# RLBWT Tricks

Nathaniel K. Brown[1]    Travis Gagie[1]    Massimiliano Rossi[2]

[1]Faculty of Computer Science
Dalhousie University
Halifax, NS, Canada

[2]Department of Computer and Information Science and Engineering
University of Florida
Gainesville, FL, USA

Data Compression Conference, March 2022

# Table of Contents

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction
Table Lookup
Mapping Runs
Table Lookup
Implementation
Experiments
Thanks

# Introduction

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

- **String Indexing**: Support sub-string queries on text
- **FM-Index**: basis for key tools in computational genomics
    - Short read aligners such as BWA and Bowtie
    - Application of Burrows-Wheeler Transform (BWT)
- **Computational Pan-Genomics**:
    - Want to index many genomes in reasonable space
    - *Solution*: Versions of FM-Index based on run-length compressed BWT (RLBWT)

- **Using Burrows-Wheeler Transform (BWT)**
  - Leverage last-to-first (LF) mapping

- **Pan-Genomic Indexes on run-length BWT (RLBWT)**
  - Conventionally, cannot compute LF steps in constant time

- **Nishimoto and Tabei's OptBWTR** (ICALP '21)
  - New, simple and constant-time implementation

We show experimentally that their approach can be made
practical for LF even without theoretical guarantees

# LF Permutation

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

**Mapping Runs**

Table Lookup

Implementation

Experiments

Thanks

$T = \text{GATTAGATACAT}$

| $L$ | $F$ |
|---|---|
| $T_1$ | $\$$ |
| $T_2$ | $A_1$ |
| $T_3$ | $A_2$ |
| $C_1$ | $A_3$ |
| $G_1$ | $A_4$ |
| $G_2$ | $A_5$ |
| $A_1$ | $C_1$ |
| $A_2$ | $G_1$ |
| $\$$ | $G_2$ |
| $A_3$ | $T_1$ |
| $A_4$ | $T_2$ |
| $T_4$ | $T_3$ |
| $A_5$ | $T_4$ |

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

# LF Permutation

$T = \text{GATTAGATACAT}$



*L*

$T_1$
$T_2$
$T_3$
$C_1$
$G_1$
$G_2$
$A_1$
$A_2$
$\$$
$A_3$
$A_4$
$T_4$
$A_5$

*F*

$\$$
$A_1$
$A_2$
$A_3$
$A_4$
$A_5$
$C_1$
$G_1$
$G_2$
$T_1$
$T_2$
$T_3$
$T_4$

# LF Runs

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

*L*

|   |       |
|---|-------|
|   | $T_1$ |
| 0 | $T_2$ |
|   | $T_3$ |
| 1 | $C_1$ |
| 2 | $G_1$ |
|   | $G_2$ |
|   | $A_1$ |
| 3 | $A_2$ |
| 4 | $\$$ |
| 5 | $A_3$ |
|   | $A_4$ |
| 6 | $T_4$ |
| 7 | $A_5$ |

*F*

|       |   |
|-------|---|
| $\$$ | 4 |
| $A_1$ |   |
| $A_2$ | 3 |
| $A_3$ |   |
| $A_4$ | 5 |
| $A_5$ | 7 |
| $C_1$ | 1 |
| $G_1$ |   |
| $G_2$ | 2 |
| $T_1$ |   |
| $T_2$ | 0 |
| $T_3$ |   |
| $T_4$ | 6 |

# LF Runs

RLBWT
Tricks

Brown, Gagie,
Rossi

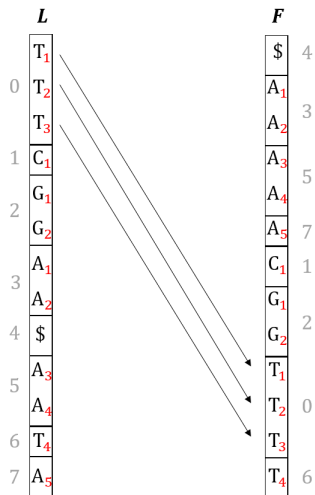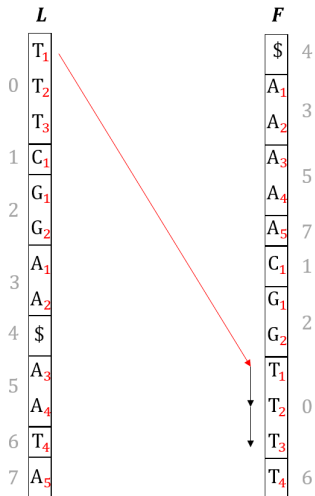Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

# Any Permutation

RLBWT
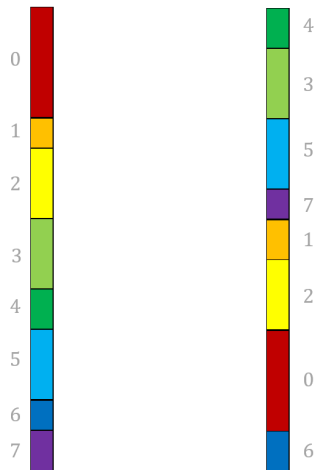Tricks

Brown, Gagie,
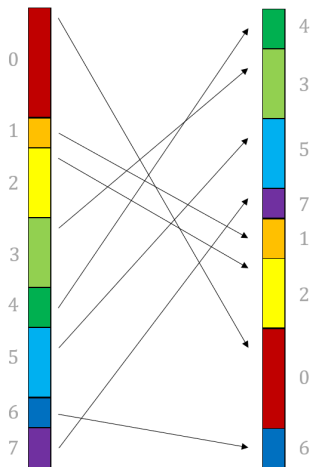Rossi

Introduction

Table Lookup

**Mapping Runs**

Table Lookup

Implementation

Experiments

Thanks

# Any Permutation

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

# LF Table

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

**Table Lookup**

Implementation

Experiments

Thanks

| L | | F | | run | offset |
|---|---|---|---|-----|--------|
| $T_1$ | | $\$$ | | 0 | 0 |
| $T_2$ | | $A_1$ | | 0 | 1 |
| $T_3$ | | $A_2$ | | 0 | 2 |
| $C_1$ | | $A_3$ | | 1 | 0 |
| $G_1$ | | $A_4$ | | 2 | 0 |
| $G_2$ | | $A_5$ | | 2 | 1 |
| $A_1$ | | $C_1$ | | 3 | 0 |
| $A_2$ | | $G_1$ | | 3 | 1 |
| $\$$ | | $G_2$ | | 4 | 0 |
| $A_3$ | | $T_1$ | | 5 | 0 |
| $A_4$ | | $T_2$ | | 5 | 1 |
| $T_4$ | | $T_3$ | | 6 | 0 |
| $A_5$ | | $T_4$ | | 7 | 0 |

| | character | length | destination | offset |
|---|-----------|--------|-------------|--------|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |
| 4 | | | | |
| 5 | | | | |
| 6 | | | | |
| 7 | | | | |

# LF Table

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

| | L | | F | run | offset |
|---|---|---|---|---|---|
| | $T_1$ | | $ | 0 | 0 |
| | $T_2$ | | $A_1$ | 0 | 1 |
| | $T_3$ | | $A_2$ | 0 | 2 |
| | $C_1$ | | $A_3$ | 1 | 0 |
| | $G_1$ | | $A_4$ | 2 | 0 |
| | $G_2$ | | $A_5$ | 2 | 1 |
| | $A_1$ | | $C_1$ | 3 | 0 |
| | $A_2$ | | $G_1$ | 3 | 1 |
| | $ | | $G_2$ | 4 | 0 |
| | $A_3$ | | $T_1$ | 5 | 0 |
| | $A_4$ | | $T_2$ | 5 | 1 |
| | $T_4$ | | $T_3$ | 6 | 0 |
| | $A_5$ | | $T_4$ | 7 | 0 |

| | character | length | destination | offset |
|---|---|---|---|---|
| 0 | T | 3 | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |
| 4 | | | | |
| 5 | | | | |
| 6 | | | | |
| 7 | | | | |

# LF Table

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

# LF Table

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

| | $L$ | | $F$ | run | offset |
|---|---|---|---|---|---|
| | $T_1$ | | $\$$ | 0 | 0 |
| | $T_2$ | | $A_1$ | 0 | 1 |
| | $T_3$ | | $A_2$ | 0 | 2 |
| | $C_1$ | | $A_3$ | 1 | 0 |
| | $G_1$ | | $A_4$ | 2 | 0 |
| | $G_2$ | | $A_5$ | 2 | 1 |
| | $A_1$ | | $C_1$ | 3 | 0 |
| | $A_2$ | | $G_1$ | 3 | 1 |
| | $\$$ | | $G_2$ | 4 | 0 |
| | $A_3$ | | $T_1$ | 5 | 0 |
| | $A_4$ | | $T_2$ | 5 | 1 |
| | $T_4$ | | $T_3$ | 6 | 0 |
| | $A_5$ | | $T_4$ | 7 | 0 |

| | character | length | destination | offset |
|---|---|---|---|---|
| 0 | T | 3 | 5 | 0 |
| 1 | C | 1 | 3 | 0 |
| 2 | G | 2 | 3 | 1 |
| 3 | A | 2 | 0 | 1 |
| 4 | $ | 1 | 0 | 0 |
| 5 | A | 2 | 1 | 0 |
| 6 | T | 1 | 7 | 0 |
| 7 | A | 1 | 2 | 1 |

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

# Boundaries in Practice

- Nishimoto and Tabei limit crossings (additional space)
- 98% cross less than 5 boundaries



Runs Scanned for Sampled LF Steps of Chromosome-19

# Compression

RLBWT
Tricks
Brown, Gagie,
Rossi

Introduction
Table Lookup
Mapping Runs
Table Lookup
Implementation
Experiments
Thanks

- Preliminary results vs. conventional approach:
    - LF steps $\approx$ 6 times faster
    - Table $\approx$ 14 times larger
- We devise a compression scheme specific to LF
    - To perform column-wise compression, partition into blocks to mitigate locality concerns
    - For alphabet size $\sigma$, LF mapping of run-heads forms $\sigma$ non-decreasing subsequences

- Randomly sample 10000 patterns of length 100 and compute count queries
- Query against chromosome-19 genomes of 128, 256, 512 and 1000 copies
- **Data Structures**:
  - **sparse bv**: The sparse bitvector component of $r$-index
  - **wt_fbb** Fixed block boosting wavelet tree
  - **table** Our implementation of LF using Nishimoto and Tabei's approach
  - **RLCSA** BWT component of run-length encoded compressed suffix array

RLBWT
Tricks

Brown, Gagie,
Rossi

Introduction

Table Lookup

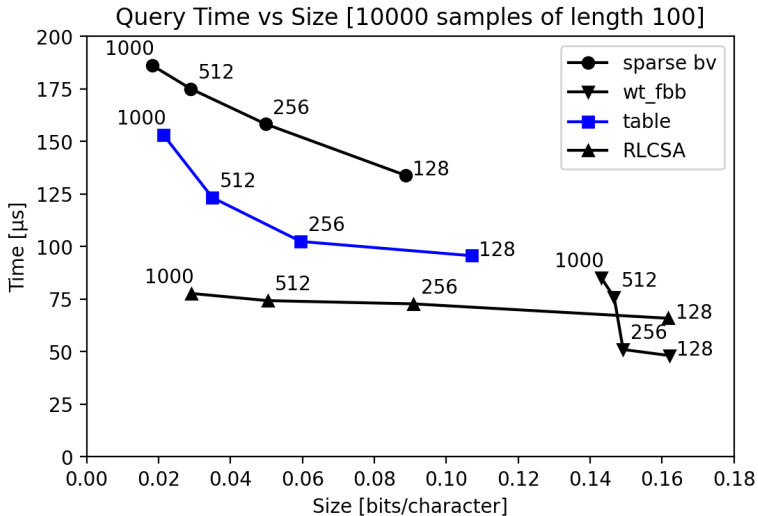Mapping Runs

Table Lookup

Implementation

Experiments

Thanks

# Results



Query Time vs Size [10000 samples of length 100]

## Thanks

- Email: nathaniel.brown@dal.ca
- Full Paper: https://arxiv.org/abs/2112.04271