

Exact Alignment with FM-Index on the Intel Xeon Phi Knights Landing Processor

Jose M. Herruzo¹ Sonia González-Navarro¹ Pablo Ibáñez²
Víctor Viñals² Jesús Alastruey-Benedé² Oscar Plata¹

¹Departamento de Arquitectura de Computadores
Universidad de Málaga

²Grupo de Arquitectura de Computadores
Universidad de Zaragoza

Accelerator Architecture in Computational Biology and
Bioinformatics, 2018

Motivation

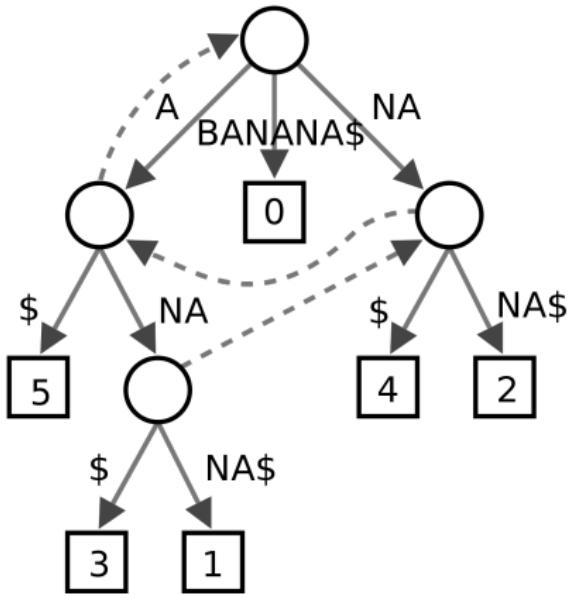
- Genomic Sequencing
 - New Sequencing Technologies

- Sequence alignment



Motivation

Indices



• Suffix Tree

• Hash Tables

• FM-Index

Motivation

BWT & FM-Index

- Burrows Wheeler Transform (BWT)

- FM-Index
 - Bowtie
 - BWA
 - SOAP2

banana\$

0	banana\$
1	anana\$b
2	nana\$ba
3	ana\$ban
4	na\$bana
5	a\$banan
6	\$banana

Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

FM-Index

- Data structures
 - C array
 - Occ table
- Random memory accesses
- Memory bound algorithm
- Backward Search algorithm

Algorithm: Backward Search Based on FM-index

Input: FM-index of T text (C & Occ), Q query, $n:|T|$, $p:|Q|$

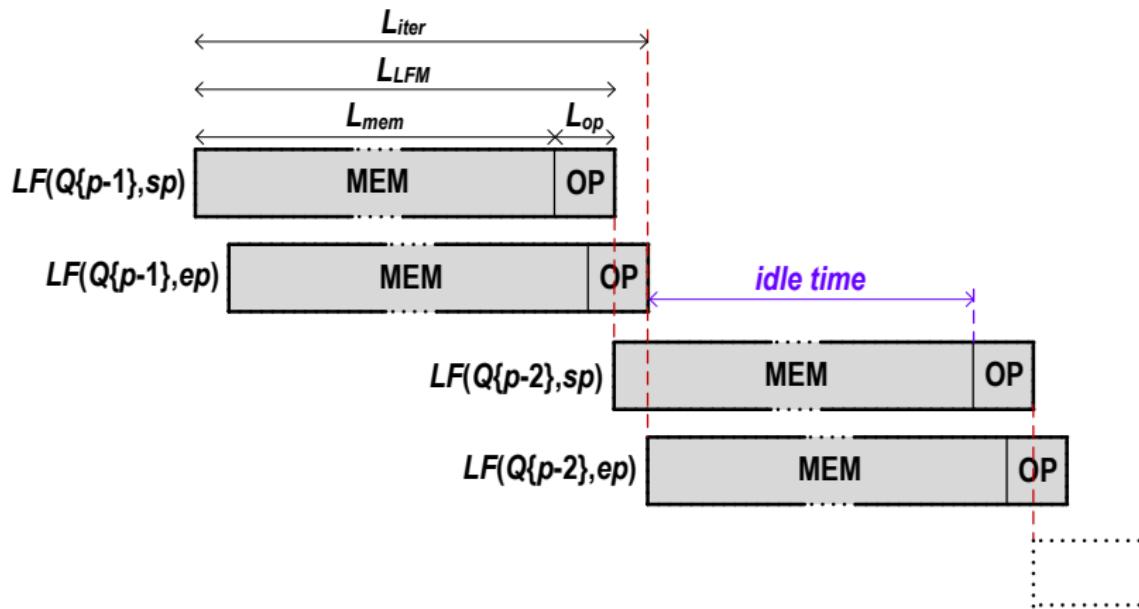
Ouput: (sp,ep) : Interval pointers of Q in T

begin

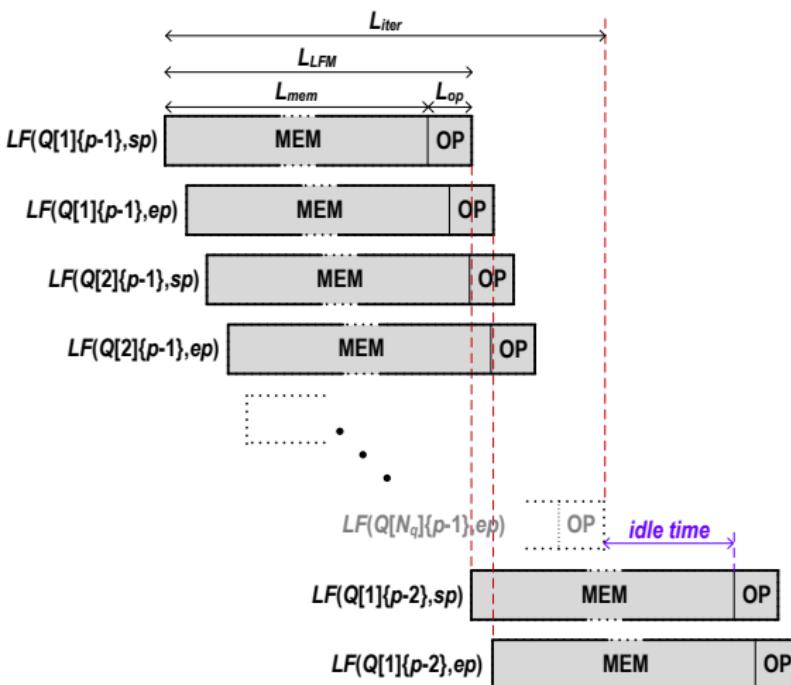
```
1:  $sp = C[Q[p]]$ 
2:  $ep = C[Q[p]+1]$ 
3: for  $i$  from  $p-1$  to  $1$  step  $-1$ 
4:    $sp = LF(Q[i],sp)$ 
5:    $ep = LF(Q[i],ep)$ 
6: end for
7: return  $(sp+1,ep)$ 
end
```

2 LFop-chains

FM-Index



FM-Index



Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

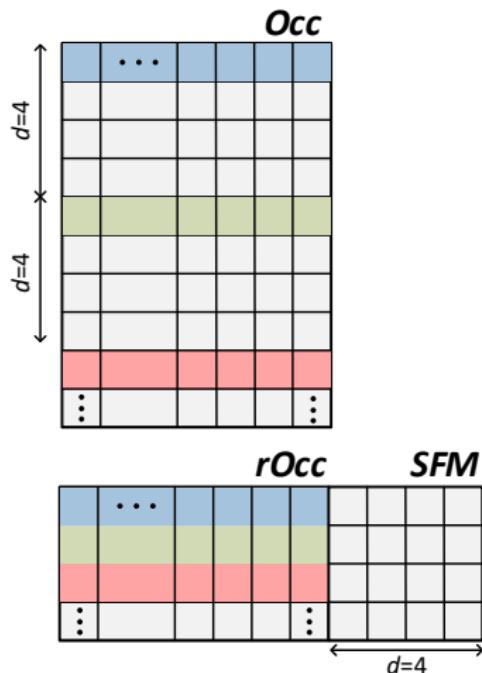
2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

Sampled FM-Index



Reduce Memory Footprint
 ↓
 Increase computing
 requirements

From Ferragina, P. and Manzini, G.: "Opportunistic Data Structures with Applications" (2000)

Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

3 Results

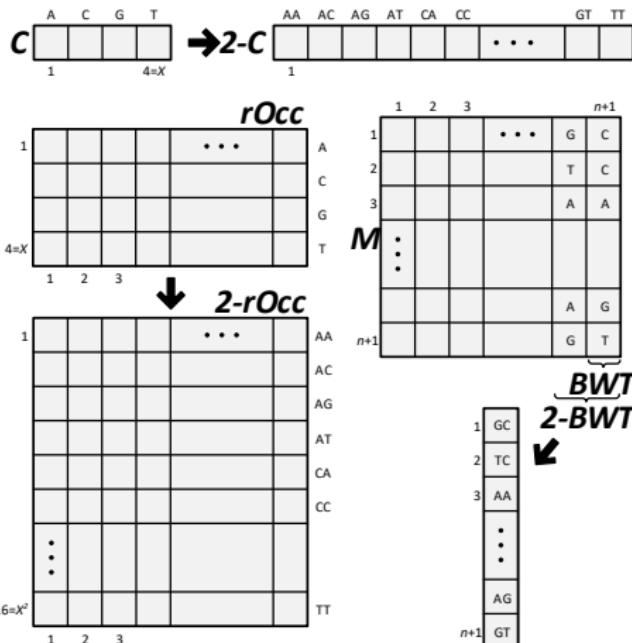
- RANDOM Benchmark
- Throughput results

4 Conclusions

K-Step FM-Index

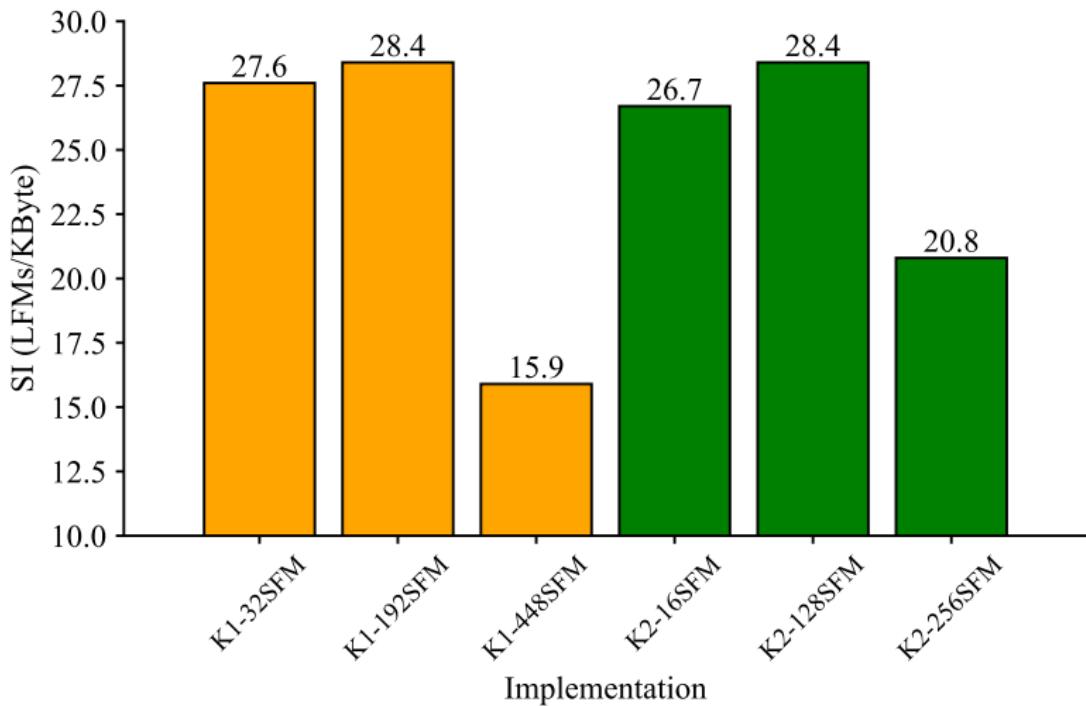
Searching several symbols per iteration

- Increased memory footprint
- Reduced number of LF operations



From Chacon et al. (2015)

FM-Index Comparison



Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

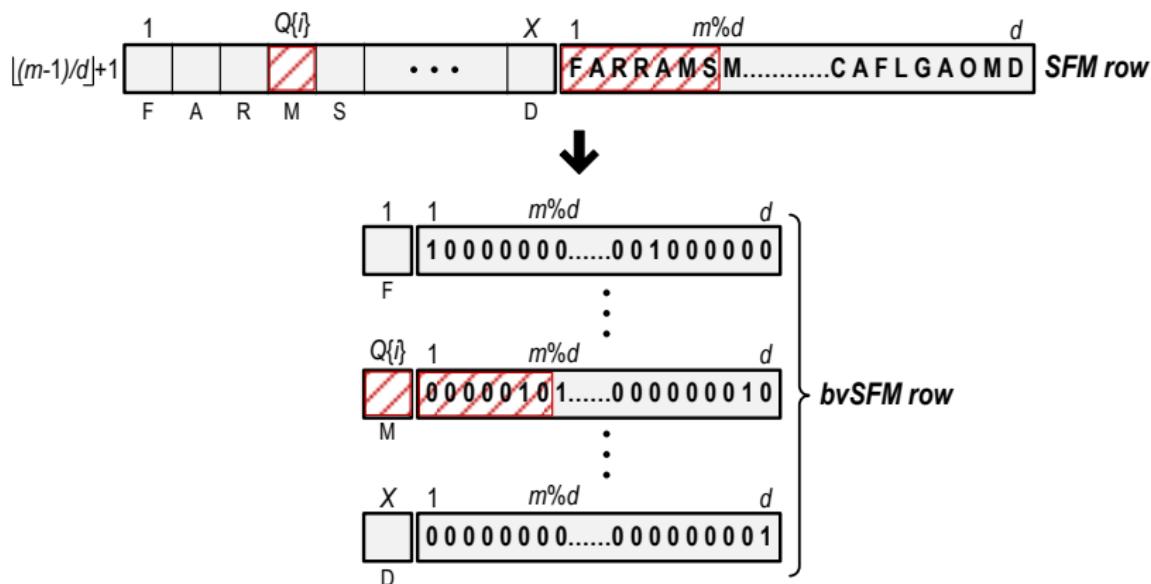
2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

Bit-Vector FM-Index



Bit-Vector FM-Index

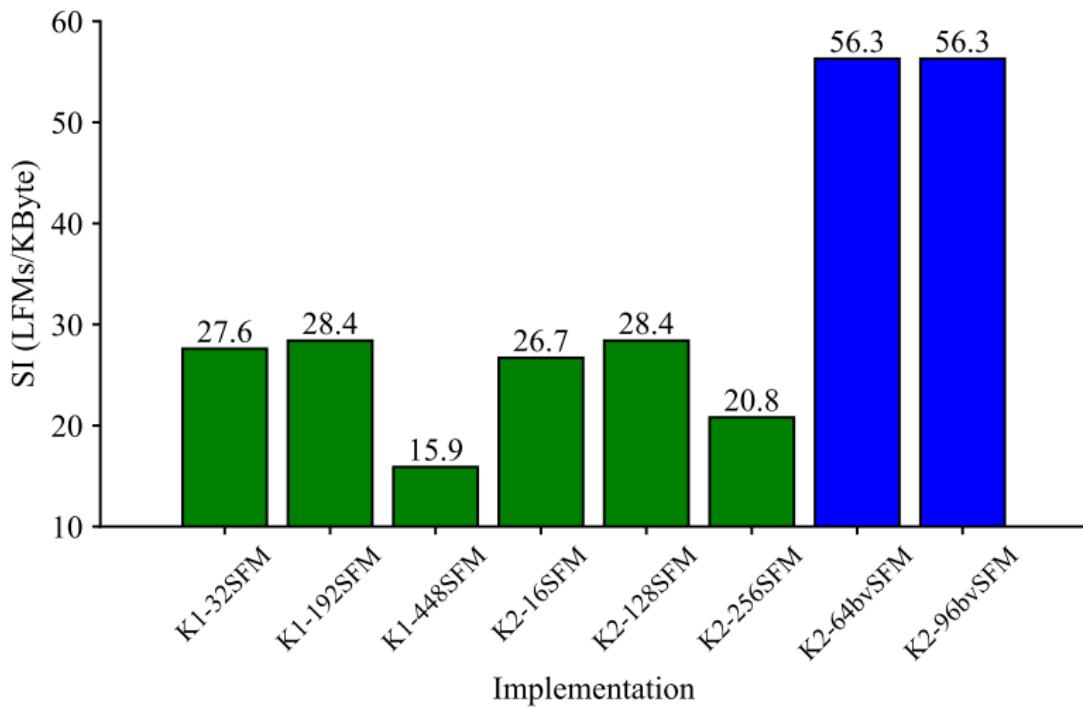
Advantages

- Reduced data movement
- Reduced computing requirements

Disadvantages

- Increased memory footprint

FM-Index Comparison



Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

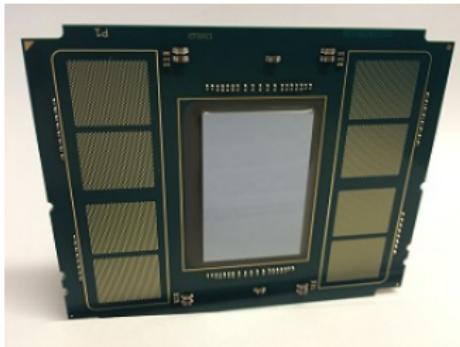
3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

Hardware Resources

Xeon Phi 7210 (KNL)	Xeon E5-2630V4 (Broadwell)
64 cores @ 1.3 GHz	10 cores @ 2.2 GHz
4 threads per core	2 threads per core
400 GB/s (MCDRAM)	68 GB/s (DDR4)



Intel Xeon Phi

- AVX 512 Vectorial Processing Units
- High Bandwidth Memory (HBM)
 - Cache / Hybrid / Flat

Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

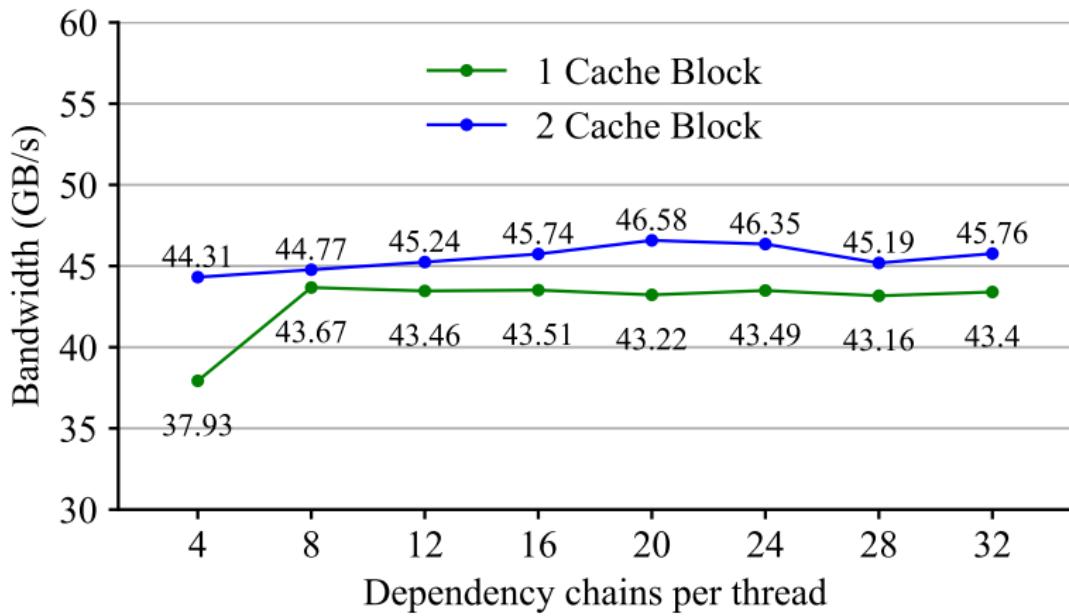
3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

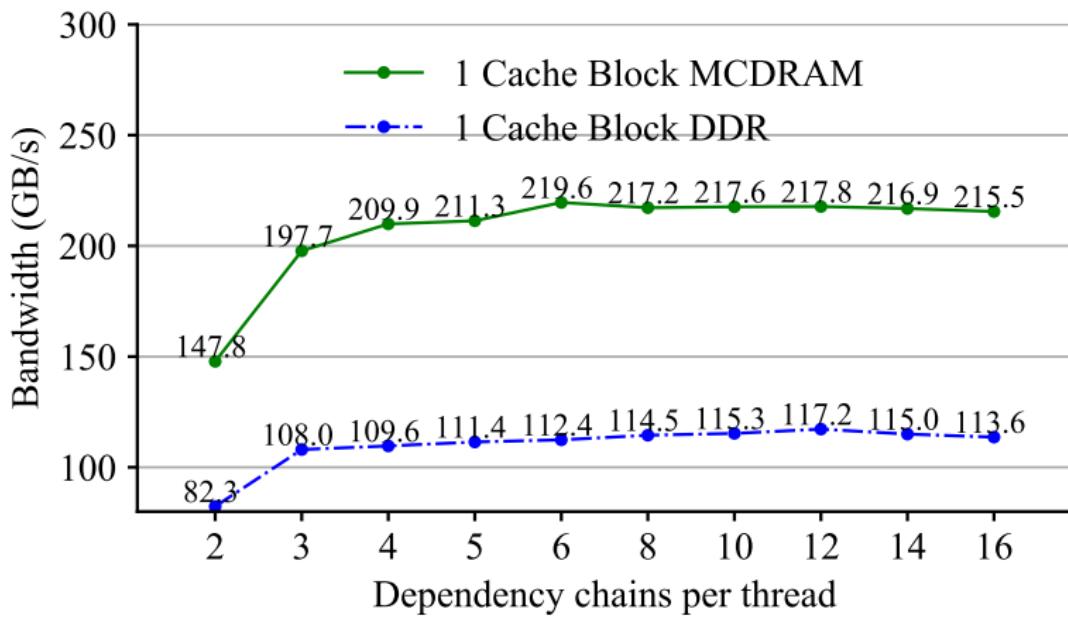
Results

Xeon Broadwell RANDOM Benchmark



Results

Xeon Phi KNL RANDOM Benchmark



Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

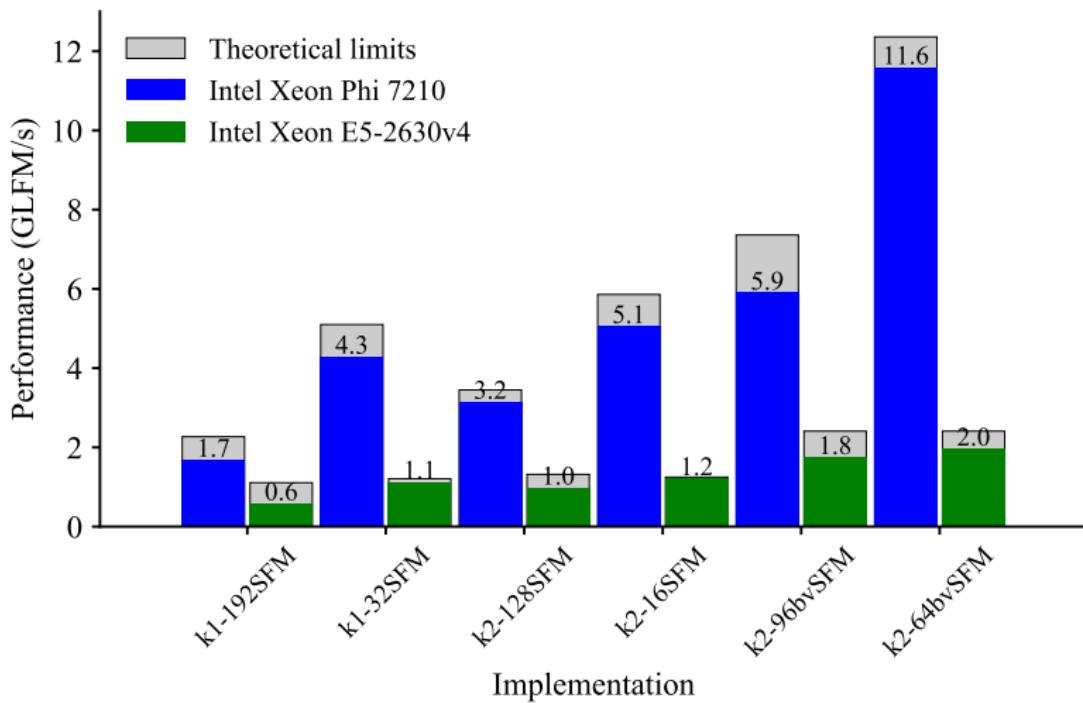
4 Conclusions

Results

Experimental Setup

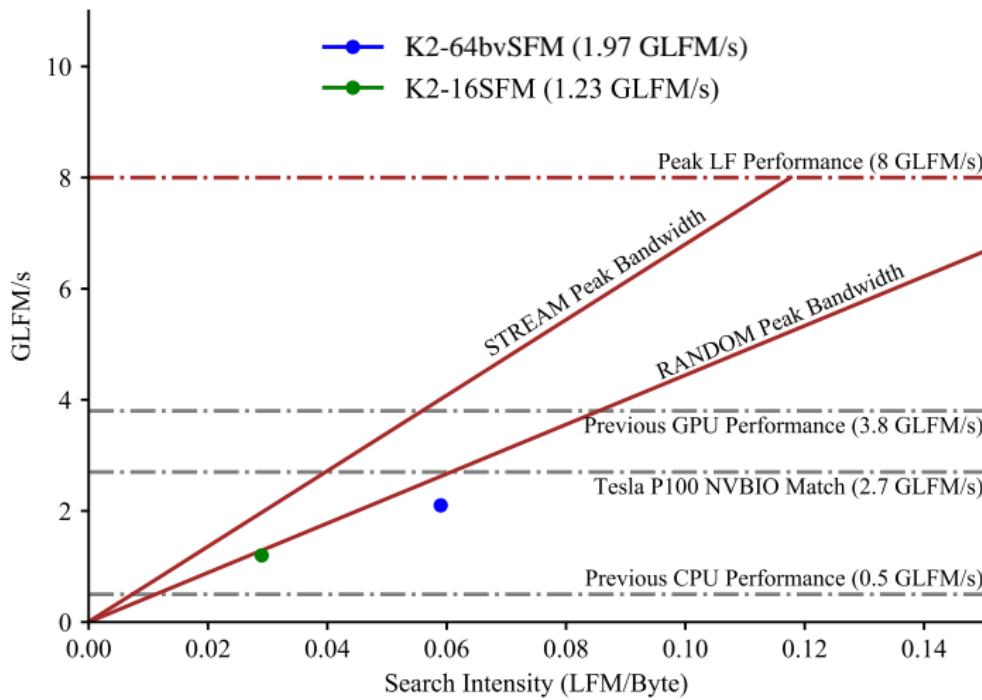
- Human Genome (Around 3GBases).
- 20 million input queries generated by Mason.
- 200 symbols per Sequence.
- Measurements started after loading the sequences into main memory.
- LFM/s as performance evaluation metric

Results



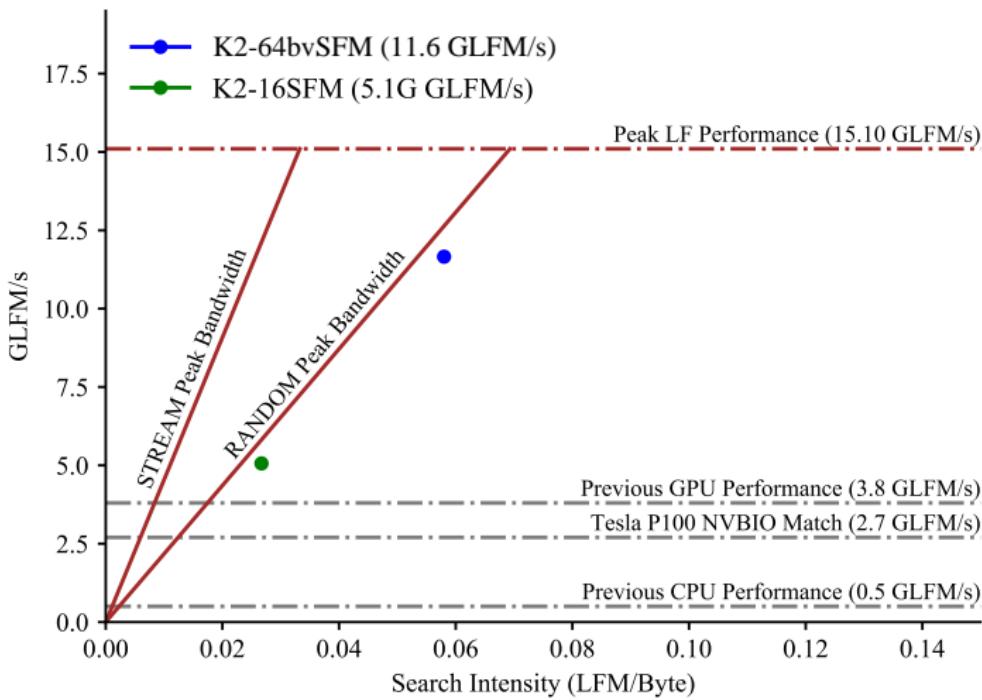
Results

Roofline (Broadwell)



Results

Roofline (Knights Landing)



Outline

1 FM-Index Search Algorithm

- Sampled FM-Index
- K-Step FM-Index
- Bit-Vector FM-Index

2 Hardware resources

3 Results

- RANDOM Benchmark
- Throughput results

4 Conclusions

Conclusions

- Reduced data movement
- Efficiently used memory bandwidth
- Great performance improvement
 - Up to 11.7 G LFM/s
 - 3x faster than previous GPU versions

Exact Alignment with FM-Index on the Intel Xeon Phi Knights Landing Processor

Jose M. Herruzo¹ Sonia González-Navarro¹ Pablo Ibáñez²
Víctor Viñals² Jesús Alastruey-Benedé² Oscar Plata¹

¹Departamento de Arquitectura de Computadores
Universidad de Málaga

²Grupo de Arquitectura de Computadores
Universidad de Zaragoza

Accelerator Architecture in Computational Biology and
Bioinformatics, 2018