



UNIVERSITÀ
di **VERONA**

On different variants of the Burrows-Wheeler- Transform of string collections

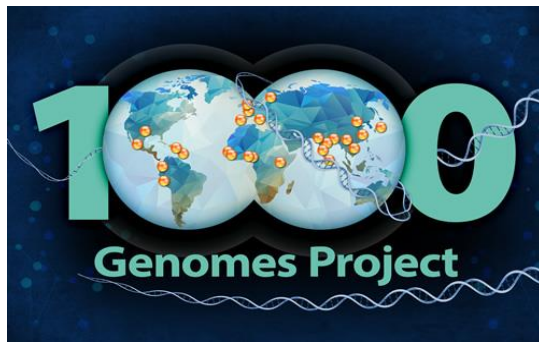
Davide Cenzato and Zsuzsanna Lipták

University of Verona, Department of Computer Science.

DCC 2022, March 22 – 25, 2022 - Snowbird, Utah, US

Large string collections are highly abundant

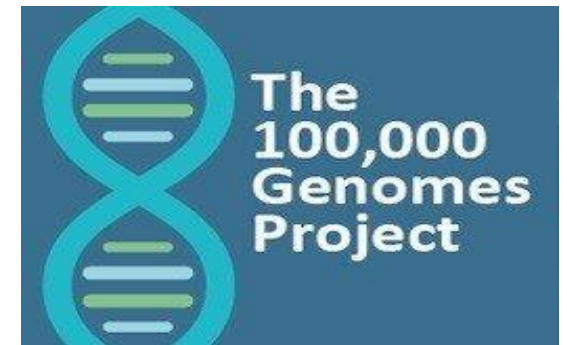
- number of sequenced genomes is growing at unprecedented pace
- focus has moved from single strings to string collections



2008



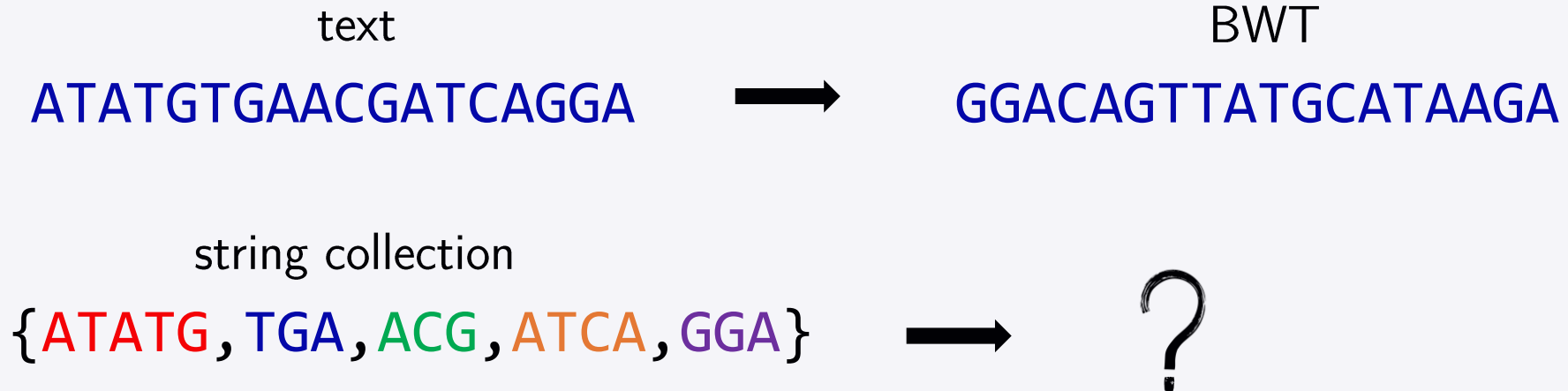
2009



2013

The Burrows-Wheeler-Transform for string collections

- originally defined on single strings
- not all tools compute the same transform



Each variant uses a **different order** for the characters in the interesting intervals.

BWT variant	criteria for characters ordering
eBWT(\mathcal{M})	omega-order of strings
doIEBWT(\mathcal{M})	lexicographic order of strings
mdoIBWT(\mathcal{M})	input order of strings
concBWT(\mathcal{M})	lexicographic order of subsequent strings in input
colexBWT(\mathcal{M})	colexicographic order (“reverse lex. order”)

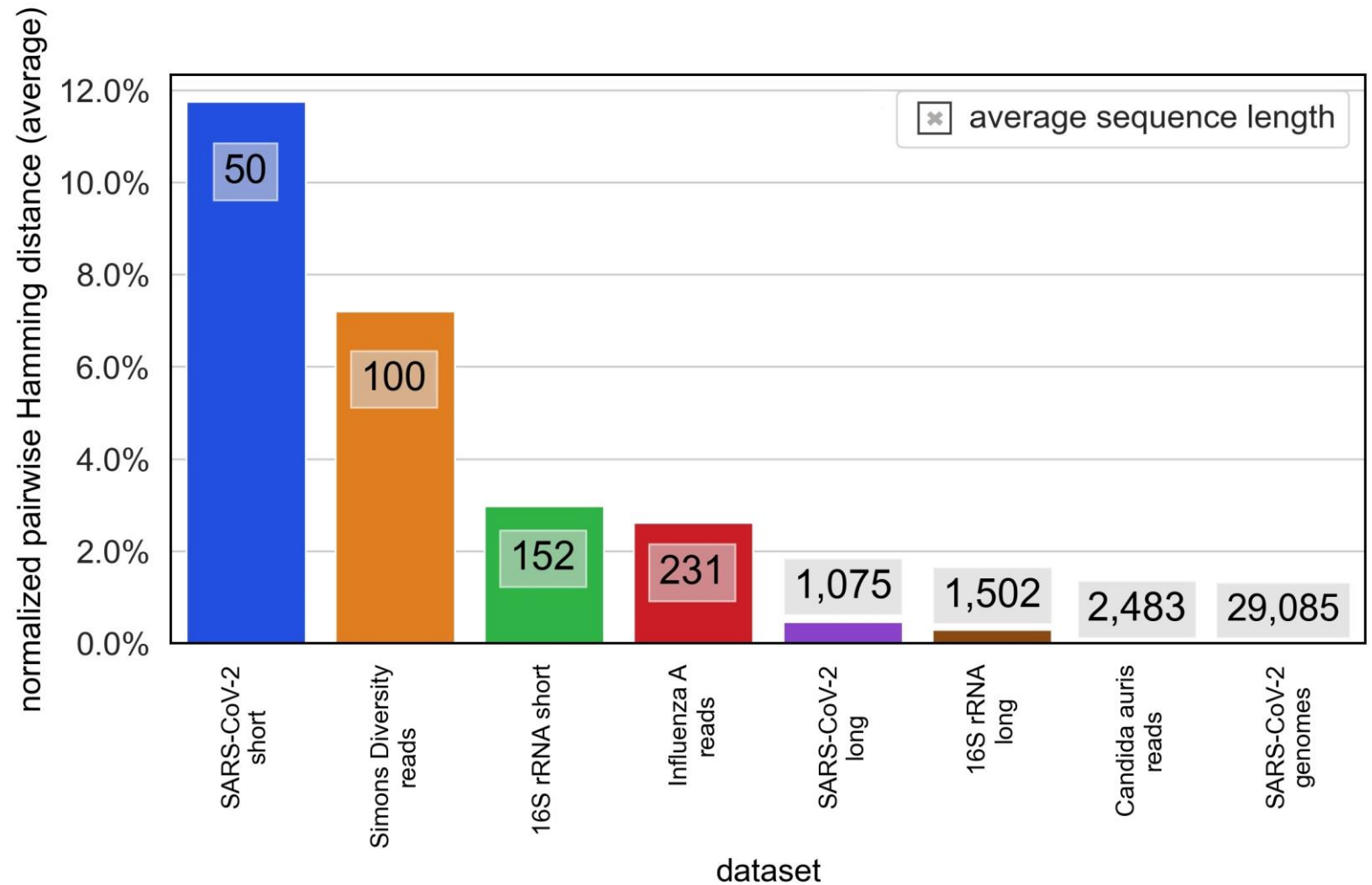
● depends on the **input order**:

$$\begin{aligned} \text{mdoIBWT}(\mathcal{M}_1) &= \text{GAGAAGCG} \text{ $$$TTATCTGAAA} & \mathcal{M}_1 &= \{\text{ATATG, TGA, ACG, ATCA, GGA}\} \\ & \quad \color{red}{\times\times\times} \quad \color{red}{\times\times} \quad \color{red}{\times\times\times\times} \\ \text{mdoIBWT}(\mathcal{M}_2) &= \text{GGAAAGGC} \text{ $$$TACTGTAAA} & \mathcal{M}_2 &= \{\text{ACG, ATATG, GGA, TGA, ATCA}\} \end{aligned}$$

Experimental results 1: Hamming distance

We conducted experiments on 8 real-life datasets to determine how much these differences matter.

- strongly depends on **sequence length**
- SARS-Cov-2 short: **500,000** sequences of length **50**
- on average **almost 12%** different positions



Experimental results 2: number of runs

These differences extend to the number of runs (r) of all BWT variants.

- **number of runs** is more variable on short sequence datasets
- **average runlength** (n/r) is an equivalent measure

no. runs SARSCov2short dataset		
	r	n/r
eBWT	1,902,148	13.143
doIBWT	1,868,581	13.647
mdoIBWT	3,113,818	8.189
concBWT	3,402,513	7.494
colexBWT	808,906	31.524

Different tools compute different BWT variants

- identified **five** different BWT variants computed by these tools
- some of these are sensitive to **input order**

Differences between BWT variants are not negligible

- more relevant on **short sequences**
- extends to the **number of runs** of the BWT variants



UNIVERSITÀ
di **VERONA**

Thank you for your attention

contact: `davide.cenzato@univr.it`

GitHub: `https://github.com/davidecenzato/BWT-variants-of-string-collections`

arXiv: `http://arxiv.org/abs/2202.13235`