



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Compression of k-mers set with counters

Leonardo Gemin, Enrico Rossignolo

November 10, 2022



- Introduction
- k-mer set representations
- Methods
- Results
- Conclusions



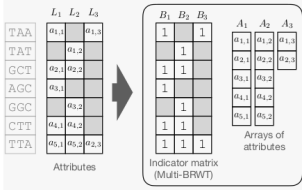
- k-mer \longrightarrow substring of length k
- Huge list (4^k) \implies bottleneck
- Wide range of application:
 - genome assembly
 - metagenomics
 - database searching
 - ...



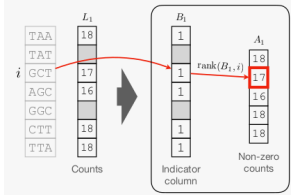
- (counting) de Bruijn graphs \longrightarrow Metagraph
- contigs \longrightarrow Metagraph
- unitigs \longrightarrow BCALM2
- simplitigs \longrightarrow prophAsm
- SPSSs \longrightarrow UST

- **dBG(K):** $G = (V, A)$ where
 - $V = K$
 - $A = \{(u, v) \in K^2 \mid \text{pref}_{k-1} u = \text{pref}_{k-1} v\}$
- counting: external data structures

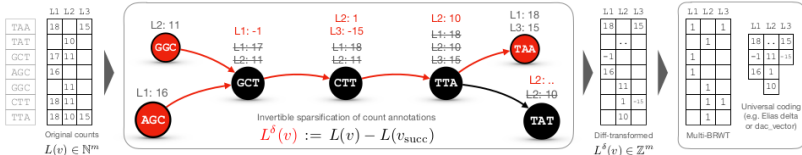
A. General scheme for multiple columns



B. Representation of one column with integers



- RRR and Multi-BRWT for binary matrix
- row diff transform for annotations arrays
 - successors assigned arbitrarily
 - anchors are unchanged and break the recursion
 - \implies labels represented with less bits!
 - \implies binary matrix sparsified!





- **contigs**: any path in a dBG
- **unitigs**: any path with non-branching internal nodes
- **simplitigs**: ?

- **simplitigs**: any path on a disjoint path cover in a dBG
- $CL = n_k + (k - 1) \cdot NS \implies CL \geq n_k + k - 1$
- not tight: $CL^* = 9 + (3 - 1) = 11 < 15$
- ProphAsm is sub-optimal with a greedy approach

Individual k-mers: ACG, CGA, GAA, AAG, AGC, CGT, GTA, TAG, CGG



Maximal unitigs: ACG, CGAAG, CGTAG, AGC, CGG



Maximal simplitigs: ACGAAGC, CGTAG, CGG



- **k-spectrum** of a string s :

$$sp^k(s) \triangleq \{\{\text{canonized k-mers of } s\}\}$$

- for a set S :

$$sp^k(S) \triangleq \bigcup_{s \in S} sp^k(s)$$

- S is a **Spectrum Preserving String Set** of a k-mer set K iff

- 1 $sp^k(S) = sp^k(K)$
- 2 $|s| \geq k \quad \forall s \in S$

- Let S^{opt} be a minimum SPSS representing K and W^{opt} a minimum path cover of $cdBG(K)$. Then

$$weight(S^{opt}) = |K| + |W^{opt}| \cdot (k - 1)$$

- W is an SPSS of $K \implies$ equivalence with simplitigs
- Lower bound (not tight)

$$weight(S^{opt}) \geq |K| + \left(\left\lceil \frac{n_{dead} + n_{sp}}{2} \right\rceil + n_{iso} \right) \cdot (k - 1)$$

where $n_{sp} = \sum_{(u,su) \in (V, \{0,1\})} \max\{0, |B_{u,su}| - 1\}$

- UST is sub-optimal with a greedy approach

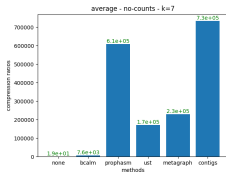


Accession	Description
SRR000001	Human haplotype map
SRR21394969	E. Coli
SRR21394970	Coronavirus 2
SRR21284212	Salmonella
SRR21073883	N. Gonorrhea
RND1664714109	Random sequence

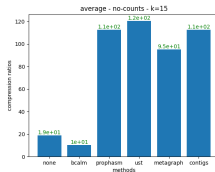
Workflow:

- from FASTA to k-mer
 - help: bash script
 - k-mer length: 7, 15, 23, 31
- with vs without counts
- compression
 - lzma
 - MFCompress
 - gZip

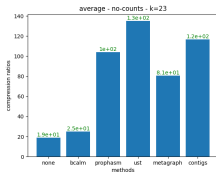
Results (without counts)



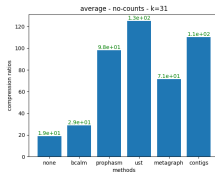
(a) $k = 7$



(b) $k = 15$

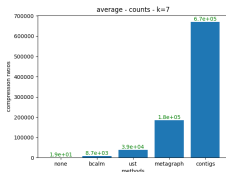


(c) $k = 23$

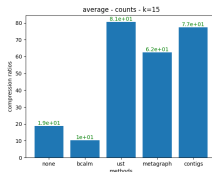


(d) $k = 31$

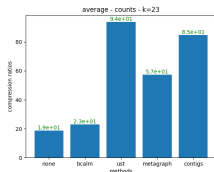
Results (with counts)



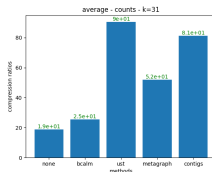
(a) $k = 7$



(b) $k = 15$



(c) $k = 23$



(d) $k = 31$

SRR21073883 (N. **Gonorrhea**): really high compression ratios!

k-mer lenght	bcalm	prophasm	ust	metagraph	contigs
15	34.37	352.40	393.78	308.85	366.26
23	80.034	323.00	442.52	256.55	369.77
31	90.89	304.78	409.91	224.03	346.93

(a) without counts

k-mer length	bcalm	ust	metagraph	contigs
15	34.12	260.98	201.39	252.74
23	75.41	308.95	181.29	274.31
31	83.90	300.16	162.52	262.04

(b) with counts

- First observations:
 - tools return compression ratios higher than with a direct compression
 - low k-mer length \rightarrow high compression ratio
- Bcalm \rightarrow no advantage, depends on repetitiveness
- ProphAsm (only without counts) \rightarrow better with $k = 7$
- UST \rightarrow the best among all
- Metagraph \rightarrow not particularly efficient
 - contigs \rightarrow rivals UST



Thanks for your attention!



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

Compression of k-mers set with counters

Leonardo Gemin, Enrico Rossignolo

November 10, 2022