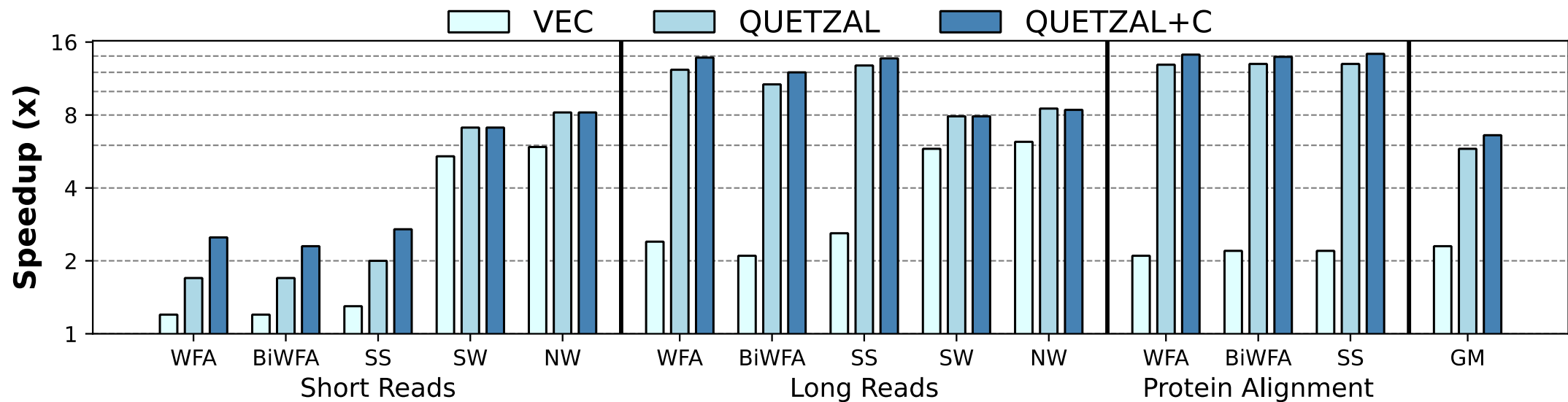
Evaluation

- QUETZAL significantly outperforms all the evaluated algorithms.
- **Speedup:** 5.7x better performance compared to other vectorized algorithms.



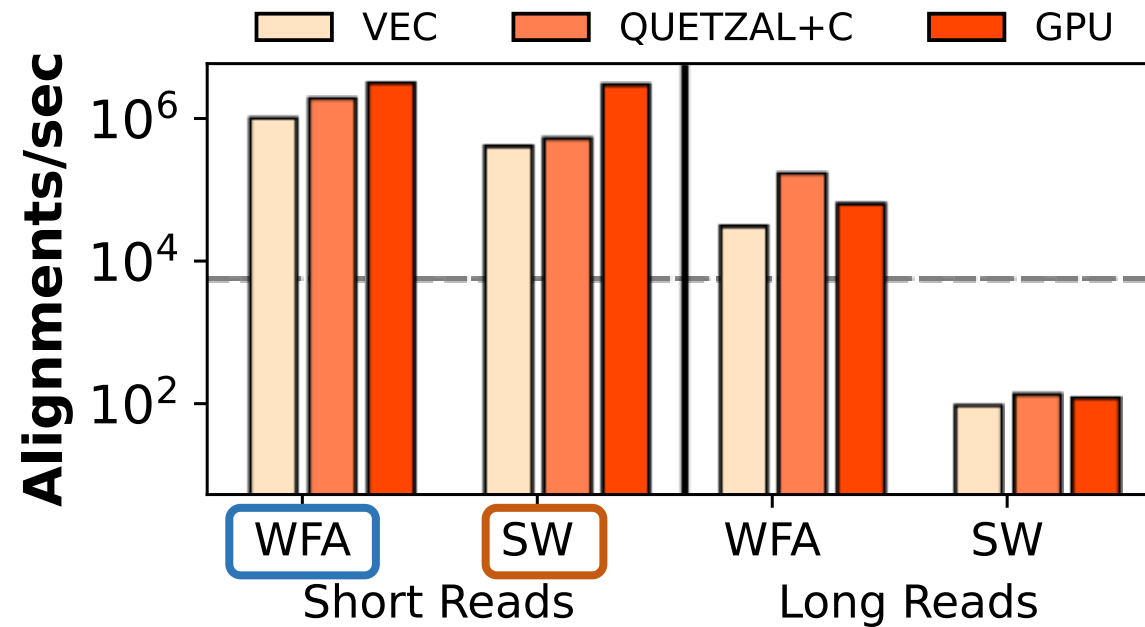
Evaluation

- QUETZAL significantly outperforms all the evaluated algorithms.
- **Speedup:** 5.7x better performance compared to other vectorized algorithms.
- QUETZAL is capable of accelerating both modern and classical genome sequence analysis algorithms.



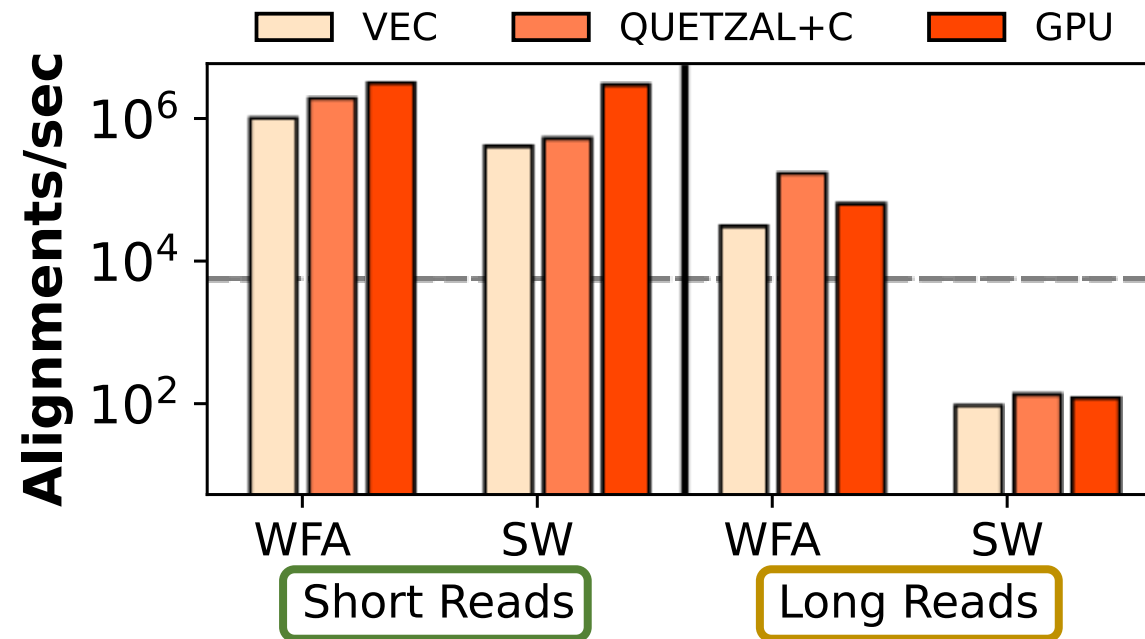
Evaluation

- Two algorithms: **one modern** and **one classical** algorithms.



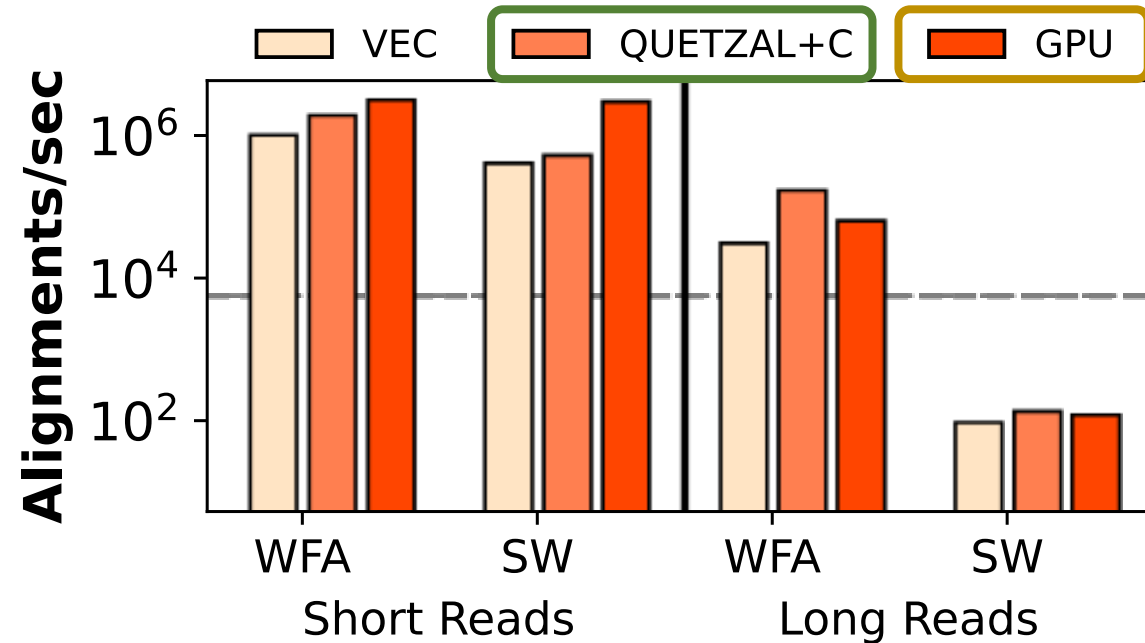
Evaluation

- Two algorithms: **one modern** and **one classical** algorithms.
- Both **short** and **long** DNA/RNA sequences.



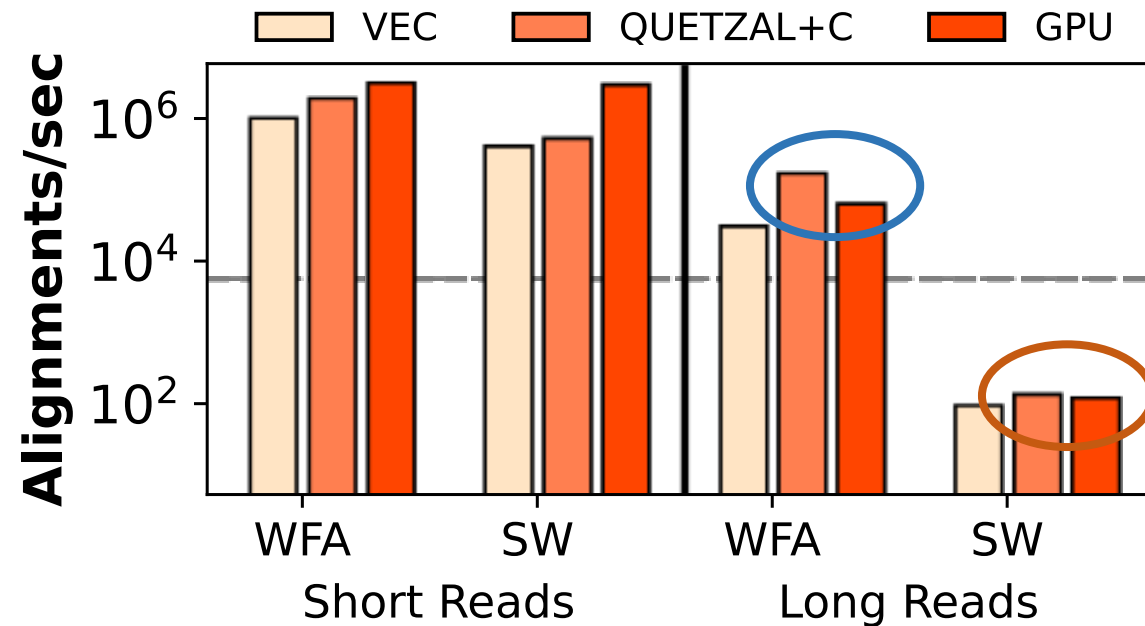
Evaluation

- Two algorithms: **one modern** and **one classical** algorithms.
- Both short and long DNA/RNA sequences.
- we compared a **16-core CPU** featuring **QUETZAL** and an **NVIDIA A40 GPU**.



Evaluation

- When processing long reads, QUETZAL outperform by 2.7x and 1.1x for WFA and Smith-Waterman (Classical algorithm), respectively.



Area overhead results

We evaluated the area overhead of QUETZAL by going all the way **from RTL to Place&Route** our design using **SystemVerilog** and the **Synopsys' ICC2 Place and Route tool** targeting **7nm** node technology and **2GHz** frequency.

Our evaluation shows a small area overhead of

1.4%

compared to a Fujitsu A64FX SoC.

Other content of the paper

1. Detailed explanation of all the hardware components in QUETZAL.
2. Examples to integrate QUETZAL in a C++ code.
3. Single- and multi-core scalability.
4. Design Space Exploration to right-size the QBUFFERS.
5. QUETZAL performance for non-genomic applications.
6. Comparison to ASIC accelerators.

Outline

1. Introduction and motivation
2. Background
 - a) Modern Genome Analysis Algorithms
 - b) Bottlenecks
3. QUETZAL
 - a) Insights and Functionality
4. Evaluation
5. Conclusion

Conclusion

- Striking a balance between high performance and generality is critical
- CPU's vector data path offers an attractive design choice for data parallel applications
- QUETZAL – hardware-software co-design for genome sequencing with 5.7x compared to baseline CPU

QUETZAL: Vector Acceleration Framework for Modern Genome Sequence Analysis Algorithms

Lead authors: Julian Pavon, Ivan Vargas Valdivieso,
Carlos Rojas, Cesar Hernandez

Co-authors: Mehmet Aslan, Roger Figueras, Yichao
Yuan, Joel Lindegger, Mohammed Alser, Francesc Moll,
Santiago Marco Sola, Oguz Ergin, Nishil Talati, Onur
Mutlu, Osman Unsal, Mateo Valero and Adrian Cristal

ETH zürich

