# Compressed Direct-Address Table - practical index for genomic sequences

Katarzyna Jabłonowska

Michał Sabaciński
Norbert Dojer

Segovia
10 October 2019

# The problem

ATTATATAGTTGTAA
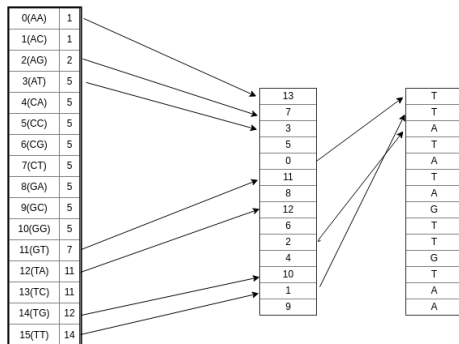
✓ GTTGT
✓ TAA
× *CTA*

Our goal is to find an algorithm that:

- scans the whole text only once
- is flexible
- is fast and memory efficient

# The basic idea

Two steps:

1. Building an index
2. Searching for patterns

Step 1

Step 2
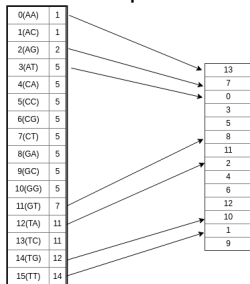
Step 3

Step 4

Prefix length = 2
Substring length = 3

# Searching for patterns



**ATAGT** → **AT**AGT

ATTATATAGTTGTAA

**Binary Search AT**A**GT**

**Verification**

| | |
|---|---|
| 0(AA) | 1 |
| 1(AC) | 1 |
| 2(AG) | 2 |
| 3(AT) | 5 |
| 4(CA) | 5 |
| 5(CC) | 5 |
| 6(CG) | 5 |
| 7(CT) | 5 |
| 8(GA) | 5 |
| 9(GC) | 5 |
| 10(GG) | 5 |
| 11(GT) | 7 |
| 12(TA) | 11 |
| 13(TC) | 11 |
| 14(TG) | 12 |
| 15(TT) | 14 |

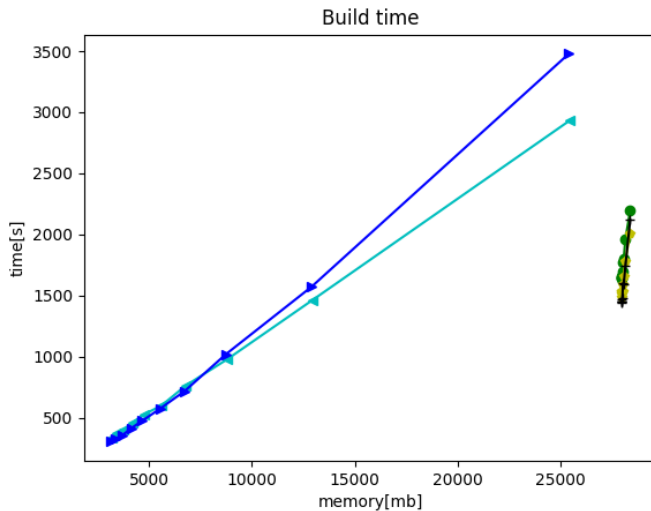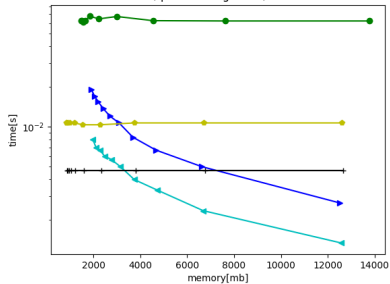| |
|---|
| 13 |
| 7 |
| 3 |
| 5 |
| 0 |
| 11 |
| 8 |
| 12 |
| 6 |
| 2 |
| 4 |
| 10 |
| 1 |
| 9 |

# Compression

1. 'GTCCAT'
2. 'GTCCAT'
3. 'GTCCAT*', $* \in \Sigma$
4. '*GTCCAT', $* \in \Sigma$
5. 'GTCCAT**', $* \in \Sigma$

## Experiments

- human genome
- C++
- prefixLength $= 10$ and 12
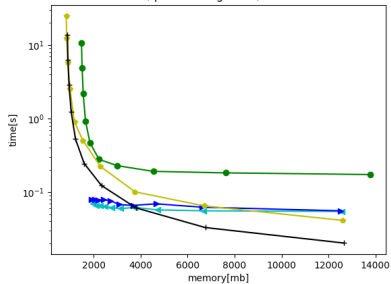- patterns of length 40 or 200. All or none of them appearing in the text

Build time

# Thank you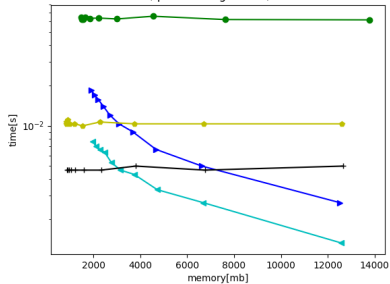