



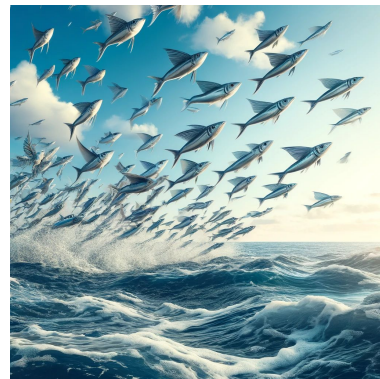
# FFIndex: Compressed Indexing for Parallel FAST-Q Parsing

## Team Members:

Siddhant Bharti

Prajwal Singhanian

Rakrish Dhakal



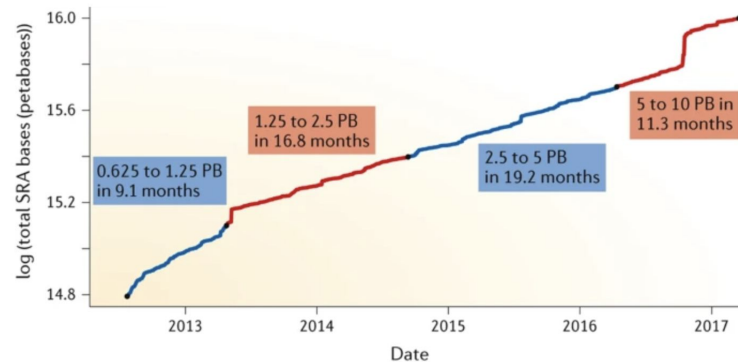
# Project Description

- Our project's aim is to address the challenge of efficiently parsing compressed FASTQ files. This is critical for genomic data processing pipelines but is hindered by the sequential nature of current decompression and parsing methods.
- The goal is to develop a tool and accompanying library that facilitates parallel processing by creating and utilizing an index over compressed FASTQ files.

# Motivation

- Sequence Read Archive (SRA) is predicted to double every 12-18 months. A solution that allows for parallel reading of compressed FASTQ files would facilitate research that relies on large-scale genomic datasets<sup>[1]</sup>.
- Most files hosted on SRA are not compressed in blocks<sup>[2]</sup>.
- If multiple files are to be read in an application, a file-parallel solution can work.

Figure 1: Increase in storage of next-generation sequencing data.



- What about a single file though?

# Motivation

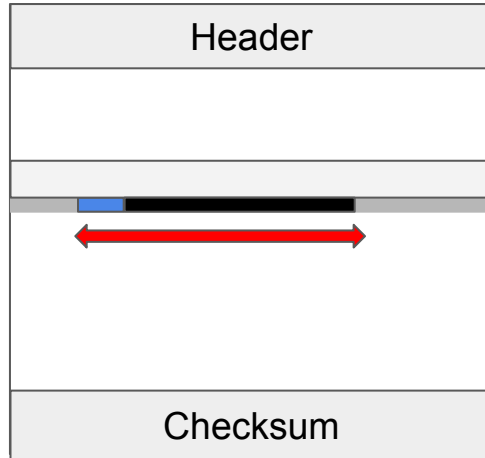
- Single FASTQ files used in large sequencing data experiments can be extremely large. For example, the size of the BGISEQ (MGISEQ-2000RS)<sup>[2]</sup> FASTQ file is around 37GBs.
- gunzip can read compressed data at around 30-50 MB/sec which is order of magnitude less than read throughput of current SATA/NVMe solid-state drives<sup>[2]</sup>.

```
[sbharti@sfeiziwks01 CMSC701-Project]$ ./countbases.out data/nematode.fq.gz | grep -E "Time"
Time taken (to read): 16106 milliseconds
Time taken (total): 25594 milliseconds
[sbharti@sfeiziwks01 CMSC701-Project]$
```

- At this rate, doing analysis for large gzip files can be time and cost-prohibitive.
- **A clever solution to decompress gzip files in parallel is needed.**

# Random Access to GZIP files?

- Gzip achieves best size reduction for FASTQ files. However, gzip's DEFLATE (mLZ77 & huffman coding) algorithm is **sequential**. This makes **random access** in compressed FASTQ files **hard**.



files<sup>[1]</sup>

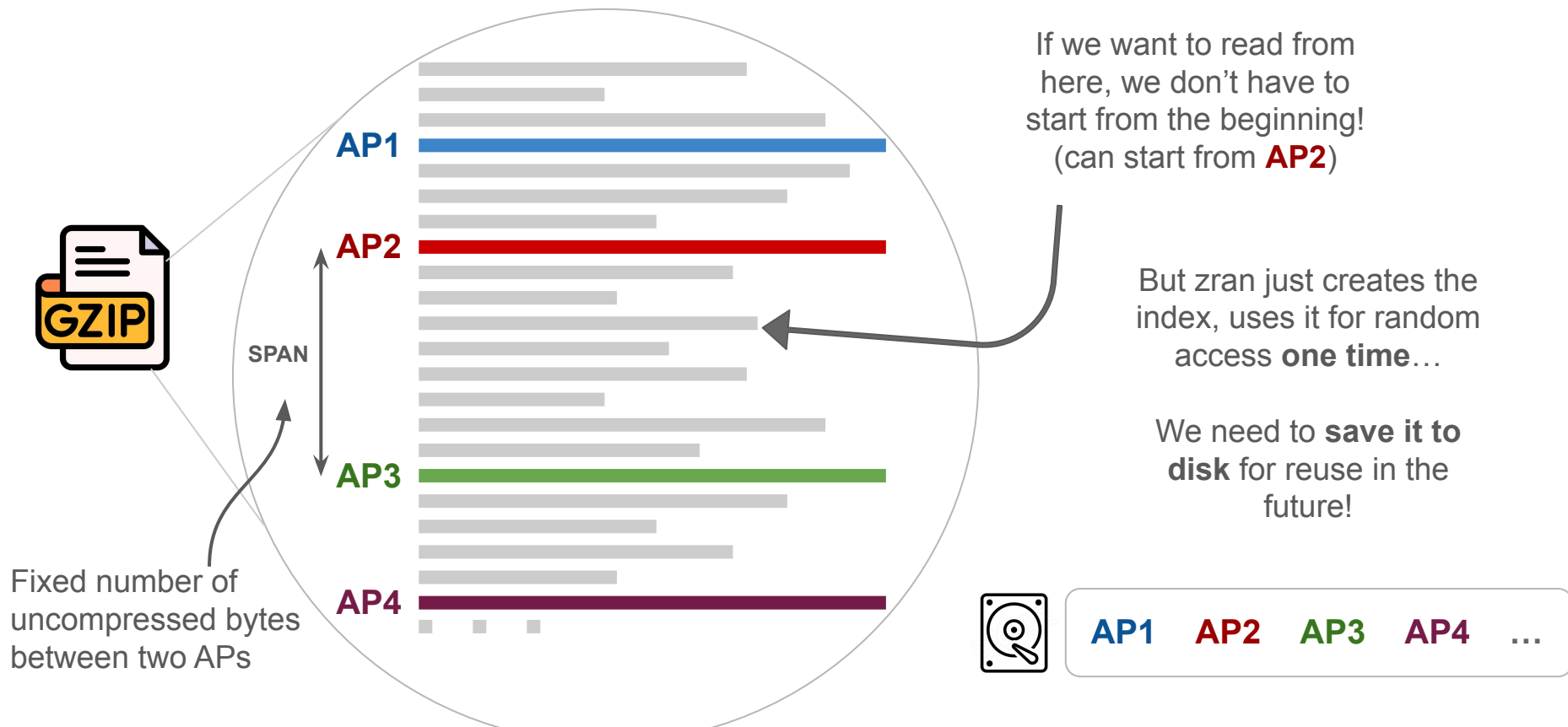
Context window for mLZ77.

Replace with (offset, length) pair from preceding context if possible. Max context size=32KB.

Huffman coding after mLZ77 compression. Reset per block.

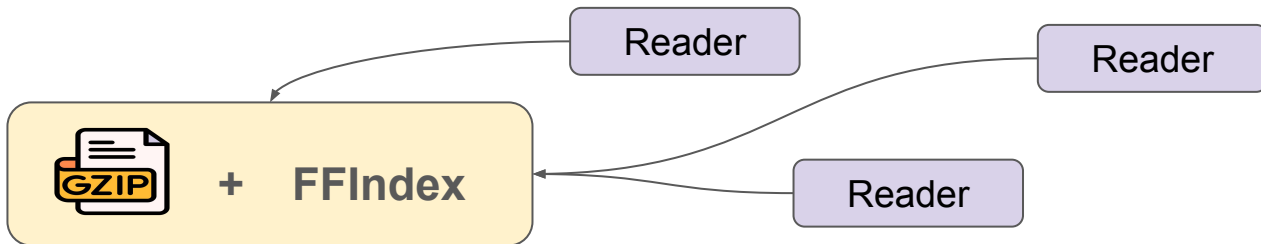
- But is it impossible? NO!

zran.c: [\[3\]](#) Example in zlib that creates “access points” for fast random access to gzip



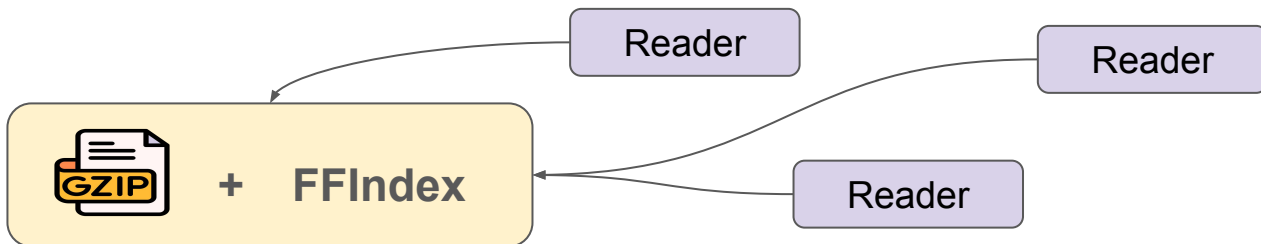
# Our approach

- **Index Building:**
  - Modify zran's index specifically for FASTQ files.
  - Save these access points and record boundaries to disk to supplement the original GZIP.
- **Parallel Parser Implementation:**
  - Implement producer-consumer model for parsing single FAST-Q files in parallel
  - Use saved index for parallel reading, using zran's random access algorithm and information on record boundaries.
  - Benchmark speedup and memory tradeoffs for various SPANs and other parameters.



# Our approach

- **Index Building:**
  - Modify zran's index specifically for FASTQ files.
  - Save these access points and record boundaries to disk to supplement the original GZIP.
- **Parallel Parser Implementation:**
  - Implement producer-consumer model for parsing single FAST-Q files in parallel
  - Use saved index for parallel reading, using zran's random access algorithm and information on record boundaries.
  - Benchmark speedup and memory tradeoffs for various SPANs and other parameters.





# Our approach: FFIndex



Read 1	Offset 1
Read 2	Offset 2
Read 3	Offset 3
...	...
Read N	Offset N



Zran Index

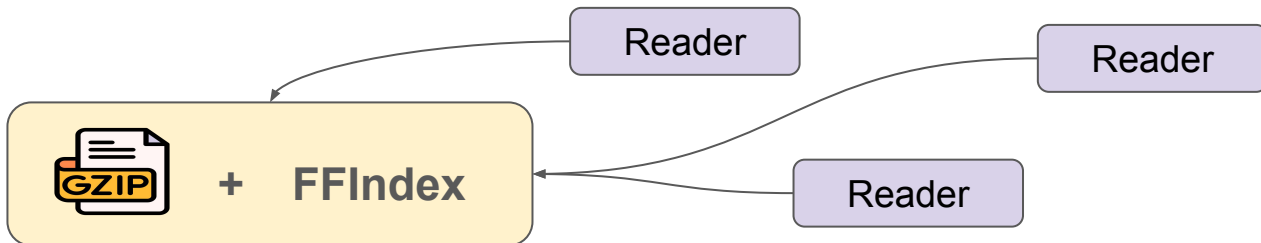
AP1 AP2 AP3 AP4 ...

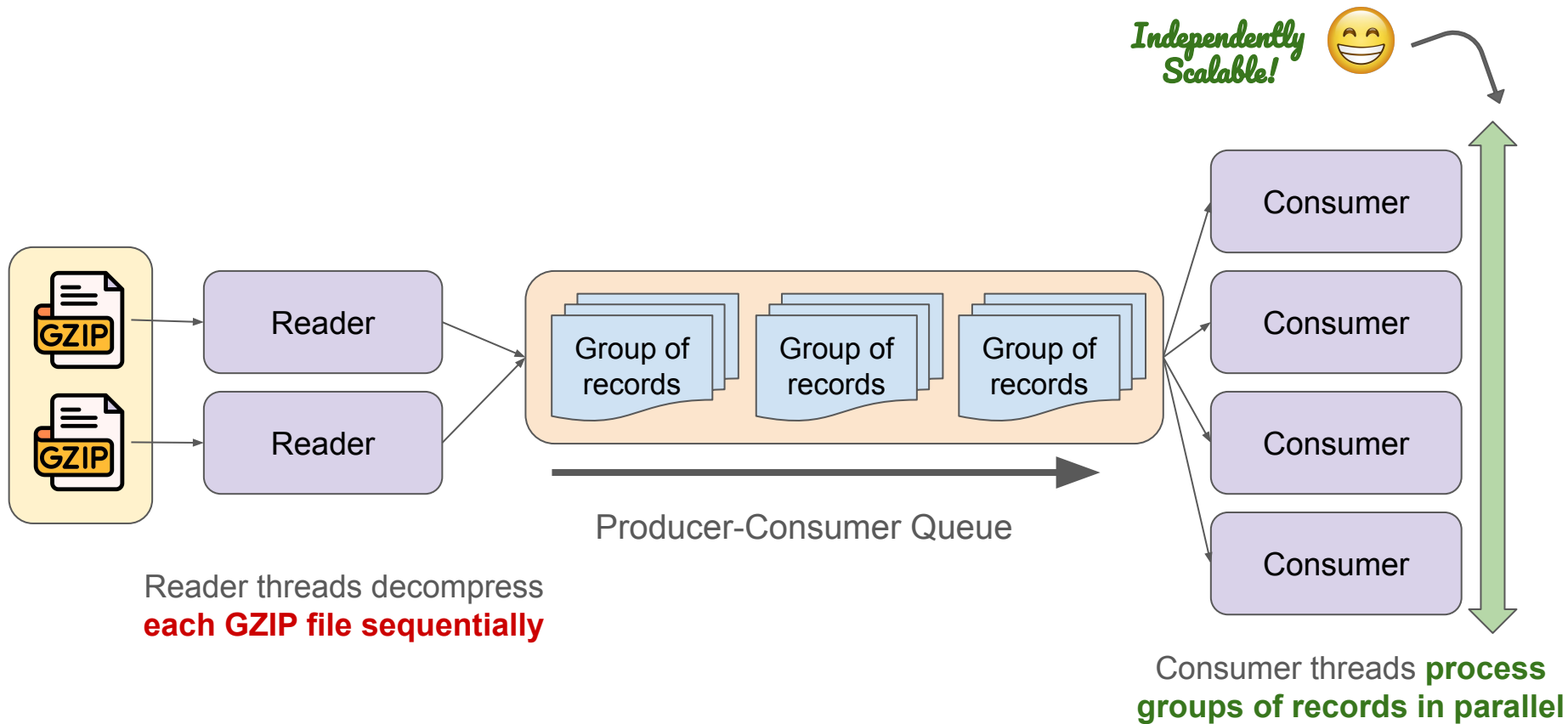
Binary  
Search

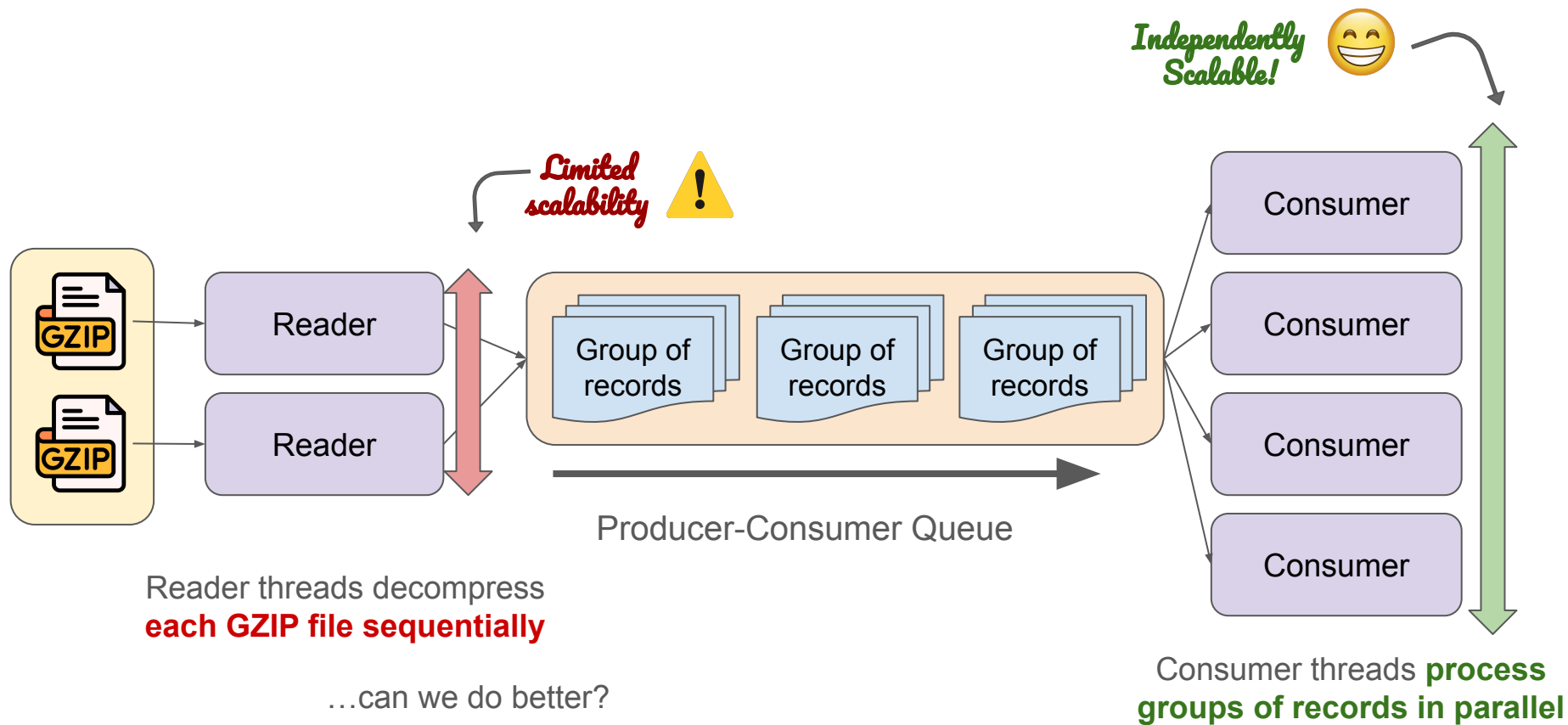


# Our approach

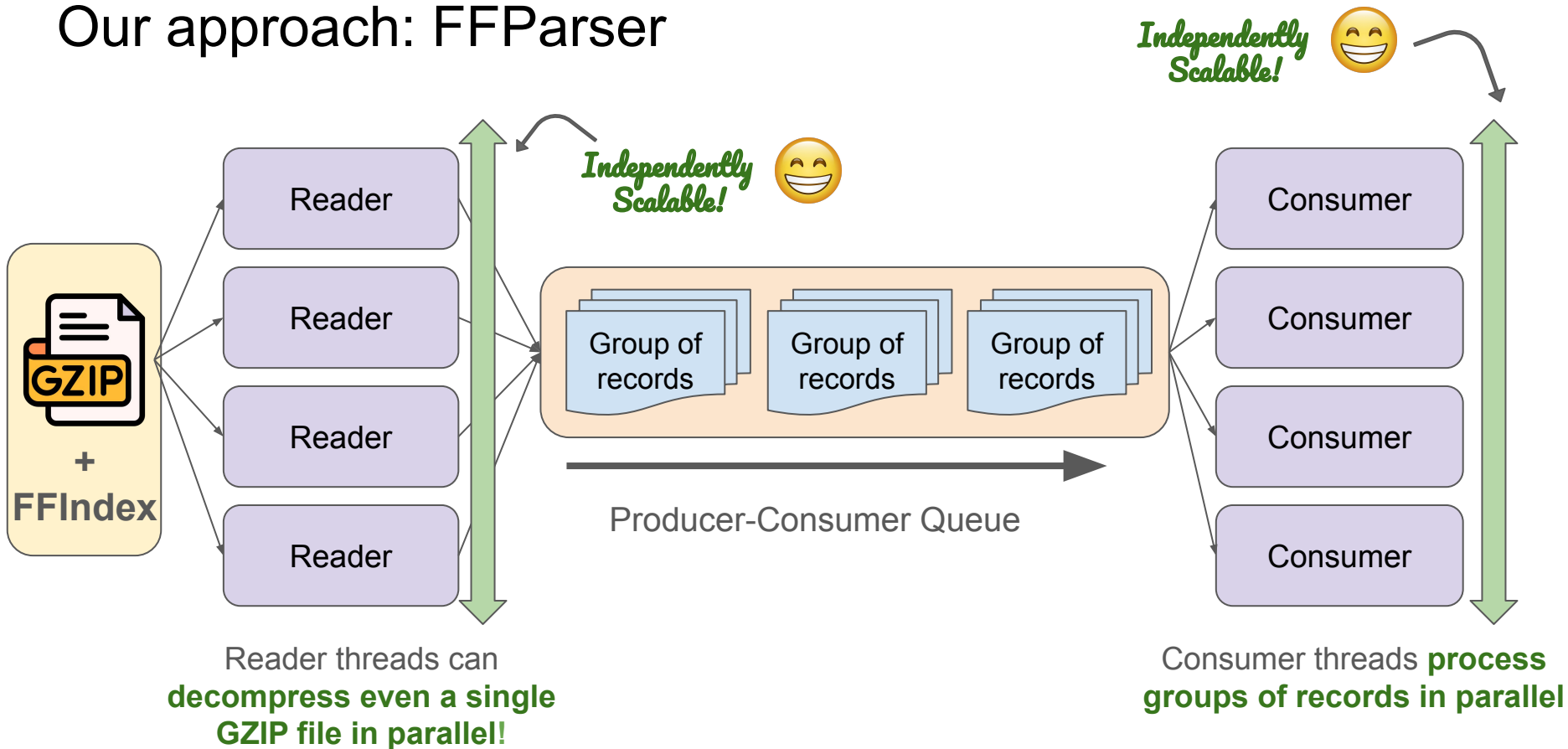
- **Index Building:**
  - Modify zran's index specifically for FASTQ files.
  - Save these access points and record boundaries to disk to supplement the original GZIP.
- **Parallel Parser Implementation:**
  - Implement producer-consumer model for parsing single FAST-Q files in parallel
  - Use saved index for parallel reading, using zran's random access algorithm and information on record boundaries.
  - Benchmark speedup and memory tradeoffs for various SPANs and other parameters.







# Our approach: FFParser



# Evaluation Plan

- **How lightweight is the index?**
  - Assess memory usage, storage needs for the index compared to the original files
- **Parallel Performance**
  - Evaluate strong scaling performance to see how much parallelism helps in reducing the time for parsing

# Preliminary Statistics on Index<sup>\*</sup>

SPAN = 1048576

	Compressed FASTQ file size	Compressed Index file size	Time to create the Index	Time to read the Index
Salmonella	1.4 MB	39 KB	56 ms	< 1ms
Ecoli	39 MB	5 MB	2819 ms	58 ms
Dataset C - Illumina <a href="#">(seqkit benchmark dataset)</a>	520 MB	38 MB	18846 ms	315 ms
Nematode	655 MB	48 MB	25816 ms	400 ms

# Preliminary Statistics on Index<sup>\*</sup>

SPAN = 1048576

	Compressed FASTQ file size	Compressed Index file size	Time to create the Index	Time to read the Index
Salmonella	1.4 MB	39 KB	56 ms	< 1ms
Ecoli	39 MB	5 MB	2819 ms	58 ms
Dataset C - Illumina <a href="#">(seqkit benchmark dataset)</a>	520 MB	38 MB	18846 ms	315 ms
Nematode	655 MB	48 MB	25816 ms	400 ms

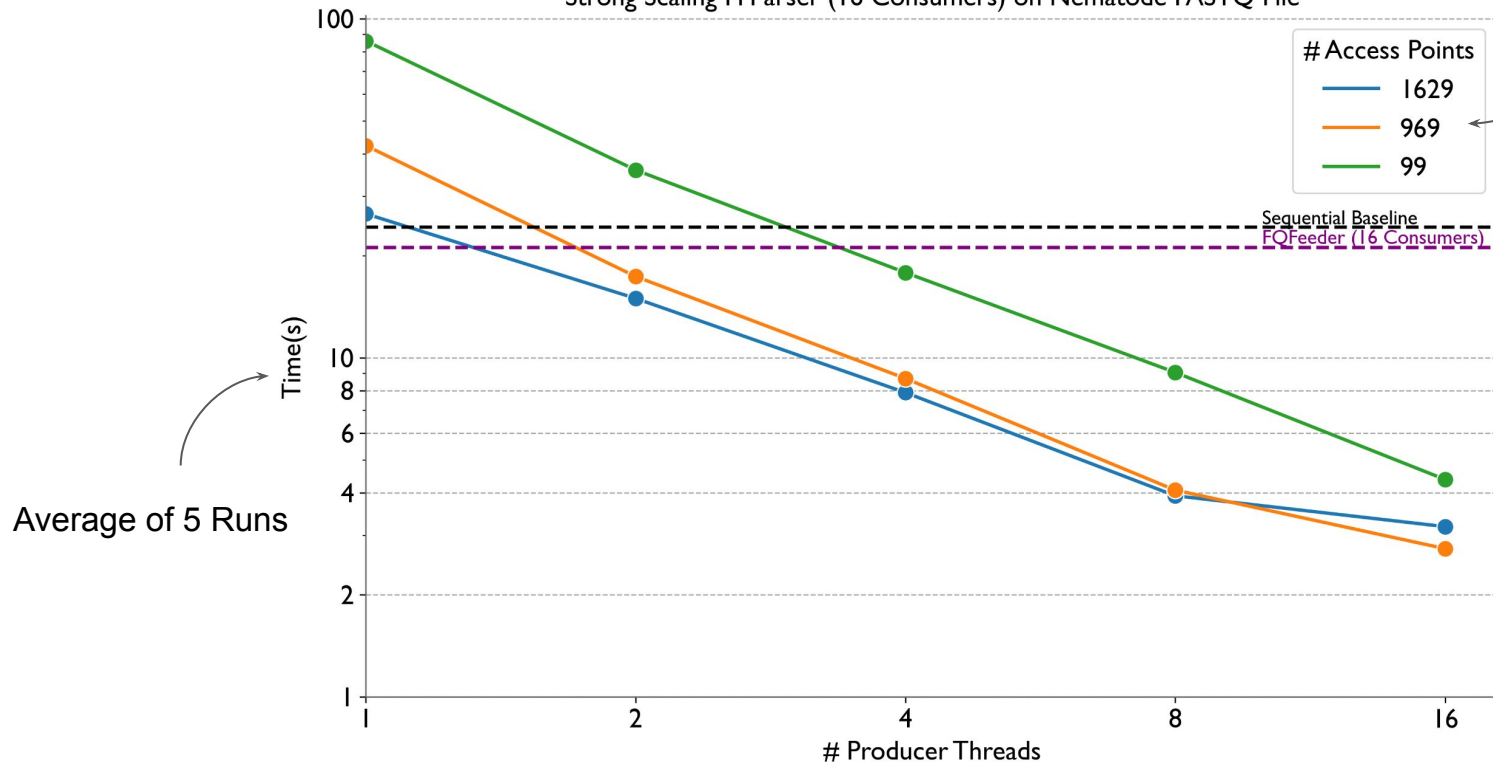
 8.8%!



# Preliminary Parallel Performance Results\*

# Access Points  $\propto 1/\text{SPAN}$

Strong Scaling FFParse (16 Consumers) on Nematode FASTQ File



\* run on Zaratán's Compute Node

# Further Evaluation and Optimizations

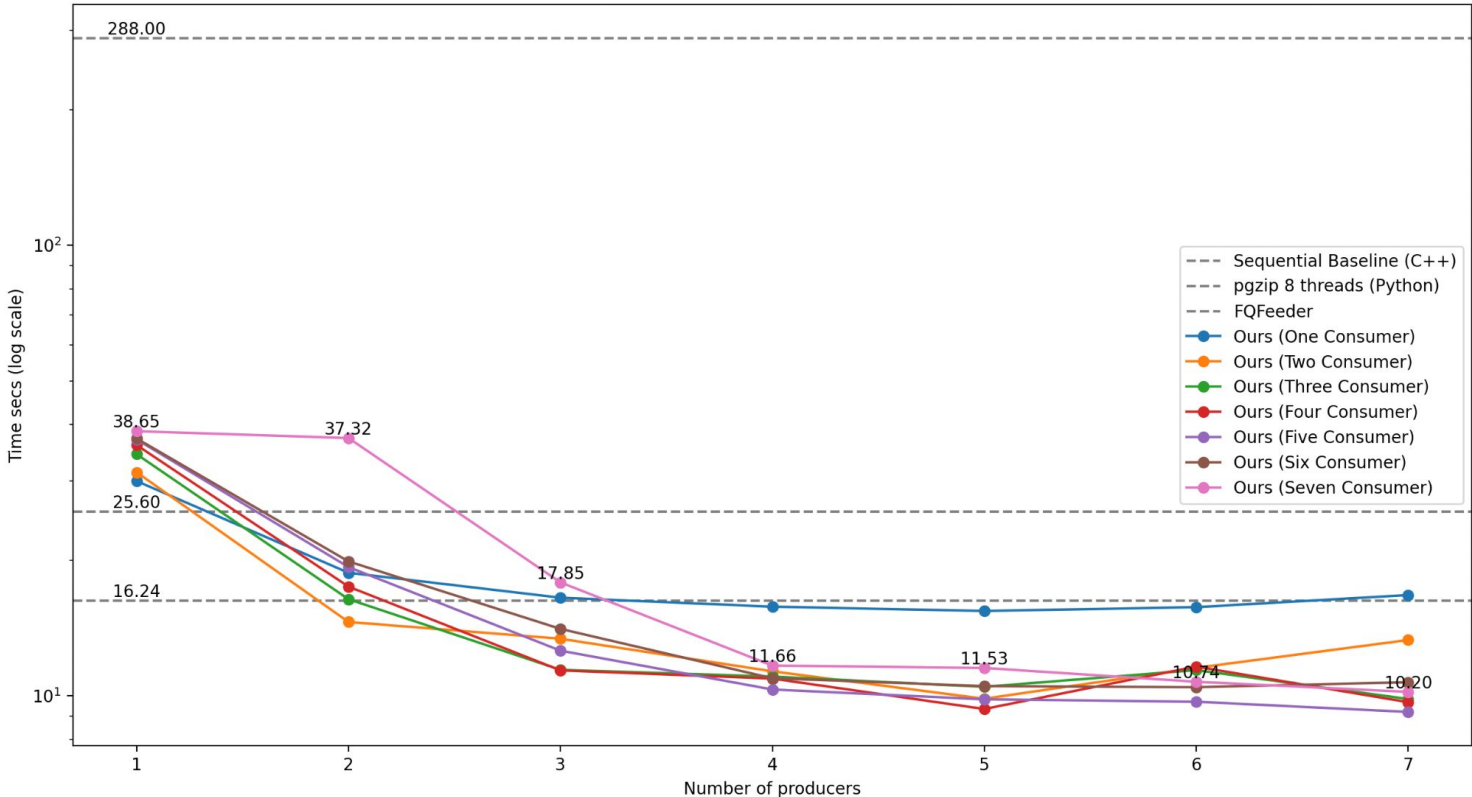
- **Datasets:**
  - Repeat experiments for bigger FAST-Q files
- **Component-wise Breakdown:**
  - Analyze which parts of the index are most space-consuming and target optimizations
  - Analyse what part of the codes can be optimised via Trace Analysis to identify bottlenecks and improving scaling efficiency
- **Optimizations Possible:**
  - Currently we store uncompressed byte offset for all FASTQ records. Probably only need to store for blocks of N records

Questions?

# Evaluation Plan

- **Lightweight Index:**
  - Assess memory usage, storage needs for the index.
- **Parallel Performance:**
  - Evaluate concurrency, throughput, and resource efficiency.
- **Datasets:**
  - Try with bigger dataset.
- **Component Breakdown:**
  - Analyze which parts of the index are most space-consuming and target optimizations.
- **Trace Analysis:**
  - Identify bottlenecks to enhance parsing speed and efficiency.

## Nematode



# Nematode

