

Helicase: Vectorized parsing and bitpacking of genomic sequences

Igor MARTAYAN, Loup LOBET, Camille MARCHET, Charles PAPERMAN

RECOMB-Arch 2026, May 24th, Thessaloniki



Context: FASTA and FASTQ formats

```
1 >header1 fasta
2 GATTAATCAC
3 >a longer header <0.0>
4 TTAGCTGTGC
5 GANNNTAGAT
6 CTA
7 >header3 ...
```

```
1 @header1 fastq
2 GATTAATCAC
3 +
4 whXa~~~k5h
5 @header2
6 TTAGCTGTGCGANNNTAGATCTA
7 +
8 @This+is+a+valid+score!
```

Context: FASTA and FASTQ formats

```
1 >header1 fasta  
2 GATTAATCAC  
3 >a longer header <0.0>  
4 TTAGCTGTGC  
5 GANNNTAGAT  
6 CTA  
7 >header3 ...
```

```
1 @header1 fastq  
2 GATTAATCAC  
3 +  
4 whXa~~~k5h  
5 @header2  
6 TTAGCTGTGCGANNNTAGATCTA  
7 +  
8 @This+is+a+valid+score!
```

Context: FASTA and FASTQ formats

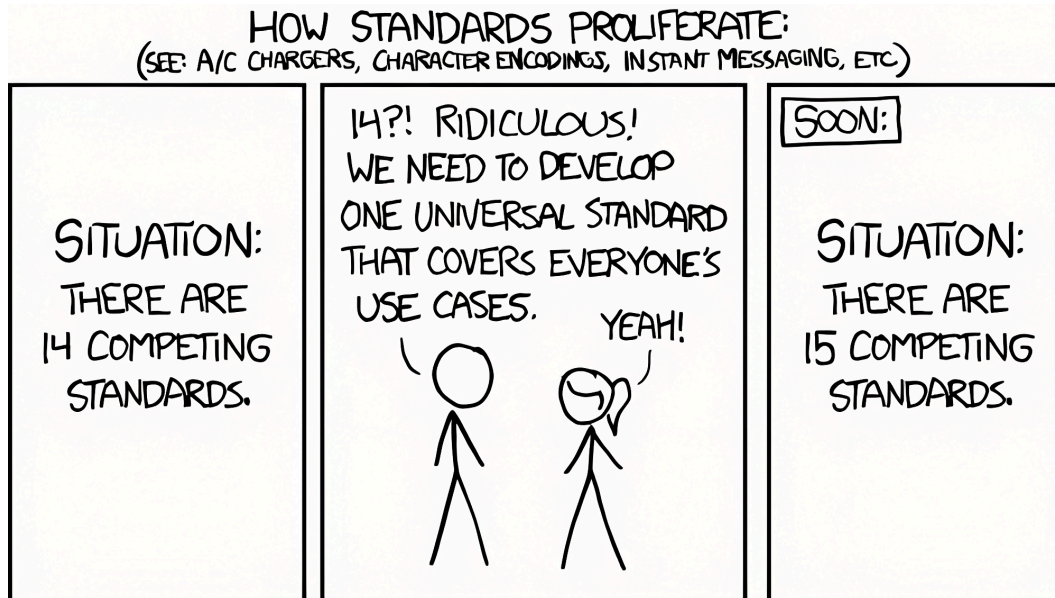
```
1 >header1 fasta  
2 GATTAATCAC  
3 >a longer header <0.0>  
4 TTAGCTGTGC  
5 GANNNTAGAT  
6 CTA  
7 >header3 ...
```

non-ACTG bases (IUPAC)



```
1 @header1 fastq  
2 GATTAATCAC  
3 +  
4 whXa~~~k5h  
5 @header2  
6 TTAGCTGTGCGANNNTAGATCTA  
7 +  
8 @This+is+a+valid+score!
```

Why another FASTA/Q parser?



xkcd.com/927

Popular parsers:

kseq, needletail, paraseq...

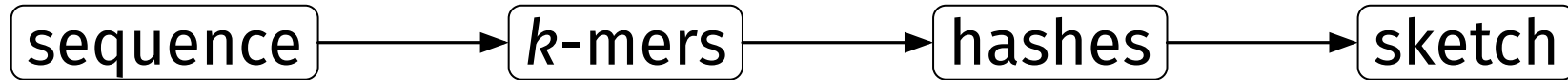
→ all rely on memchr,
no advanced vectorization

→ especially slow on
multiline FASTA

→ only output ASCII
sequences, no bitpacking

From text-based formats to bitpacked streams

Many simple tasks can be formulated as streams, e.g.



We want to reduce the overhead of intermediate representations

→ take advantage of modern hardware for high throughput

→ bitpack data for better cache efficiency

Marking FASTA headers with carry propagation

Goal: mark everything between > and ← *branchlessly*

input	>	c	h	r	l	←	C	G	G	A	C	←	A	C	G	T	←	>	c	>	h	←	C
←	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0
>	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0
!←	1	1	1	1	1	0	1	1	1	1	1	0	1	1	1	1	0	1	1	1	1	0	1

Marking FASTA headers with carry propagation

Goal: mark everything between > and ← *branchlessly*

input	>	c	h	r	l	←	C	G	G	A	C	←	A	C	G	T	←	>	c	>	h	←	C
←	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0
>	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0
!←	1	1	1	1	1	0	1	1	1	1	1	0	1	1	1	1	0	1	1	1	1	0	1
> + !←	0	0	0	0	0	1	1	1	1	1	1	0	1	1	1	1	0	0	0	1	0	1	1

Marking FASTA headers with carry propagation

Goal: mark everything between > and ← *branchlessly*

input	>	c	h	r	l	←	C	G	G	A	C	←	A	C	G	T	←	>	c	>	h	←	C
←	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0
>	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0
!←	1	1	1	1	1	0	1	1	1	1	1	0	1	1	1	1	0	1	1	1	1	0	1
> + !←	0	0	0	0	0	1	1	1	1	1	1	0	1	1	1	1	0	0	0	1	0	1	1
⊕ !←	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	1	1	0	1	1	0

Marking FASTA headers with carry propagation

Goal: mark everything between > and ← *branchlessly*

input	>	c	h	r	l	←	C	G	G	A	C	←	A	C	G	T	←	>	c	>	h	←	C
←	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0
>	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0
!←	1	1	1	1	1	0	1	1	1	1	1	0	1	1	1	1	0	1	1	1	1	0	1
> + !←	0	0	0	0	0	1	1	1	1	1	1	0	1	1	1	1	0	0	0	1	0	1	1
⊕ !←	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	1	1	0	1	1	0
∨ >	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	0

→ we produce a stream of annotation bits

Bitpacked representations of DNA

A	C	T	G
01000001	01000011	01010100	01000111

Packed and columnar representations:

sequence	G	A	T	T	A	C	A	T
packed	11	00	10	10	00	01	00	10
col hi	1	0	1	1	0	0	0	1
col lo	1	0	0	0	0	1	0	0

Bitpacked representations of DNA

A	C	T	G
01000001	01000011	01010100	01000111

Packed and columnar representations:

sequence	G	A	T	T	A	C	A	T
packed	11	00	10	10	00	01	00	10
col hi	1	0	1	1	0	0	0	1
col lo	1	0	0	0	0	1	0	0

What about non-ACTG characters?

Detecting non-ACTG characters

- 1 $LUT \leftarrow ['A', _, 'C', _, 'T', _, 'G', _]$ ▷ constant lookup table mapping bits to ACTG
- 2 $v_{\text{bits}} \leftarrow v_{\text{data}} \& 0b110$ ▷ lossy 2-bit encoding
- 3 $v_{\text{lookup}} \leftarrow \text{LOOKUP}^*(LUT, v_{\text{bits}})$ ▷ map bits to ACTG using the lookup table
- 4 $v_{\text{upper}} \leftarrow v_{\text{data}} \& !0b00100000$ ▷ uppercase the input
- 5 $v_{\text{eq}} \leftarrow \text{CMPEQ}(v_{\text{lookup}}, v_{\text{upper}})$
- 6 **return** $\text{MOVEMASK}^{**}(v_{\text{eq}})$ ▷ 1s indicating ACTG

* implemented with a shuffle

** on ARM: implemented with interleaving and byte-wise packing

Handling non-ACTGs

2 different strategies available:

G A T T N N A C A

→ *split* sequences into chunks of ACTGs

[11 00 10 10] [00 01 00]

→ *lossy encoding* of non-ACTGs + *bitmask* marking them

[11 00 10 10 11 11 00 01 00]

[0 0 0 0 1 1 0 0 0]

Specifying what we want to parse at compile time

```
1  const CONFIG: Config = ParserOptions::default()  
2      .ignore_headers()  
3      .dna_columnar()  
4      .split_non_actg()  
5      .config();
```



Unused fields / representations are eliminated by the compiler

A Rust API similar to needletail

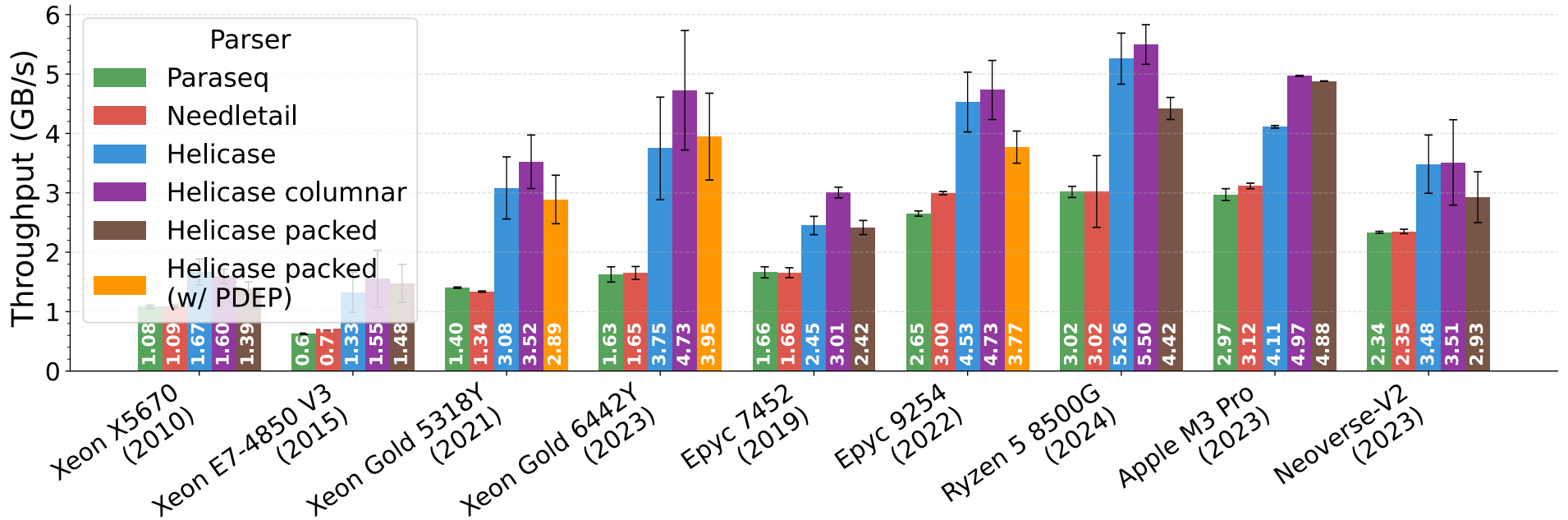
```
1 let mut parser =  
  FastxParser::::from_file(&path)?;  
2 while let Some(_) = parser.next() {  
3   let header = parser.get_header();  
4   let seq = parser.get_dna_columnar();  
5   // ...  
6 }
```



→ automatic format detection, transparent decompression

Benchmarks: parsing multiline FASTA from disk

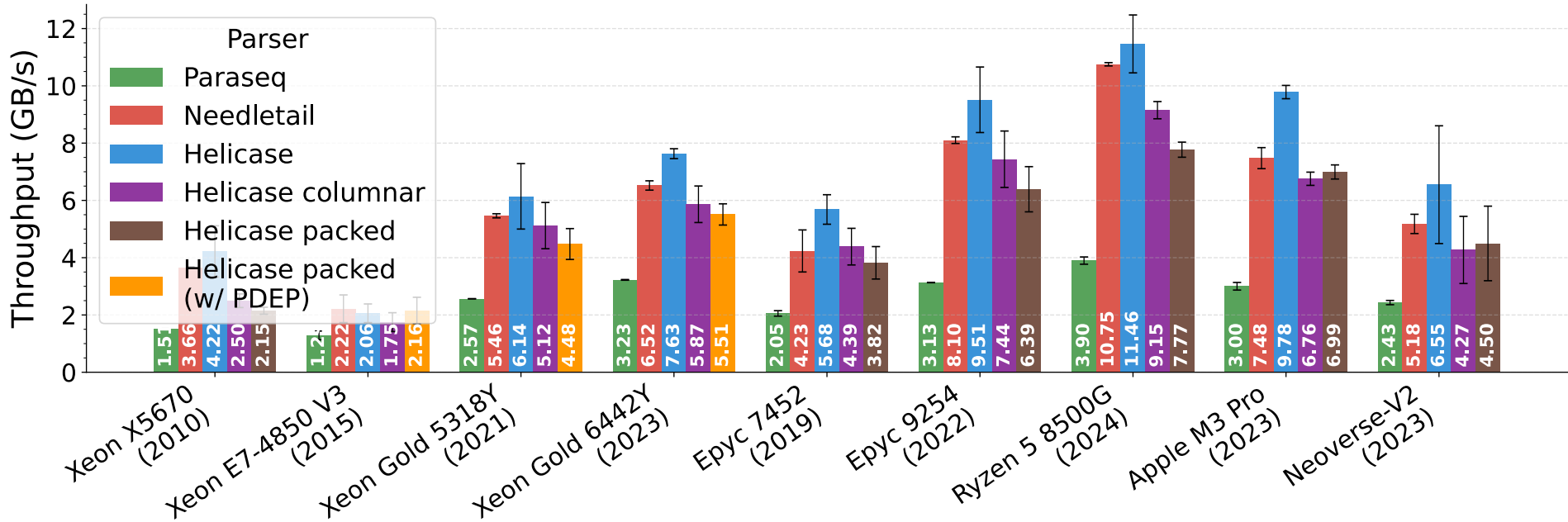
Throughput — Human genome (FASTA)



→ 50-100% faster, columnar is the fastest

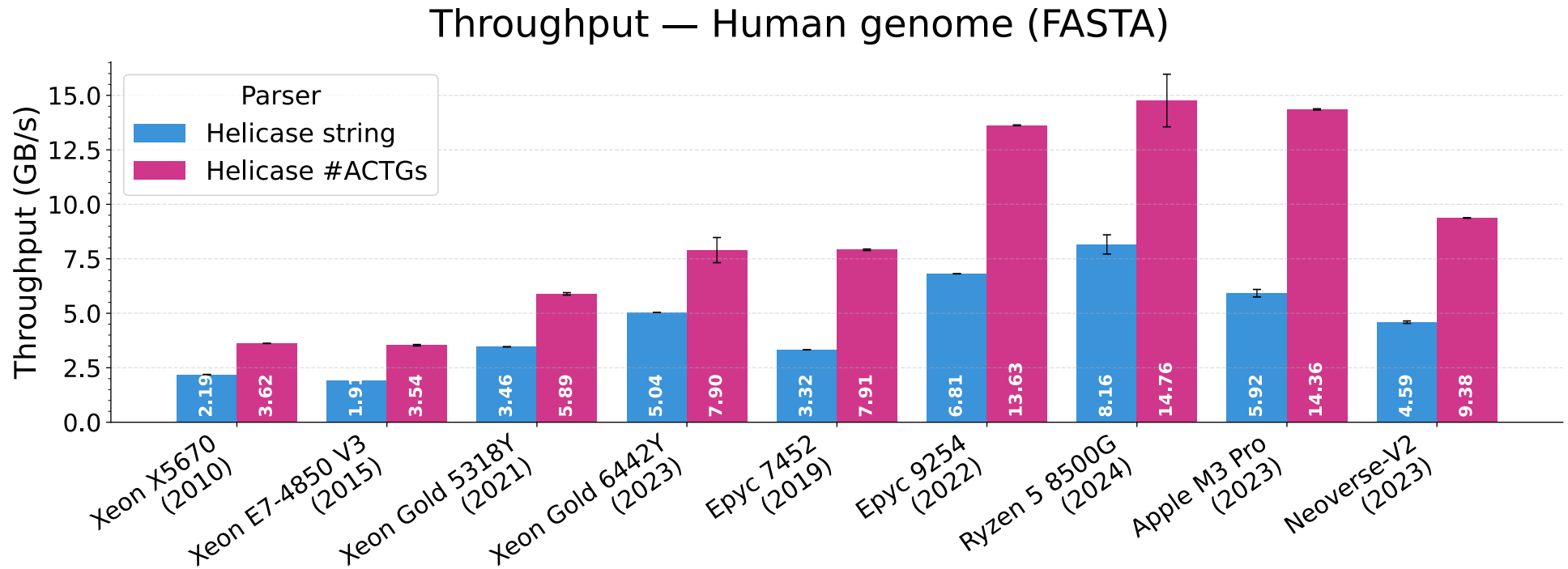
Benchmarks: parsing short reads FASTQ from disk

Throughput — Short reads (FASTQ)



→ 10-20% faster, bitpacking 20-30% more expensive

Benchmarks: counting ACTG bases in RAM



→ specialization can speed up compute-limited tasks

What's next?

Helicase improvements:

- API for parallelization
(without producer/consumer)
- built-in k -mer iterators
- Python/C/C++ bindings

github.com/imartayan/helicase

Avoiding (de)compression
bottlenecks:

- binseq ([Teyssier et al., 2025](#))
- mim ([Patro et al., 2025](#))



What's next?

Helicase improvements:

- API for parallelization (without producer/consumer)
- built-in k -mer iterators
- Python/C/C++ bindings

github.com/imartayan/helicase

Thank you!

Avoiding (de)compression bottlenecks:

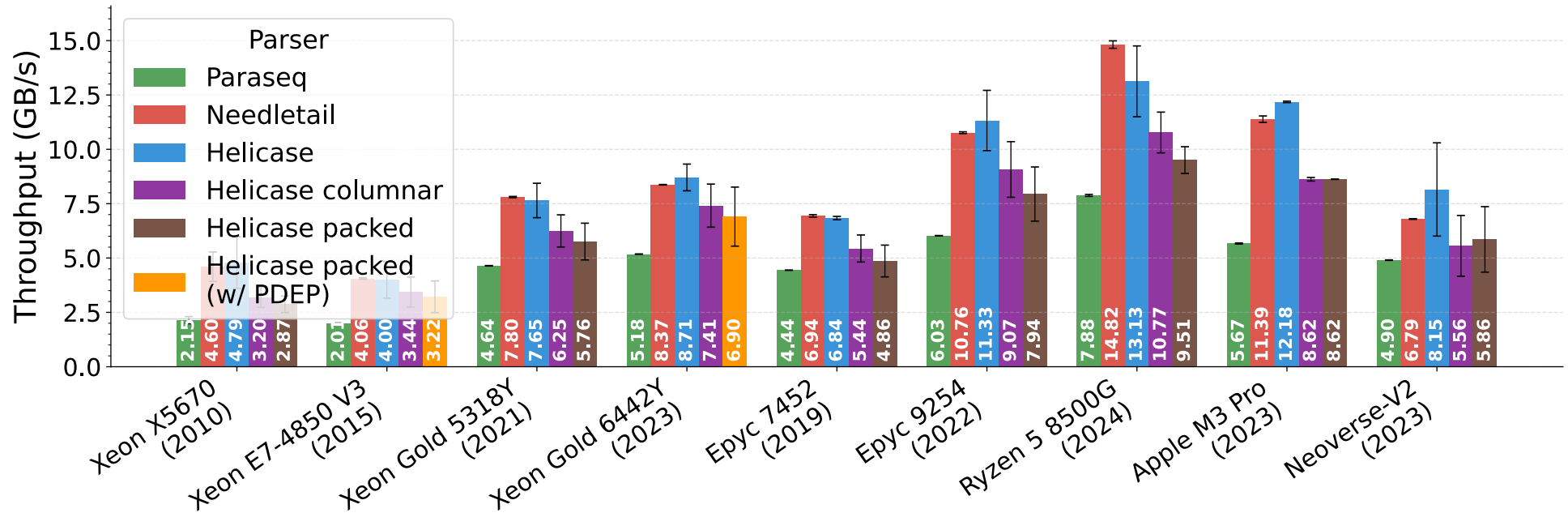
- binseq ([Teyssier et al., 2025](#))
- mim ([Patro et al., 2025](#))



Appendix

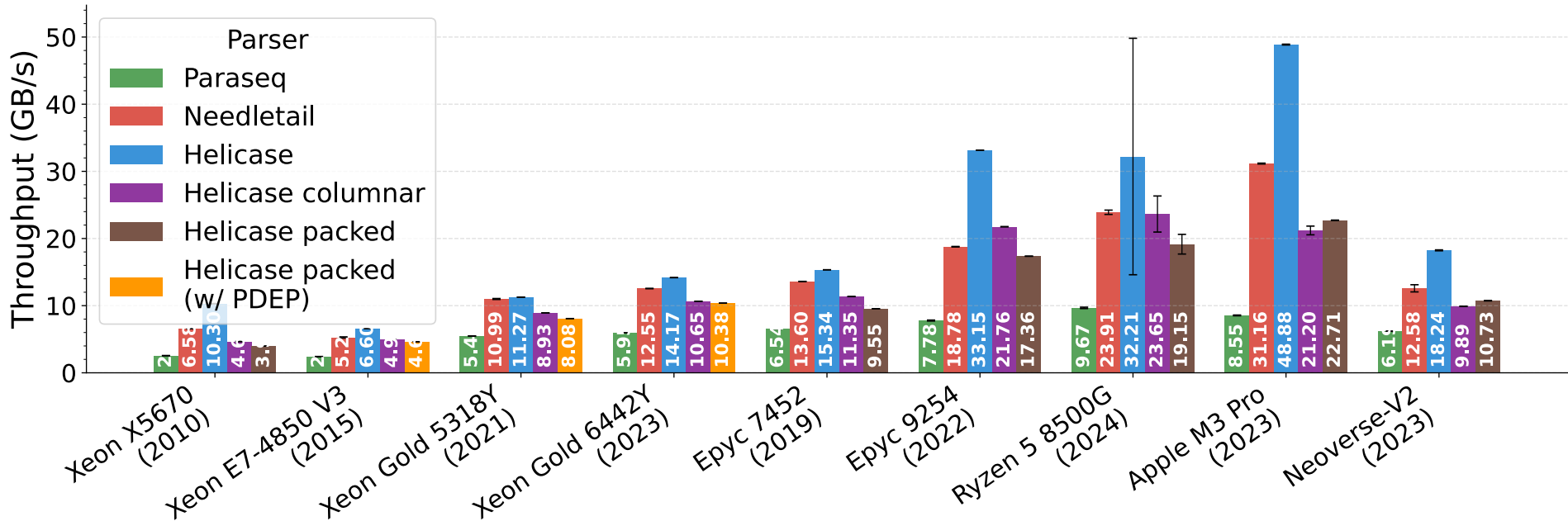
Benchmarks: parsing long reads FASTQ from disk

Throughput — Long reads (FASTQ)



Benchmarks: parsing long reads FASTQ in RAM

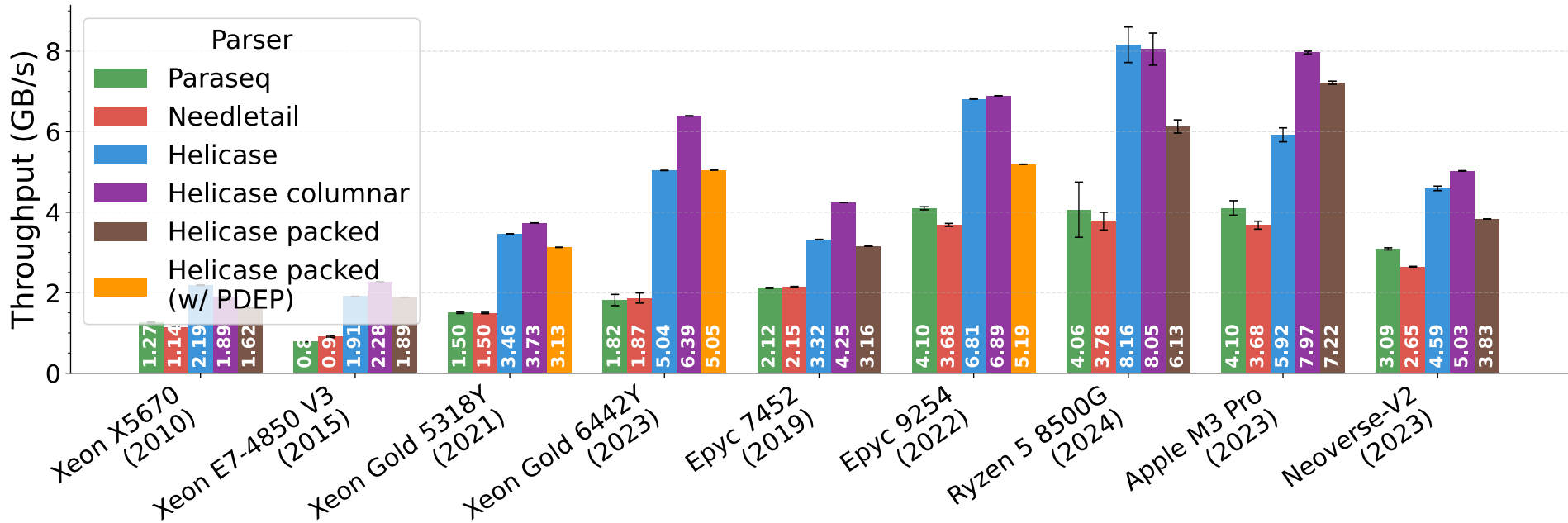
Throughput — Long reads (FASTQ)



→ zero-copy, bandwidth-limited!

Benchmarks: parsing multiline FASTA in RAM

Throughput — Human genome (FASTA)



→ slower throughput than FASTQ, compute-limited