

Augmented Thresholds for MONI

César Martínez-Guardiola¹ **Nathaniel K. Brown**²
Fernando Silva-Coira¹ Dominik Köppl³ Travis Gagie²
Susana Ladra¹

¹CITIC
Universidade da Coruña
A Coruña, Spain

²FCS
Dalhousie University
Halifax, Canada

³M&D Data Science Center
TMDU
Tokyo, Japan

Data Compression Conference, March 2023

Introduction

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **Computational Pan-Genomics:**
 - Want to index many genomes in reasonable space
 - Pan-genome graphs one method
 - *Solution:* Versions of FM-Index based on run-length compressed BWT (RLBWT)

- **r-index:** efficient exact matching in runs-bounded space
 - Computes LF^F steps, samples SA at run boundaries

Maximal Exact Matches

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **Approximate Matching:** Exact matching is not ideal for most use cases
 - Instead, want to support approximate matching
- **Maximal Exact Matches (MEMs):** The longest matching sub strings of a pattern P with a text T
 - $P[i..j]$ occurs in T
 - Neither $P[i - 1..j]$ or $P[i..j + 1]$ occurs in T
 - Supports approximate pattern matching in tools like BWA-MEM

Pan-genomics MEMs

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented Thresholds

Experiments

Thanks

- **MONI**: Rossi et al.'s r -index variant supporting MEM-finding
 - Computes MEMs using matching statistics of text
 - Two passes over the pattern, stores 'threshold' values at run boundaries
 - Can be made online without thresholds by using Longest Common Extension (LCE) queries
 - Online one-pass variant named **PHONI** (Boucher et al.)

- Online can be used in targeted approaches
 - Ahmed et al.'s **SPUMONI** for metagenomic classification using nanopore sequencing

MEMs

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented Thresholds

Experiments

Thanks

We show how to use both thresholds and LCE queries to compute MEMs online quickly by storing additional LCE information with thresholds.

Suffix Array

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **Suffix Array (SA):** $SA[i]$ is the starting position of the lexicographically i th suffix in a text $T[0..n-1]$
 - MONI samples at start/end of BWT-runs
- $BWT[i]$ is character preceding $SA[i]$
- $LF(i)$ is position of $SA[i-1]$ in SA

Longest Common Prefixes

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **LCP Array:** $LCP[i]$ stores the LCP between suffixes at $SA[i]$ and $SA[i - 1]$
- **Longest Common Extension:** $LCE(i, j)$ returns the LCP between $T[i..n - 1]$ and $T[j..n - 1]$
 - Given SA samples: $LCE(SA[i], SA[j])$
 - Range-minimum-query (RMQ) over the LCP array

Thresholds

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

	k	SA[k]	BWT[k]	T [SA[k .. n]	LCP[k]
	\vdots	\vdots	\vdots	\vdots	\vdots
	1234	8765	A	GAGACATCA...	-
$e_1 =$	1235	1519	A	GATACATTA...	-
	1236	5450	C	GATAGATTA...	4
$j =$	1237	1004	G	GATATAGAA...	4
	1238	4242	G	GATCCAATA...	3
$t =$	1239	3110	G	GATTACATA...	3
	1240	1102	T	GATTACTTA...	6
	1241	1978	T	GATTAGATA...	5
$s_2 =$	1242	2505	A	GATTATCAT...	5
	1243	2022	A	GATTATGAA...	-
	\vdots	\vdots	\vdots	\vdots	\vdots

Matching Statistics

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **Matching Statistics (MS)**: For pattern $P[0..m - 1]$ and text $T[0..n - 1]$, $MS[0..m - 1]$ defined as pair
 - $MS[i].pos$ is the starting text position of the longest prefix of $P[i..m - 1]$ that occurs in T
 - $MS[i].len$ is the length of that prefix

MONI Process

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **Matching Statistics for $P[0..m-1]$**
 - $j = \text{BWT.select}_{P[m-1]}(1)$
 - $\text{MS}[m-1].\text{pos} = \text{SA}[j]$
 - $j = \text{LF}(j)$, continue with $i = m-2, \dots, 0$
 - **Case 1:** $\text{BWT}[j] = P[i]$, $\text{MS}[i].\text{pos} = \text{MS}[i+1].\text{pos} - 1$
 - **Case 2:** $\text{BWT}[j] \neq P[i]$
 - $\text{BWT}[s_1..e_1] = \text{BWT}[s_2..e_2] = P[i]$ where $e_1 < j < s_2$
 - Threshold t ; if $j < t$ take $\text{SA}[e_1]$, if $j \geq t$ take $\text{SA}[s_2]$
 - Set j , $\text{MS}[i].\text{pos}$ corresponding with choice
 - LF step and continue
- **Second left-to-right pass to recover lengths**
 - Using compressed straight line program (SLP) of T supporting random access

Threshold LCE Example

$$MS[i + 1].pos = SA[j], P[i] = A$$

	k	SA[k]	BWT[k]	$T[SA[k]..n]$	LCP[k]
	\vdots	\vdots	\vdots	\vdots	\vdots
	1234	8765	A	GAGACATCA...	-
$e_1 =$	1235	1519	A	GATACATTA ...	-
	1236	5450	C	GATAGATTA...	4
$j =$	1237	1004	G	<u>GATA</u> TAGAA...	4
	1238	4242	G	GATCCAATA...	3
$t =$	1239	3110	G	GATTACATA...	3
	1240	1102	T	GATTACTTA...	6
	1241	1978	T	GATTAGATA...	5
$s_2 =$	1242	2505	A	GATTATCAT ...	5
	1243	2022	A	GATTATGAA...	-
	\vdots	\vdots	\vdots	\vdots	\vdots

PHONI

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **PHONI**: Compute lengths with positions using LCEs
- **Start**: $MS[m - 1].len = 1$
- **Case 1 (match)**: $MS[i].len = MS[i + 1].len + 1$
- **Case 2 (mismatch)**:
 - $maxLCE = \max(LCE(SA[j], SA[e_1]), LCE(SA[j], SA[s_2]))$
 - $MS[i + 1].len = \min(MS[i + 1].len, maxLCE) + 1$

Finding Lengths

$MS[i + 1].len = 4, P[i] = A$

	k	SA[k]	BWT[k]	$T[SA[k]..n]$	LCP[k]
	⋮	⋮	⋮	⋮	⋮
	1234	8765	A	GAGACATCA...	-
$e_1 =$	1235	1519	A	GATACATTA ...	-
	1236	5450	C	GATAGATTA...	4
$j =$	1237	1004	G	<u>GATA</u> TAGAA...	4
	1238	4242	G	GATCCAATA...	3
	1239	3110	G	GATTACATA...	3
	1240	1102	T	GATTACTTA...	6
	1241	1978	T	GATTAGATA...	5
$s_2 =$	1242	2505	A	GATTATCAT ...	5
	1243	2022	A	GATTATGAA...	-
	⋮	⋮	⋮	⋮	⋮

Finding Lengths

$$MS[i + 1].len = 2, P[i] = A$$

	k	SA[k]	BWT[k]	$T[SA[k]..n]$	LCP[k]
	\vdots	\vdots	\vdots	\vdots	\vdots
	1234	8765	A	GAGACATCA...	-
$e_1 =$	1235	1519	A	GATACATTA ...	-
	1236	5450	C	GATAGATTA...	4
$j =$	1237	1004	G	<u>GAT</u> ATAGAA...	4
	1238	4242	G	GATCCAATA...	3
	1239	3110	G	GATTACATA...	3
	1240	1102	T	GATTACTTA...	6
	1241	1978	T	GATTAGATA...	5
$s_2 =$	1242	2505	A	GATTATCAT ...	5
	1243	2022	A	GATTATGAA...	-
	\vdots	\vdots	\vdots	\vdots	\vdots

Mismatches

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- **Case 2a:** Max LCE computed $> MS[i + 1].len$
- **Case 2b:** Max LCE computed $\leq MS[i + 1].len$

- Can we avoid the LCEs in 2b?

LCEs and Thresholds

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented Thresholds

Experiments

Thanks

- **LCE Queries:** With threshold, need only one LCE
 - $\text{LCE}(\text{SA}[i], \text{SA}[e_1])$ where $e_1 < i < t$
 - or $\text{LCE}(\text{SA}[i], \text{SA}[s_2])$ where $t \leq i < s_2$
- **LCEs as RMQ:**
 - Consider $i \geq t$ (other case symmetrical)
 - $\text{LCE}(\text{SA}[i], \text{SA}[s_2]) = \min(\text{LCP}[i + 1..s_2])$
 - $\text{LCE}(\text{SA}[t], \text{SA}[s_2]) = \min(\text{LCP}[t + 1..s_2])$
 - Since $t \leq i < s_2$,
 $\text{LCE}(\text{SA}[t], \text{SA}[s_2]) \leq \text{LCE}(\text{SA}[i], \text{SA}[s_2])$

Finding Lengths

$MS[i + 1].len = 2, P[i] = A$

	k	SA[k]	BWT[k]	$T[SA[k]..n]$	LCP[k]
	\vdots	\vdots	\vdots	\vdots	\vdots
	1234	8765	A	GAGACATCA...	-
$e_1 =$	1235	1519	A	GATACATTA ...	-
	1236	5450	C	GATAGATTA...	4
$j =$	1237	1004	G	GATATAGAA ...	4
	1238	4242	G	GATCCAATA...	3
$t =$	1239	3110	G	GATTACATA...	3
	1240	1102	T	GATTACTTA...	6
	1241	1978	T	GATTAGATA...	5
$s_2 =$	1242	2505	A	GATTATCAT ...	5
	1243	2022	A	GATTATGAA...	-
	\vdots	\vdots	\vdots	\vdots	\vdots

Augmenting Thresholds

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- Store $\text{LCE}(\text{SA}[t], \text{SA}[e_1])$ and $\text{LCE}(\text{SA}[t], \text{SA}[s_2])$
- e.g. $\text{LCE}(\text{SA}[t], \text{SA}[s_2]) \leq \text{MS}[i + 1].\text{len}$
 - Set $\text{MS}[i].\text{len} = \text{MS}[i + 1].\text{len} + 1$
 - Skip explicit LCE computation
- **Augmented Thresholds:**
 - Store threshold positions (one for each run)
 - With each threshold store the corresponding LCEs for before/after

Reducing LCEs

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

	<i>MONI</i>	<i>PHONI</i>	AUG
Case 1	None	None	None
Case 2a	1 SLP	1 or 2 LCEs	1 LCE
Case 2b	1 SLP	1 or 2 LCEs	None
<i>thresholds?</i>	yes	no	yes

Finding Threshold LCEs

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- MONI computes thresholds efficiently using Boucher et al's prefix free parsing (PFP)
 - Finds the minimum in the LCP array between run boundaries
- Threshold LCEs are computed similarly
 - Minimum between first boundary and threshold
 - Minimum after threshold and next boundary

Compressing Values

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- For each distinct character, threshold positions form increasing subsequences
- When $t = e_1 + 1$ or $t = s_2$, some threshold LCEs unused
 - Occurs up to roughly 1/3 of time on test data
 - Can 'null' these values
- Although LCPs can be large, expect median to be small
 - 99% fit in a byte on test data

Threshold LCE Variants

Augmented Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

- PHONI: Standard version of one-pass MONI
- Aug-Full: One-pass MONI modified with augmented thresholds
- Aug-1: One byte per threshold LCE, perform LCE on overflow
- Aug-BV-Full: Bitvector marks threshold LCEs which are non-zero and stored
- Aug-BV-1: As above, but one byte per threshold LCE (default to LCE query)
- Aug-DAC: Stores threshold LCEs using a directly addressable code (DAC)
- Aug-BV-DAC: Same as Aug-BV-Full, but substituting in a DAC

Datasets

- Patterns: 10 distinct chromosome-19 genomes
- References: chromosome-19 genomes of 16, 32, 64, 128, 256, 512 and 1000 copies
- Try compressed and plain SLPs

#	$n/10^6$	n/r	SLP _{comp} [MB]	SLP _{plain} [MB]
16	946.01	29.20	36.10	70.54
32	1892.01	57.64	37.80	74.75
64	3784.01	113.50	39.48	79.84
128	7568.01	222.24	42.11	88.89
256	15136.04	424.93	47.43	102.52
512	30272.08	771.54	58.00	131.09
1,000	59125.12	1287.38	80.63	186.98

Results

Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

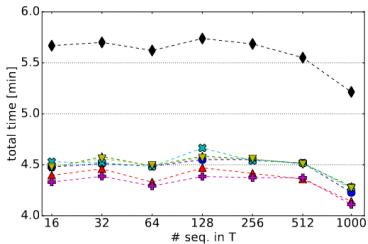
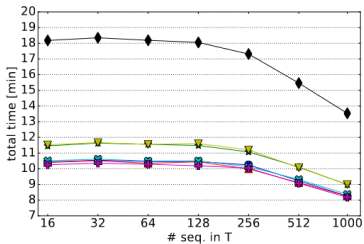
MONI

Augmented
Thresholds

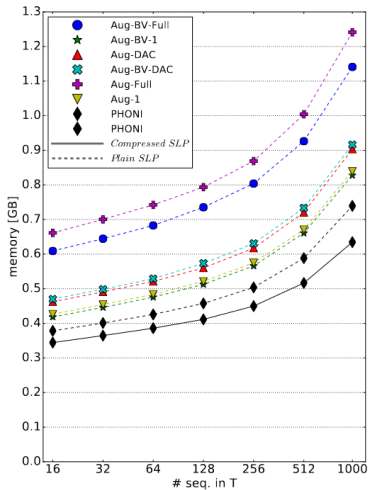
Experiments

Thanks

Average Query Time



Disk Size



Augmented
Thresholds

Guardiola,
Brown, Coira
et al.

Introduction

Preliminaries

MONI

Augmented
Thresholds

Experiments

Thanks

Thanks

- Email: `nathaniel.brown@dal.ca`
- Full-paper: <https://arxiv.org/abs/2211.07794>