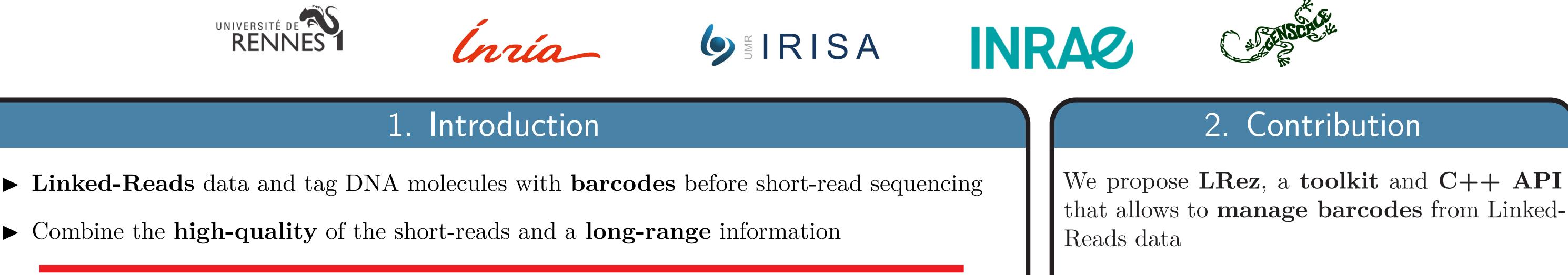
LRez: C++ API and toolkit for analyzing and managing Linked-Reads data

Pierre Morisse¹, Claire Lemaitre¹, Fabrice Legeai^{1,2}

¹Univ Rennes, Inria, CNRS, IRISA, F-35000 Rennes, France ²IGEPP, INRAE, Institut Agro, Univ Rennes, 35000, Rennes, France



► Various functionalities (indexing, querying, comparison, ...)

2. Contribution

► Allows to process **BAM** and FASTQ / gzipped FASTQ files

- 10x Genomics (2016) [1]

Multiple sequencing technologies:

50kbp

- stLFR (2019) [2]
- Haplotagging (2020) [3]
- TELL-Seq (2020) [4]
- ► Useful in a broad range of applications: assembly, phasing, scaffolding and SV calling
- ▶ No tool dedicated to Linked-Reads / barcodes management was previously available
 - 3. Functionalities
- ► Compatible with all currently available Linked-Reads technologies
- ► Easily usable in any external tool or pipeline to **improve perfor**mances

Toolkit command	Description	API module
compare	Compute the number of common barcodes between pairs of regions or between pairs of contig ends	BarcodesComparison
extract	Extract the barcodes from a given region of a BAM file	BarcodesExtraction
index bam	Index the BAM offsets or genomic positions of the barcodes contained in a BAM file	IndexManagementBam
index fastq	Index by barcode the offsets of the sequences contained in a FASTQ or gripped FASTQ file	${\it IndexManagementFastq}$
query bam	Query the index to retrieve alignments in a BAM file given a barcode or list of barcodes	${\it Alignments} {\it Retrieval}$
query fastq	Query the index to retrieve sequences in a FASTQ or gripped FASTQ file given a barcode or list of barcodes	ReadsRetrieval

4. Methods

► Index is represented as a map

– Keys: barcodes

5. Indexing performances

► **Datasets** from **all** Linked-Reads sequencing **technologies** and various species

- Values: list of occurrences positions	$- E. \ coli$ $- H. \ sapiens$					
ACGTAGCTGTAGTTAG: 0,3512,5340,,576948	- H. erato					
TTAGTTACGATTGAGG: $440,6598,9549,\ldots,657483$	Detect	\mathbf{BAM}	# Donadoa	Duntingo	\mathbf{RAM}	Disk
	Dataset	size (GB)	# Barcodes	Runtime	(MB)	(MB)
GGCCTAAAGCGATTCG: 842,4560,8756,,458765	TELL-Seq $(E. \ coli)$	1	$634,\!133$	1 min	293	340
► Query the index and browse the BAM /	10x Genomics (<i>H. sapiens</i>)	61	$609,\!058$	$52 \min$	$9,\!320$	$15,\!062$
FASTQ file to retrieve reads or align-	Haplotagging $(H. \ erato)$	70	$36,\!645,\!651$	1 h 09 min	10,751	$10,\!125$
ments	stLFR (<i>H. sapiens</i>)	206	38,779,362	3 h 06 min	26,769	34,256

6. Results: querying with LRez vs. samtools

- ▶ Querying experiments on the previous datasets from *E. coli*
- Comparison against a **naive methods** based on **samtools**
- ► Reported statistics from a **thousand queries**

	Overall runtime			Runtime per query		
Dataset	Samtools	LRez	LRez	Samtools	LRez	
		(index + query)	(query)	Samoois	LNez	
TELL-Seq $(E. \ coli)$	8 h 18 min	1 min	$11 \mathrm{sec}$	$30 \mathrm{sec}$	$4 \mathrm{ms}$	
10x Genomics (<i>H. sapiens</i>)	11.8 days	1 h 02 min	$10 \min$	$17 \min$	$290 \mathrm{\ ms}$	
Haplotagging $(H. erato)$	7.6 days	1 h 14 min	$5 \min$	11 min	$11 \mathrm{ms}$	
stLFR (H. sapiens)	41.6 days	3 h 14 min	$8 \min$	1 h	$15 \mathrm{\ ms}$	

7. Conclusion

- ► Novel and open-source toolkit and C++ API for processing Linked-Reads barcodes
- ► **Compatible** with **all** Linked-Reads technologies
- ► Various functionalities (indexing, querying, comparison, extraction)

- ▶ LRez is much faster, even when counting indexing time
- Runtime per query varies according to the number of reads / alignments
- LRez reaches a runtime of 15 ms per query on a 206 GB BAM file

Brock Medsker et al. Haplotyping germline and cancer genomes using high- throughput linked-read sequencing. Nature Biotechnology, 34(3):303-311, 2016. 1

- Ou Wang et al. Efficient and unique co-barcoding of second-generation sequencing reads from long dna molecules enabling cost effective and accurate sequencing, |2|haplotyping, and de novo assembly. Genome Research, 2019.
- Joana I. Meier et al. Haplotype tagging reveals parallel formation of hybrid races in two butterfly species. bioRxiv, pages 1–27, 2020. |3|
- Zhoutao Chen et al. Ultra-low input single tube linked-read library method enables short-read second-generation sequencing systems to generate highly accurate and economical long-range sequencing information routinely. Genome Research, 2020.

usable Easily external in projects (used in a SV calling tool and a gap-filling pipeline)

Efficient: Time-saving and scaling barcode processing



github.com/morispi/LRez Contact: pierre.morisse@inria.fr