

Reviews in Quantitative Biology

Applications of kmer analysis in quantitative biology

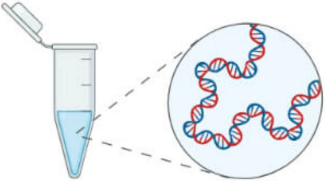
Sina Majidian

November 18th, 2022

Outline

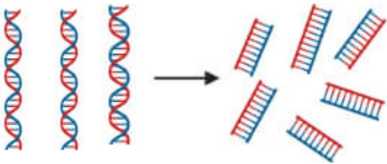
- **Introduction**
 - DNA sequencing
 - Kmer definition
 - Minimizer definition
- **Applications**
 - Metagenomics
 - Genome assembly
 - Kmer counting
- **Discussion**

DNA Sequencing



DNA extraction

GGCGTCTATATCTCGGCTCTAGGCCCTCATT TTTT



Amplification

GGCGTCTATATCTCGGCTCTAGGCCCTCATT TTTT
GGCGTCTATATCTCGGCTCTAGGCCCTCATT TTTT
GGCGTCTATATCTCGGCTCTAGGCCCTCATT TTTT
GGCGTCTATATCTCGGCTCTAGGCCCTCATT TTTT



Sequencing

GGCGTCTA TATCTCGG CTCTAGGCCCTC ATTTTTT
GGC GTCTATAT CTCGGCTCTAGGCCCTCA TTTTTT
GGCGTC TATATCT CGGCTCTAGGCCCT CATT TTTT
GGCGTCTAT ATCTCGGCTCTAG GCCCTCA TTTTTT

Genomics pipelines

```
CATCGACCGAGCGCGATGCTAGCTAGGTGATCGT  
TGCCGCATCGACCGAGCGCGATGCTAGCTAGGTGATCGT  
GCATGCCGCATCGACCGAGCGCGATGCTAGCTAGGTGATCGT  
GTGCATGCCGCATCGACCGAGCGCGATGCTAGCTAGGTGATC  
.AGGTGCATGCCGCATCGATCGAGCGCGATGCTAGCTAGCTGATCGT
```



1) De novo genome assembly

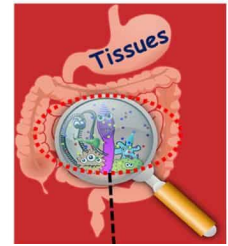
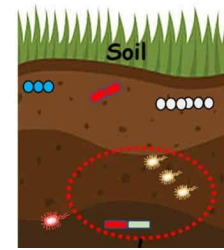
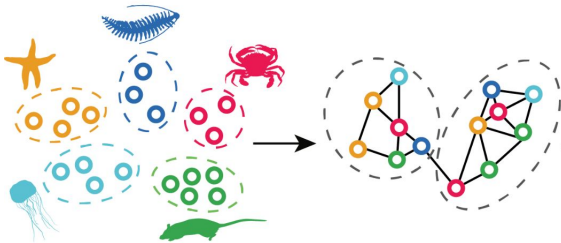
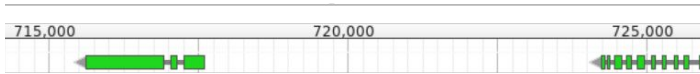
4) Read alignment

5) Variant calling & genotyping

2) Gene annotation

6) Metagenomics

3) Orthology inference

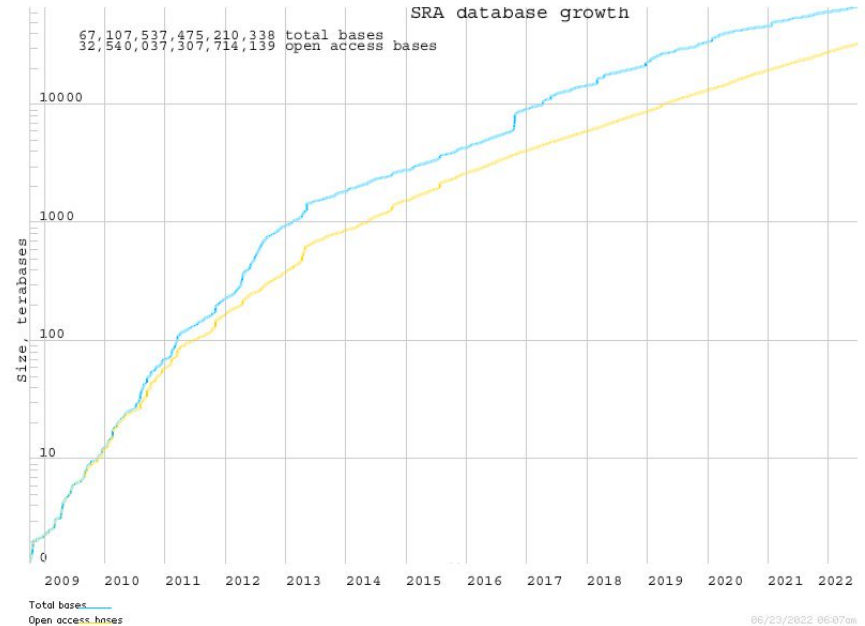
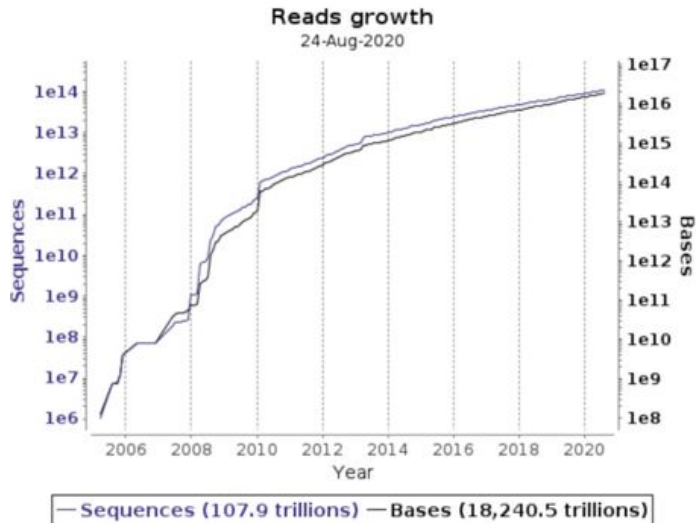


Genome assembly: importance of scalable methods

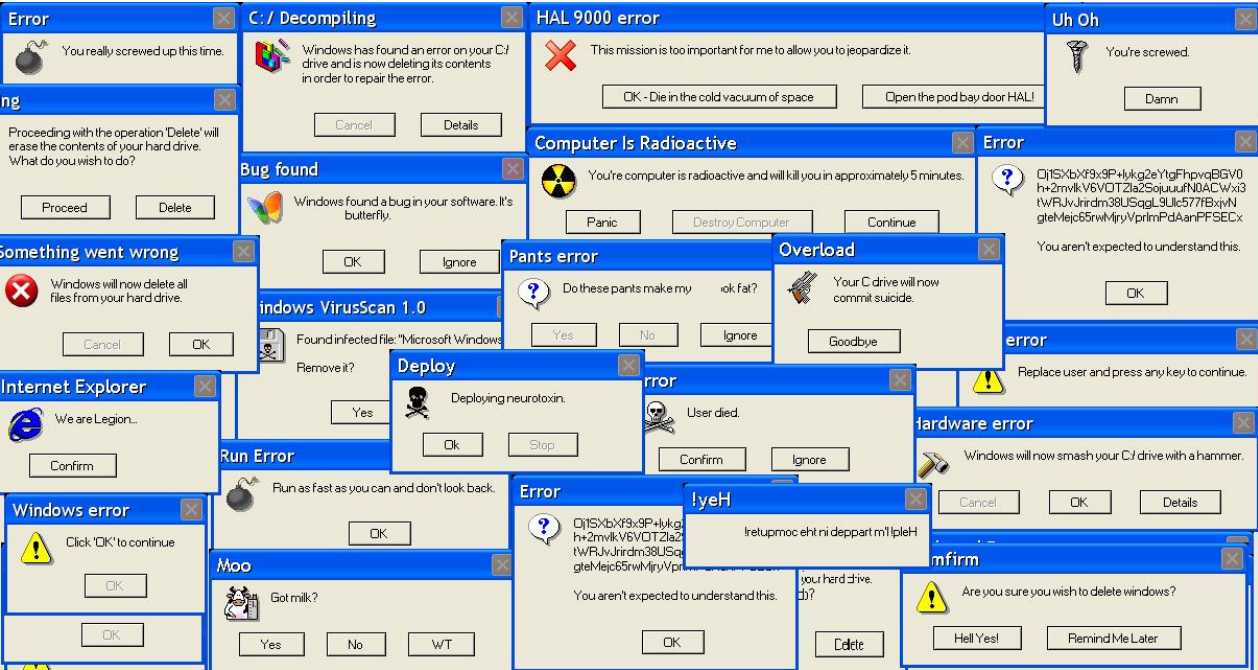
	<i>D. mel</i> 100× real HiFi reads				Human real HiFi reads		
Tool	Peregrine	HiCanu	Hifiasm	Rust-mdbg	Peregrine	Hifiasm	Rust-mdbg
Time	40 min 11s	7 h 43min	5 h 17min	1 min 9 s	14 h 8 min	58 h 41min	10 min 23 s
Memory	12 GB	12 GB	21 GB	1.5 GB	188 GB	195 GB	10 GB

We need efficient methods in terms of **memory**, **speed** and **storage**.

Fast growing genomics data



No Space Left on Device



Segmentation fault

Segmentation fault

Core dumped

Too many open files

ERROR ERROR ERROR

Basics on K-mer & Minimizers

K-mer definition

- a contiguous substring of length k.
- alphabet= {A,C,G,T}
- string (sequence)= GATTACA
- Example: k=2, 2mers
- GA, AT, TT, TA, AC, CA
- Number of 2-mers = 6

GA**TT**ACA

G**AT**TACA

GAT**TT**ACA

GATT**TA**CA

GATTAC**A**

GATTAC**CA**

K-mer stats

For the same sequence GATTACA (L=7)

k=3: GAT, ATT, TTA, TAC, ACA (5 kmers)

k=4: GATT, ATTA, TTAC, TACA (4 kmers)

GATTACA

GATTACA

GATTACA

GATTACA

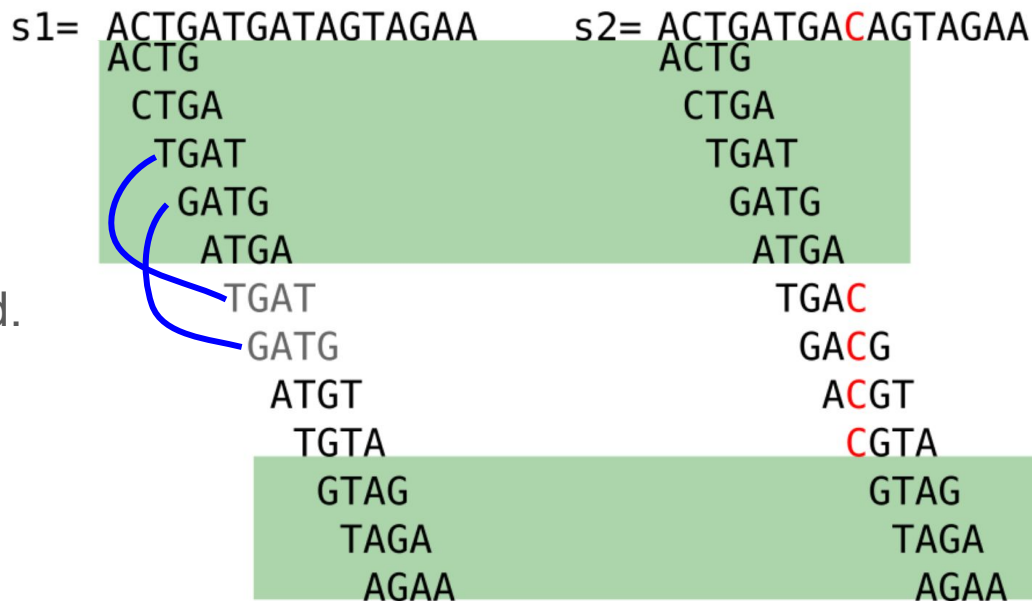
GATTACA

- How many 3-mers do exist for a sequence of length L=100?
 - Number of 3-mers = $100-3+1 = 98$ (L-k+1)
- How many different 3-mers do exist?
 - $4^3=64$ (4^k)

--	--	--

Sequence comparison

Two similar sequences **with one base error**:



- Many shared 4-mers (green)
- Neighbors of C/T are not shared.
- Repeated kmers

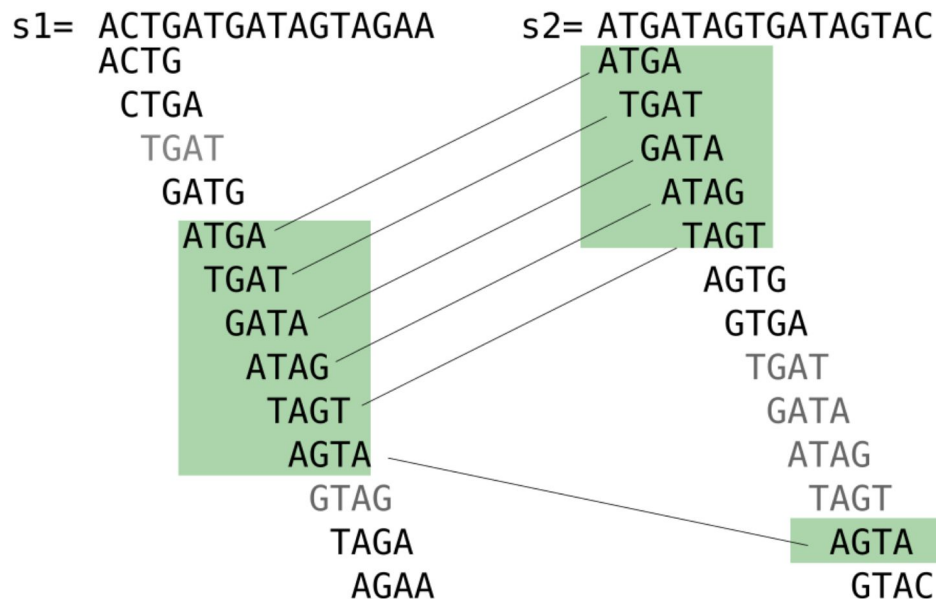
Impact of k on specificity

By increasing k, we can keep the differences.

s1=	ACTGATGATAGTAGAA	s2=	ACTGATGACAGTAGAA
	ACTGAT		ACTGAT
	CTGATG		CTGATG
	TGATGA		TGATGA
	GATGAT		GATGAC
	ATGATG		ATGACG
	TGATGT		TGACGT
	GATGTA		GACGTA
	ATGTAG		ACGTAG
	TGTAGA		CGTAGA
	GTAGAA		GTAGAA

Impact of k on specificity

- Two not so similar sequences:
- Still a lot of shared 4-mers
- Number 5-mers = $4^5 = 1024$
- Every 1k, one shared kmer by chance. (if distributed uniformly)



K-mer stats

- Billions (10^9) of k-mers exist in a metagenomic datasets,
- resulting in computational challenges

Dataset	k	Total	Distinct
Cow rumen	28	7.39×10^9	5.09×10^9
	55	4.69×10^9	3.73×10^9
Marine	28	7.26×10^9	3.73×10^9
	55	4.57×10^9	3.04×10^9

- Each k-mer contains little information compared to the nearest one.
- Idea: sampling (also referred to as “sketching”)

GATTACA

GATTACA

GATTACA

GATTACA

GATTACA

Minimizer

- The smallest kmer in a window.
- Alphabetical order: AA < AT
- Window of size 5, 3-mers.

1st window GATTA : ATT

2nd window ATTAC : ATT

3rd window TTACA : ACA

Minimizers for GATTACA are ATT, ACA.

(2 fingerprints instead of 5 kmers)

Similar sequences ~ shared minimizers.

GATTA

GATTA

GATTA

ATTAC

ATTAC

ATTAC

TTACA

TTACA

TTACA

GATTACA

GATTACA

1st window

GATTACA

2nd window

GATTACA

3rd window

Different ordering

1. Alphabetical order: AA < TT
2. Defining a function on a string

A=1, C=2, G=3, T=4

TAAT = 4114

CTTT = 2444

TAAT > CTTT

3. Hash function

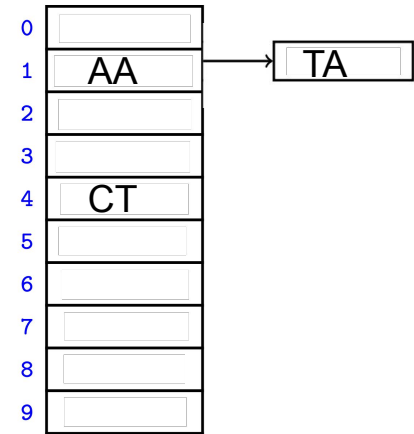
$$h(x) = x \bmod 11$$

“remainder of a division by 10”

$$AA=11 \rightarrow \{AA:1\}$$

$$CT: 24 = 2 \cdot 10 + 4$$

$$TA: 41 = 4 \cdot 10 + 1$$

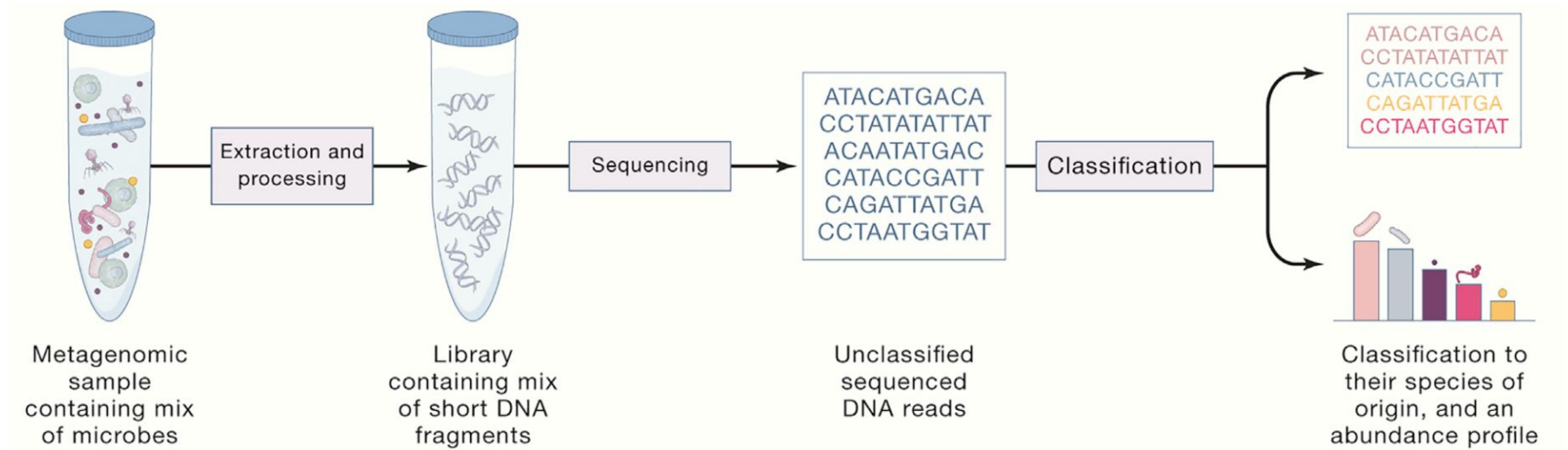


Applications

1- Metagenomics

Metagenomics

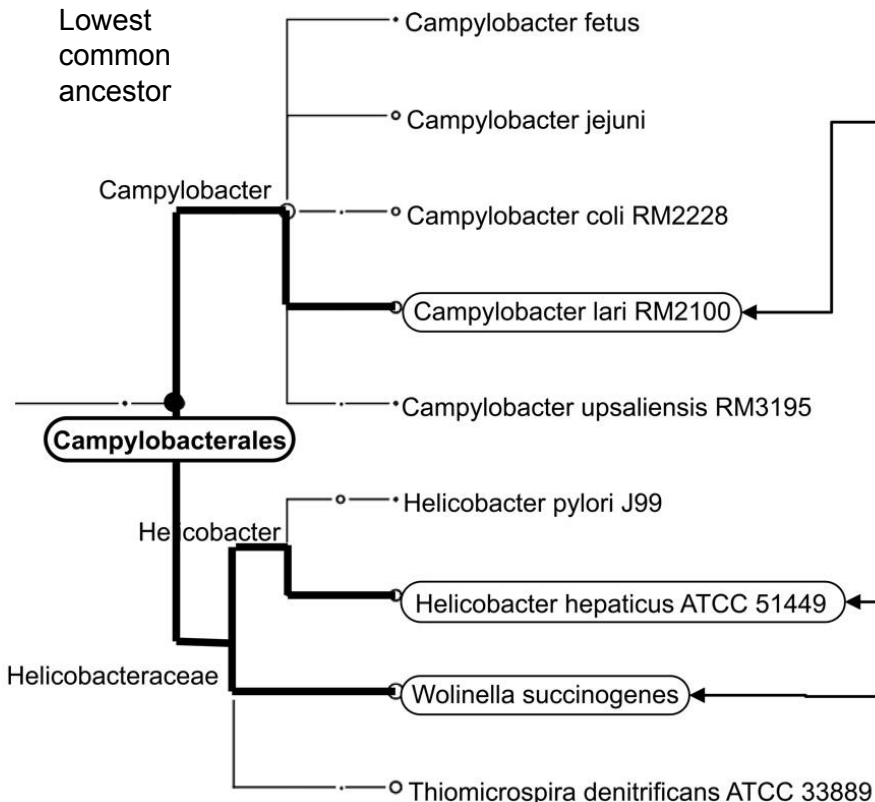
Metagenomics: the study of genomic sequences obtained directly from an environment (gut microbiome, wastewater surveillance, ..)



MEGAN: exploring taxonomical content

BLASTX [read1](#) against NCBI-NR

Lowest
common
ancestor



```
>gi|57241447|ref|ZP_00369393.1| flagellar motor switch protein FligG  
[Campylobacter lari RM2100]  
Score = 33.9 bits (76), Expect = 1.8  
Identities = 13/26 (50%), Positives = 19/26 (73%)  
  
Query: 79 LMFVFDLLATVEENGINEIREIINRADKK 2  
+MF FDD++ + N IRE++ ADK+  
Sbjct: 243 LMFTFDDISQLSTNAIREVLKAADKR 268
```

1st match
for read1

```
>gi|32262158|gb|AAP77207.1| flagellar motor switch protein FligG  
[Helicobacter hepaticus ATCC 51449]  
Score = 33.5 bits (75), Expect = 2.4  
Identities = 13/26 (50%), Positives = 20/26 (76%)  
  
Query: 79 LMFVFDLLATVEENGINEIREIINRADKK 2  
+MF F+D++ ++ N IREI+ ADKK  
Sbjct: 244 MMFTFEDISKLDNNAIREILKIADKK 269
```

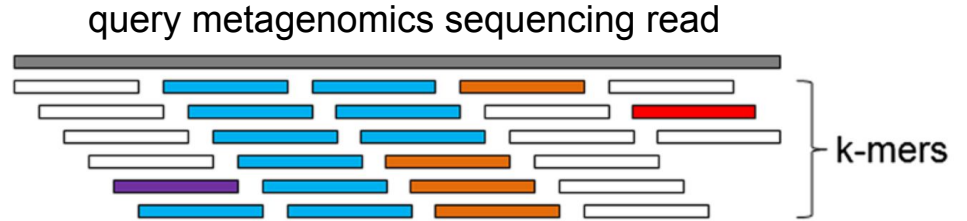
2nd match
for read1

```
>gi|34484004|emb|CAE11000.1| FLAGELLAR MOTOR SWITCH PROTEIN FLIG  
[Wolinella succinogenes]  
Score = 32.7 bits (73), Expect = 4.1  
Identities = 13/26 (50%), Positives = 19/26 (73%)  
  
Query: 79 LMFVFDLLATVEENGINEIREIINRADKK 2  
+MF F+D+ ++ N IREI+ ADKK  
Sbjct: 242 MMFTFEDIEKLDNNAIREILKVADKK 267
```

3rd match
for read1

Kraken: assigning taxonomic labels

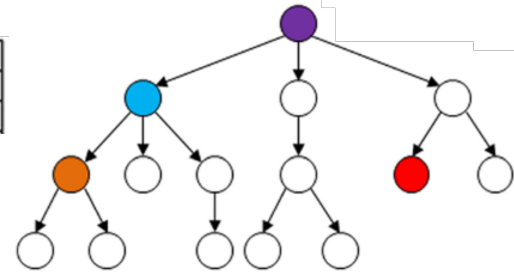
1. Extract k-mers from the input sequencing read (query)



2. Use precomputed database, extract the taxonomic ID of 31-mers

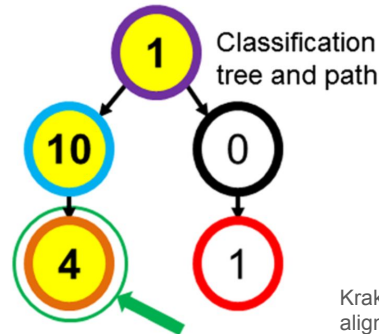
31-mers	Taxonomy ID
GACAAGGAGGCTATTGAGAAGAATATGGAAC	<i>Prevotella intermedia</i>
TAGCAGGCAGGTTACCCGGGTGGCGGTTTT	<i>Rothia mucilaginosa</i>
ATATGATTTTGAGCGAAAAATGAATAAAATT	<i>Streptococcus pneumoniae</i>

3. Count # kmers assigned to each taxonomic level



4. Find the path with maximum sum

5. Report the leaf as the taxonomic label



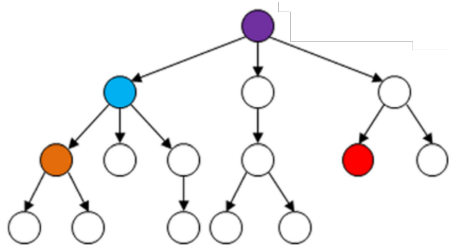
Kraken: kmer-LCA table

- Pre-computing the kmer-LCA database
- Extract all 31-mers from microbial genome in NCBI RefSeq

31-mers	Taxonomy ID
GACAAGGAGGCTATTGAGAAGAATATGGAAC	<i>Prevotella intermedia</i>
TAGCAGGCAGGTTACCCGGGTGGCGGTTTT	<i>Rothia mucilaginosa</i>
ATATGATTTTGAGCGAAAAATGAATAAAATT	<i>Streptococcus pneumoniae</i>



31-mers



- To speed up kmer search in the table
- Use 15mer minimizer for grouping similar 31mer

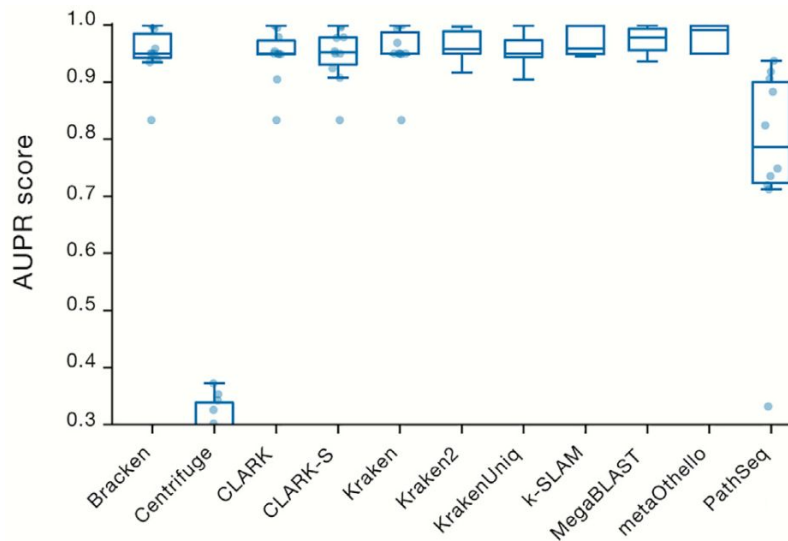
31-mers	Taxonomy ID

- Fill the table with 31-mers of shared minimizer near each other
- For a 31-mer query, start searching within its minimizer

Kraken2 results

Classification accuracy and speed comparison on a simulated dataset

Classifier	Memory Required	Time Required
Bracken	<1 Gb	<1 min
Centrifuge	20 Gb	7 min
CLARK	80 Gb	2 min
CLARK-S	170 Gb	40 min
Kraken	190 Gb	1 min
Kraken2	36 Gb	1 min
KrakenUniq	200 Gb	1 min
k-SLAM	130 Gb	2 h
MegaBLAST	61 Gb	4 h
metaOthello	30 Gb	1 min
PathSeq	140 Gb	5 min



Article | [Open Access](#) | [Published: 11 July 2019](#)

Strain-level metagenomic assignment and compositional estimation for long reads with MetaMaps

[Alexander T. Dilthey](#) , [Chirag Jain](#), [Sergey Koren](#) & [Adam M. Phillippy](#)

Nature Communications **10**, Article number: 3066 (2019) | [Cite this article](#)

15k Accesses | 49 Citations | 68 Altmetric | [Metrics](#)

Short Report | [Open Access](#) | [Published: 28 November 2019](#)

Improved metagenomic analysis with Kraken 2

[Derrick E. Wood](#), [Jennifer Lu](#) & [Ben Langmead](#) 

Genome Biology **20**, Article number: 257 (2019) | [Cite this article](#)

52k Accesses | 1081 Citations | 91 Altmetric | [Metrics](#)

Journals & Magazines > [IEEE/ACM Transactions on Comp...](#) > Volume: 19 Issue: 1 

K2Mem: Discovering Discriminative K-mers From Sequencing Data for Metagenomic Reads Classification

Publisher: **IEEE**

[Cite This](#)

[PDF](#)

[Improved kraken2](#)

[Davide Storato](#) ; [Matteo Comin](#)  [All Authors](#)

Research | [Open Access](#) | [Published: 31 October 2022](#)

Sketching and sampling approaches for fast and accurate long read classification

[Arun Das](#)  & [Michael C. Schatz](#)

BMC Bioinformatics **23**, Article number: 452 (2022) | [Cite this article](#)

19 Accesses | 3 Altmetric | [Metrics](#) [Comparing minimizer with MinHash & ..](#)



< [BIOINFORMATICS AND GENOMICS](#)

Deconvolute individual genomes from metagenome sequences through short read clustering



[Research article](#) [Bioinformatics](#) [Microbiology](#) [Data Science](#)



[Kexue Li](#) ^{1,2}, [Yakang Lu](#) ^{1,2}, [Li Deng](#) ^{1,2,3}, [Lili Wang](#)^{1,2}, [Lizhen Shi](#)⁴, [Zhong Wang](#) ^{3,5,6}

Published April 8, 2020

[Metagenome assembly](#)

MetaProb 2: Metagenomic Reads Binning Based on Assembly Using Minimizers and K-Mers Statistics

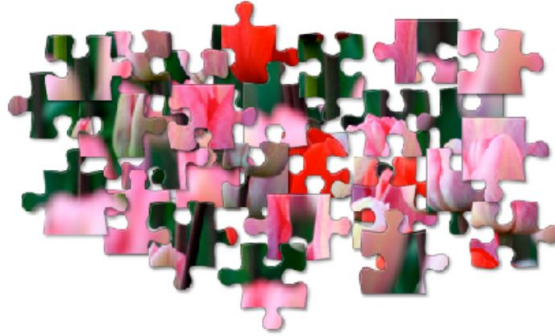
FRANCESCO ANDREACE, CINZIA PIZZI, and MATTEO COMIN[†]

[Metagenome assembly](#)

Applications

2- Genome assembly

De novo genome assembly



Sequencing reads

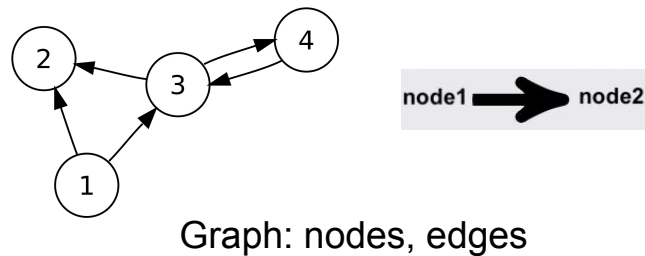
CTAGGCCCTCAATTTTT
GGCGTCTATATCT
CTCTAGGCCCTCAATTTTT
TCTATATCTCGGCTCTAGG
GGCTCTAGGCCCTCATTTTTT
CTCGGCTCTAGCCCCTCATT
TATCTCGACTCTAGGCCCTCA
GGCGTCGATATCT
TATCTCGACTCTAGGCC
GGCGTCTATATCTCG



CTAGGCCCTCAATTTT
CTCTAGGCCCTCAATTTT
GGCTCTAGGCCCTCATTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
TATCTCGACTCTAGGCC
TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCG
GGCGTCGATATCT
GGCGTCTATATCT
GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTT

- de Bruijn graphs
 - assemble the genome.
 - finding overlaps between kmers of reads

de Bruijn graphs



- Nodes: kmers
- Edge: link between k-mers that overlap by k-1 bases.

- Example:

- 3mers for the sequence **ACTG**
- 3mers: ACT, CTG

ACT → **CTG**

- Merging two nodes result in a K+1 mer (=a part of the sequence).

dBG for few reads

ACTG

CTGC

TGCC

ACTG

ACTG

CTGC

CTGC

CTGC

TGCC

TGCC

ACT → CTG → TGC → GCC

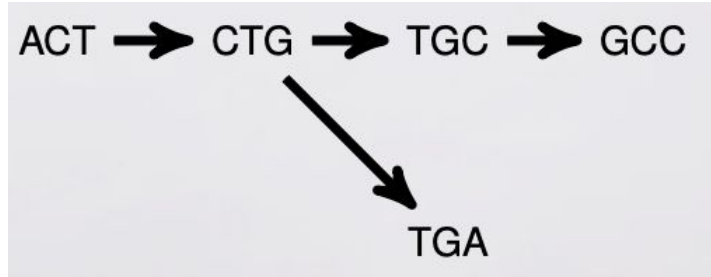
A sequencing error on de Bruijn graph

ACTG

CTGC

TGCC

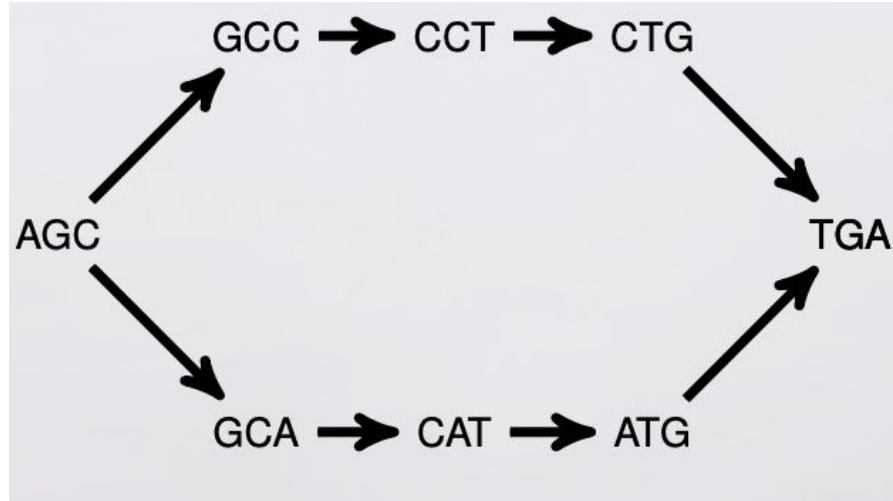
CTGA



SNPs in the graph

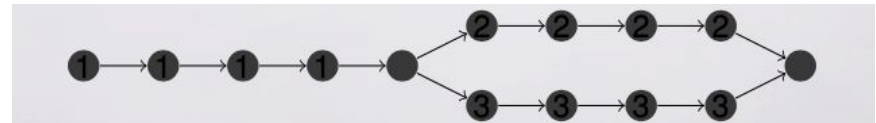
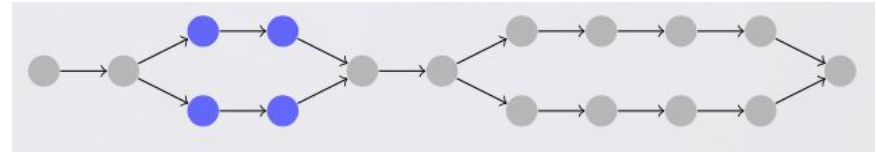
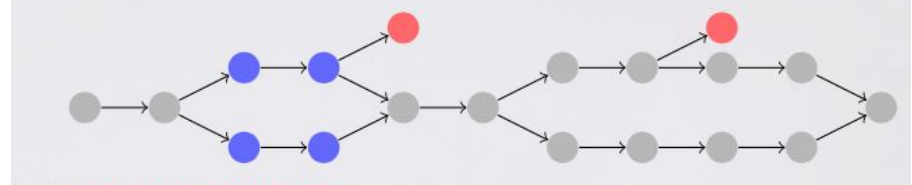
AGC**C**TGA

AGC**A**TGA



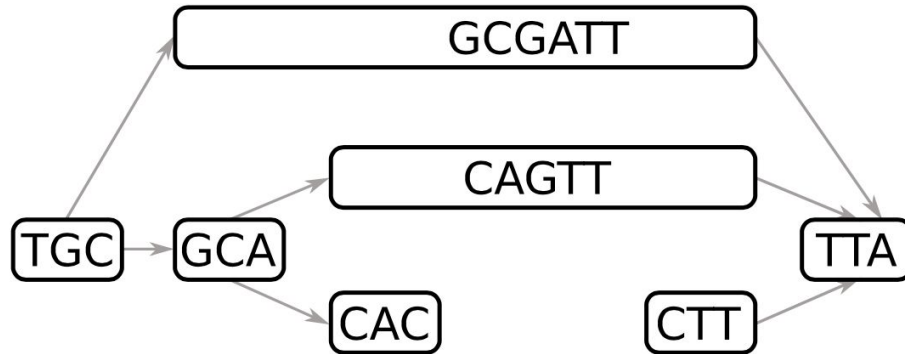
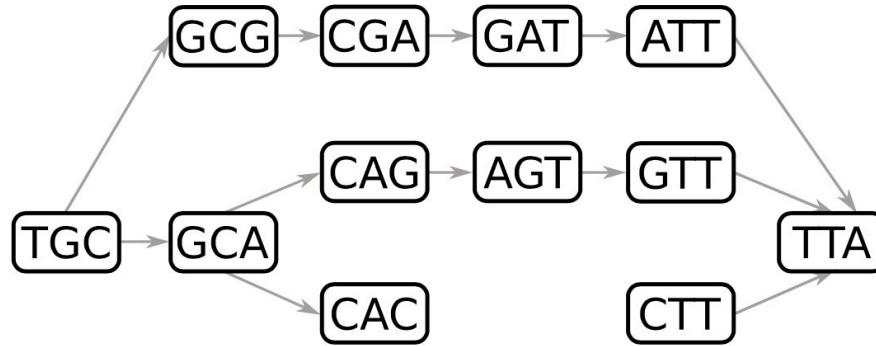
Short read assemblers

- Construct de Bruijn graph
- Remove likely sequencing errors
- Remove variations
- Return simple paths (contigs)
 - Combining all kmers in the path
- Extra steps: repeat-resolving, scaffolding



Compacting the de Bruijn graph

1 TB seq data
700 GB kmers

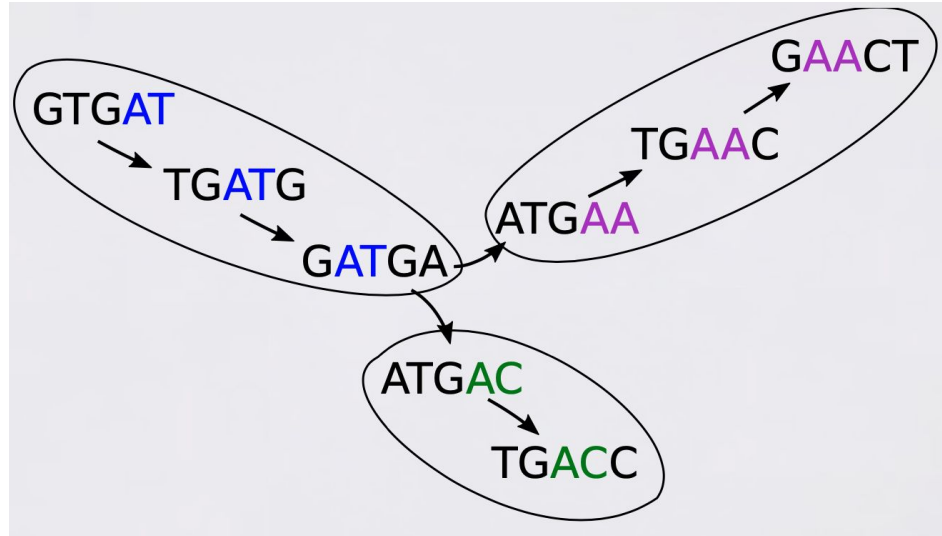


Compacting in parallel

GTG**AT**GA

ATG**ACC**

ATG**AACT**



Grouping kmers based on their minimizer to work it parallel.

BCALM2 is a tool for compacting the de Bruijn graph.

Bioinformatics, 32, 2016, 1201–1208
doi: 10.1093/bioinformatics/btw279
ISMB 2016



Compacting de Bruijn graphs from sequencing data quickly and in low memory

Bioinformatics, 37(16), 2021, 2476–2478
doi: 10.1093/bioinformatics/btab004
Advance Access Publication Date: 21 January 2021
Applications Note



Sequence analysis

MBG: Minimizer-based sparse de Bruijn Graph construction

Holley and Melsted *Genome Biology* (2020) 21:249
<https://doi.org/10.1186/s13059-020-02135-8>

Genome Biology

SOFTWARE

Open Access

Bifrost: highly parallel construction and indexing of colored and compacted de Bruijn graphs



JOURNAL OF COMPUTATIONAL BIOLOGY
Volume 22, Number 5, 2015
© Mary Ann Liebert, Inc.
Pp. 336–352
DOI: 10.1089/cmb.2014.0160

Research Articles

On the Representation of De Bruijn Graphs

Too much
math

A Preprocessor for Shotgun Assembly of Large Genomes

Ye et al. *BMC Bioinformatics* 2012, **13**(Suppl 6):S1
<http://www.biomedcentral.com/1471-2105/13/S6/S1>



PROCEEDINGS

Open Access

Exploiting sparseness in *de novo* genome assembly

Bioinformatics, 36(12), 2020, 3885–3887
doi: 10.1093/bioinformatics/btaa253
Advance Access Publication Date: 20 April 2020
Applications Note

OXFORD

Genome analysis

ntJoin: Fast and lightweight assembly-guided scaffolding using minimizer graphs

ARTICLES

<https://doi.org/10.1038/s41587-020-00747-w>

nature
biotechnology



OPEN

Efficient hybrid *de novo* assembly of human genomes with WENGAN

CellPress
OPEN ACCESS
Focus on RECOMB

Cell Systems

Article

Minimizer-space de Bruijn graphs: Whole-genome assembly of long reads in minutes on a personal computer

LJA: Assembling Long and Accurate Reads Using Multiplex de Bruijn Graphs

Anton Bankevich, Andrey Bzikadze, Mikhail Kolmogorov, Dmitry Antipov, Pavel A. Pevzner
doi: <https://doi.org/10.1101/2020.12.10.420448>

nature
biotechnology

ARTICLES

<https://doi.org/10.1038/s41587-022-01220-6>



Multiplex de Bruijn graphs enable genome assembly from long, high-fidelity reads

Applications

3- Kmer counting

K-mer counting

- Given a string, counting the number of occurrences of every kmer.
- Many applications

- A naive approach:
 - `count_dict={}`; `count_dict['ATA']=1`; `count_dict['ATC']=2`.
 - k-mers as keys and their counts as values.

- Efficient in terms of speed, space and memory
- Parallelization: locking memory

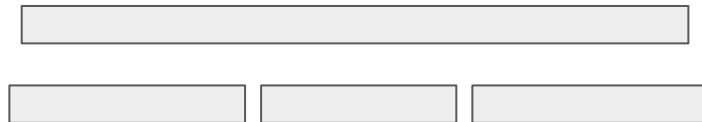
JellyFish: based on Hash table

- k-mers are stored as keys and their counts are stored as values.
- Lock-free hash table allowing parallel insertion of k-mers and frequency updates.



CAG	62
GGA	33

Gerbil: with minimizer



Evaluation

Data

- Illumina reads
- H. sapiens
- 292 GByte

System

- memory: 64GB
- Hard disk 1TB
- CPUs 16
- Intel CPU E5-2698 2.30GHz

	Time (s)	RAM (GB)	Disk (GB)
Jellyfish (2.2.6)	>15 hours (system hang)		
DSK (2.2.0)	7,722	12	133
DSK (2.2.0; gzip)	9,240	11	134
KAnalyze (2.0.0)	Failed: "IO error writing segment file: no space left on device"		
KAnalyze (2.0.0; gzip)	Failed: "IO error writing segment file: no space left on device"		
KMC3	3,725*	10	78
KMC3 (gzip)	1,964	11	79
Gerbil (1.0)	4,078	6*	66*
Gerbil (1.0; gzip)	2,849	6	66
KCMBT (1.0)	>23 hours		
MSPKmerCounter (0.1)	>15 hours (phase 2 failed: "OutOfMemoryError")		
aTurtle (0.3)	Aborted (core dumped)		

Minimizer -> not equal bins

Sequence analysis

KMC 2: fast and resource-frugal k -mer counting

Sebastian Deorowicz^{1,*}, Marek Kokot¹, Szymon Grabowski² and Agnieszka Debudaj-Grabysz¹

Sequence Analysis

kmtricks: efficient and flexible construction of Bloom filters for large sequencing data collections

Téo Lemane¹, Paul Medvedev^{2,3,4}, Rayan Chikhi⁵ and Pierre Peterlongo^{1,*}

Data Set-Adaptive Minimizer Order Reduces Memory Usage in k -Mer Counting

DAN FLOMIN, DAVID PELLOW, and RON SHAMIR

Sequence analysis

Compact and evenly distributed k -mer binning for genomic sequences

Johan Nyström-Persson^{1,2,*}, Gabriel Keeble-Gagnère³ and Niamat Zawad²

RESEARCH

Open Access



Space-efficient representation of genomic k -mer count tables

Yoshihiro Shibuya¹, Djamel Belazzougui² and Gregory Kucherov^{1,3*}

RESEARCH

Open Access



Gerbil: a fast and memory-efficient k -mer counter with GPU-support

Marius Erbert, Steffen Rechner[†] and Matthias Müller-Hannemann

ORIGINAL RESEARCH ARTICLE

Review paper

Counting Kmers for Biological Sequences at Large Scale

Jianqiu Ge¹ · Jintao Meng¹ · Ning Guo¹ · Yanjie Wei¹ · Pavan Balaji² · Shengzhong Feng¹



Review paper

REVIEW

A benchmark study of k -mer counting methods for high-throughput sequencing

Thank you!

QA?



ETHAN HAWKE

UMA THURMAN

GATTACA

HOW DO YOU HIDE WHEN YOU'RE RUNNING FROM YOURSELF?

COLUMBIA PICTURES PRESENTS A JERSEY FILMS PRODUCTION ANDREO DICCOLL ETHAN HAWKE UMA THURMAN "GATTACA" ALAN ARKIN JOE LUI LOREN DEAN ERNEST BORICONE
MUSIC BY MICHAEL NYMAN COSTUME DESIGNER COLLEEN ATWOOD EDITOR LISA ZERO CHURCHIN PRODUCTION DESIGNER JIM ROBELES EXECUTIVE PRODUCERS SARAUDOMIR IZCHIK PRODUCED BY DANIEL DEWITT MICHAEL SHAMBERG STACEY SHER
www.gattaca.com
COLUMBIA PICTURES JERSEY FILMS ANDREO DICCOLL
DOLBY DIGITAL
DOLBY DIGITAL
DOLBY DIGITAL
DOLBY DIGITAL

Available reviews papers

Rowe *Genome Biology* (2019) 20:199
<https://doi.org/10.1186/s13059-019-1809-x>

Genome Biology

REVIEW

Open Access

When the levee breaks: a practical guide to sketching algorithms for processing the flood of genomic data



Will P. M. Rowe^{1,2}

Tutorial + basics

Annual Review of Biomedical Data Science

Sketching and Sublinear Data Structures in Genomics

Guillaume Marçais,^{1,*} Brad Solomon,^{2,*} Rob Patro,³ and Carl Kingsford¹

Two pages on Minimizer

Time order $O(n)$ ACM Computing Surveys, Vol. 54, No. 1, Article 17. Publication date: March 2021.

Data Structures to Represent a Set of k -long DNA Sequences

RAYAN CHIKHI, Center of Bioinformatics and Biostatistics and Integrative Biology

JAN HOLUB, Department of Theoretical Computer Science, Czech Technical University in Prague

PAUL MEDVEDEV, Center for Computational Biology and Bioinformatics

Review

Data structures based on k -mers for querying large collections of sequencing data sets

Camille Marchet,¹ Christina Boucher,² Simon J. Puglisi,³ Paul Medvedev,^{4,5,6} Mikaël Salson,¹ and Rayan Chikhi⁷



Useful, we can write similarly describing Biological applications.
Good supplementary describing basics

A survey of mapping algorithms in the long-reads era

Intro on seeding, anchor, kmer

4- Sequence compression

SCIENTIFIC
REPORTS

nature research

Corrected: Author Correction

FQsqueezer: k -mer-based compression of sequencing data

Sebastian Deorowicz 

JOURNAL OF COMPUTATIONAL BIOLOGY
Volume 25, Number 7, 2018
© Mary Ann Liebert, Inc.
Pp. 825–836
DOI: 10.1089/cmb.2018.0068

Dynamic Alignment-Free and Reference-Free
Read Compression

PLOS COMPUTATIONAL BIOLOGY

RESEARCH ARTICLE

Hamming-shifting graph of genomic short reads: Efficient construction and its application for compression

Yuansheng Liu , Jinyan Li 

Data Science Institute, University of Technology Sydney, Sydney, Australia

Bioinformatics, 35(12), 2019, 2066–2074
doi: 10.1093/bioinformatics/bty936
Advance Access Publication Date: 8 November 2018
Original Paper

OXFORD

Sequence analysis

Index suffix–prefix overlaps by (w, k) -minimizer to generate long contigs for reads compression

Bioinformatics, 31(9), 2015, 1389–1395
doi: 10.1093/bioinformatics/btu844
Advance Access Publication Date: 22 December 2014
Original Paper

OXFORD

Sequence analysis

Disk-based compression of data from genome sequencing

Szymon Grabowski¹, Sebastian Deorowicz^{2,*} and Łukasz Roguski^{3,4}

Bioinformatics, 34(16), 2018, 2748–2756
doi: 10.1093/bioinformatics/bty205
Advance Access Publication Date: 29 March 2018
Original Paper

OXFORD

Sequence analysis

FaStore: a space-saving solution for raw sequencing data

Łukasz Roguski^{1,2}, Idoia Ochoa³, Mikel Hernandez⁴ and Sebastian Deorowicz^{5,*}

5- Sequence error correction

LaPierre et al. *BMC Bioinformatics* (2019) 20:552
<https://doi.org/10.1186/s12859-019-3103-z>

BMC Bioinformatics

SOFTWARE

Open Access

De novo Nanopore read quality improvement using deep learning



nature communications



Article

<https://doi.org/10.1038/s41467-022-34381-8>

VeChat: correcting errors in long reads using variation graphs

Simultaneous compression of multiple error-corrected short-read sets for faster data transmission and better de novo assemblies



ARTICLE



<https://doi.org/10.1038/s41467-020-20340-8>

OPEN

Error correction enables use of Oxford Nanopore technology for reference-free transcriptome analysis

Kristoffer Sahlin¹ & Paul Medvedev^{2,3,4✉}

Bioinformatics, 37(11), 2021, 1604–1606

doi: 10.1093/bioinformatics/btab915

Advance Access Publication Date: 28 October 2020

Applications Note

OXFORD

Sequence analysis

Minirmd: accurate and fast duplicate removal tool for short reads via multiple minimizers

Yuansheng Liu¹, Xiaocai Zhang², Quan Zou³ and Xiangxiang Zeng^{1,*}

Not a
priority

6- Variant calling & MSA

Dysgu: efficient structural variant calling using short or long reads

Kez Cleal^{1*} and Duncan M. Baird^{1*}

State-of-the-art structural variant calling:
What went conceptually wrong and how
to fix it?

Markus Schmidt¹ and Arne Kutzner^{1,}*

**Mapping-free variant calling using haplotype
reconstruction from k-mer frequencies**

Peter A. Audano*, Shashidhar Ravishankar and Fredrik O. Vannberg*

**FAME: fast and memory efficient multiple sequences
alignment tool through compatible chain of roots**

Etminan Naznooshadat¹, Parvinnia Elham^{1,*} and Sharifi-Zarchi Ali²

7- Read alignment (good reviews exist)

Bioinformatics, 36, 2020, i111–i118
doi: 10.1093/bioinformatics/btaz435
ISMB 2020



Weighted minimizer sampling improves long read mapping

Chirag Jain^{1,*}, Arang Rhie¹, Haowen Zhang², Claudia Chu², Brian P. Walenz¹, Sergey Koren¹ and Adam M. Phillippy¹

nature|methods

ARTICLES

<https://doi.org/10.1038/s41592-022-01457-8>



Long-read mapping to repetitive reference sequences using Winnommap2

Chirag Jain^{1,2,3}, Arang Rhie², Nancy F. Hansen³, Sergey Koren² and Adam M. Phillippy²

Rautialainen and Marschall *Genome Biology* (2020) 21:253
<https://doi.org/10.1186/s13059-020-02157-2>

Genome Biology

SOFTWARE

Open Access

GraphAligner: rapid and versatile sequence-to-graph alignment

Mikko Rautialainen^{1,2,3*} and Tobias Marschall^{4*}



Alser et al. *Genome Biology* (2021) 22:249
<https://doi.org/10.1186/s13059-021-02443-7>

Genome Biology

REVIEW

Open Access

Technology dictates algorithms: recent developments in read alignment



Mohammed Alser^{1,2,3†}, Jeremy Rotman^{4†}, Dhriti Deshpande⁵, Kody Taraszka⁴, Huwenbo Shi^{6,7}, Pelin Icer Baykal⁸, Harry Taeyun Yang^{4,9}, Victor Xue⁴, Sergey Knyazev⁸, Benjamin D. Singer^{10,11,12}, Brunilda Balliu¹³, David Koslicki^{14,15,16}, Pavel Skums⁸, Alex Zelikovsky^{8,17}, Can Alkan^{2,18}, Onur Mutlu^{1,2,3†} and Serghei Mangul^{5†}

A survey of mapping algorithms in the long-reads era

Kristoffer Sahlin

Department of Mathematics, Science for Life Laboratory, Stockholm University, 106 91, Stockholm, Sweden.

Focused on minimizer

Bioinformatics, 34(18), 2018, 3094–3100
doi: 10.1093/bioinformatics/bty191
Advance Access Publication Date: 10 May 2018
Original Paper



Sequence analysis

Minimap2: pairwise alignment for nucleotide sequences

Heng Li*

Department of Medical Population Genetics Program, Broad Institute, Cambridge, MA 02142, USA

Bioinformatics, 32(14), 2016, 2103–2110
doi: 10.1093/bioinformatics/btw152
Advance Access Publication Date: 19 March 2016
Original Paper



Sequence analysis

Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences

Heng Li

Alternative data structures (math+computer sci.)

Masked Minimizers: Unifying sequence sketching methods

A randomized parallel algorithm for efficiently finding near-optimal universal hitting sets

Sequence-specific minimizers via polar sets

Hongyu Zheng, Carl Kingsford and Guillaume Marçais*

Improved design and analysis of practical minimizers

Hongyu Zheng, Carl Kingsford and Guillaume Marçais*

BLight: efficient exact associative structure for k-mers

Camille Marchet  *, Mael Kerbiriou and Antoine Limasset  *



Syncmers are more sensitive than minimizers for selecting conserved k -mers in biological sequences

Robert Edgar

None, Corte Madera, CA, USA

Effective sequence similarity detection with strobemers

Kristoffer Sahlin

A performant bridge between fixed-size and variable-size seeding

Arne Kutzner¹, Pok-Son Kim² and Markus Schmidt^{1*} 

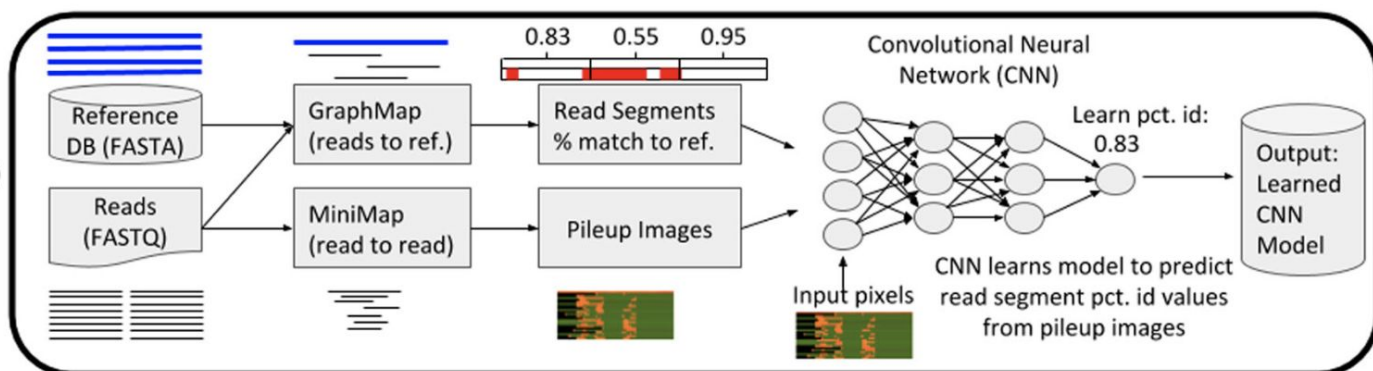
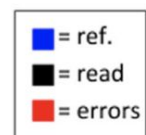
The minimizer Jaccard estimator is biased and inconsistent*

Paper outline “**Applications of Minimizer in genomics**”

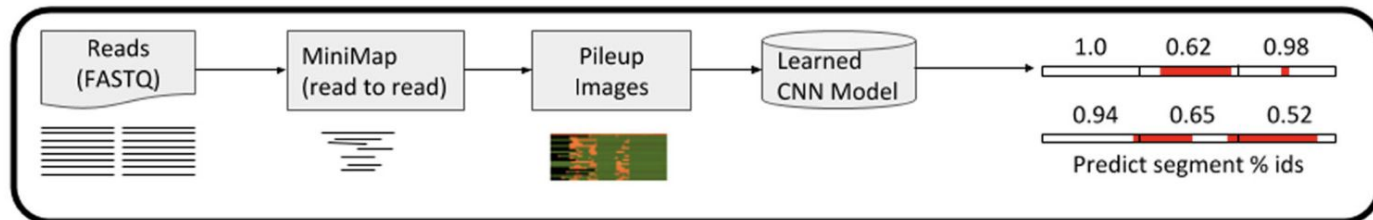
1. intro (importance, why now(biogenome project), definition)
2. metagenomics (7)
3. de Bruijn graph (4)
4. genome assembly (7)
5. kmer counting (5)
6. sequence compression (6)
7. sequence error correction (3)
8. variant calling (3)
9. Discussion and conclusion

Tables, figure.

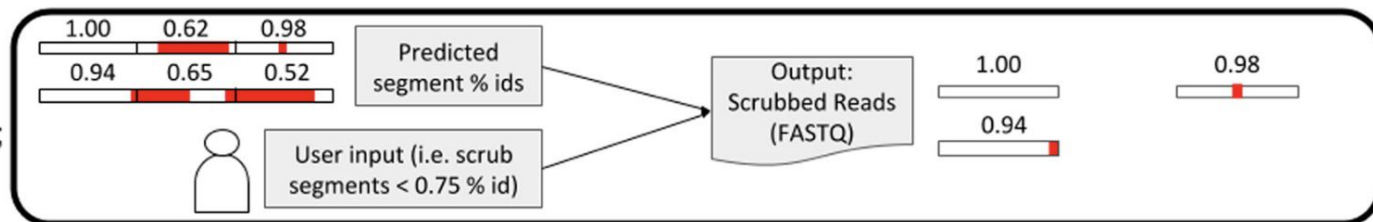
a. Model Training (pre-computed)



b. Predict read segment pct. ids (*de novo*)

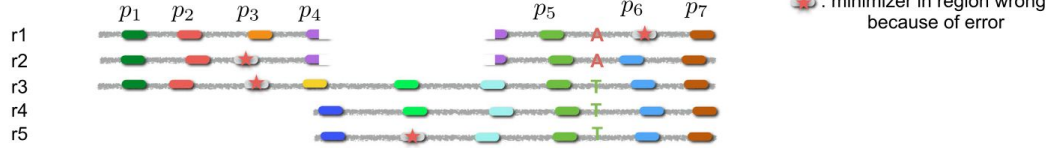


c. Read scrubbing (user side; *de novo*)



isONcorrect method

Reads from cluster instance



Correction of r1

1. Find all intervals for r1 and their weight in the cluster



2. Solve Weighted Interval Scheduling problem



3. Create spoa consensus of each segment



4. Find alternative reference sites based on kmer context



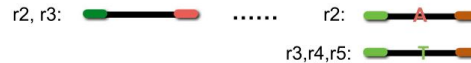
5. Correct r1 in each segment



6. Substitute corrected segments in r



- 7*. Correct and store segments for other participating reads to avoid recomputing this



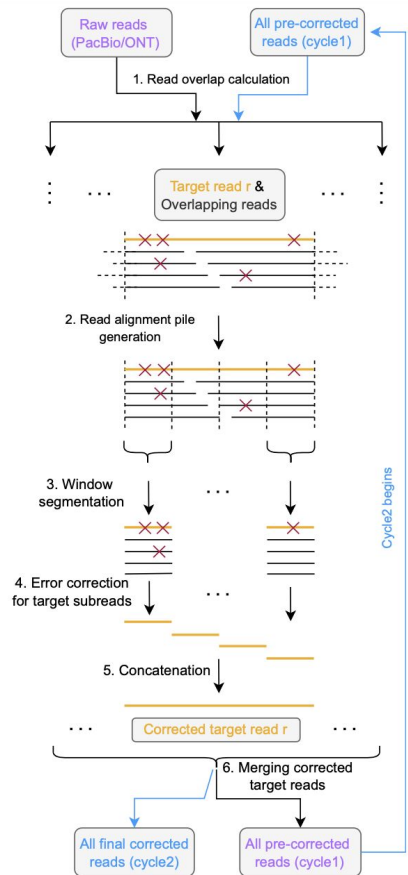


Fig. 1 | Workflow of VeChat. The input and output of cycle 1 and cycle 2 are labeled

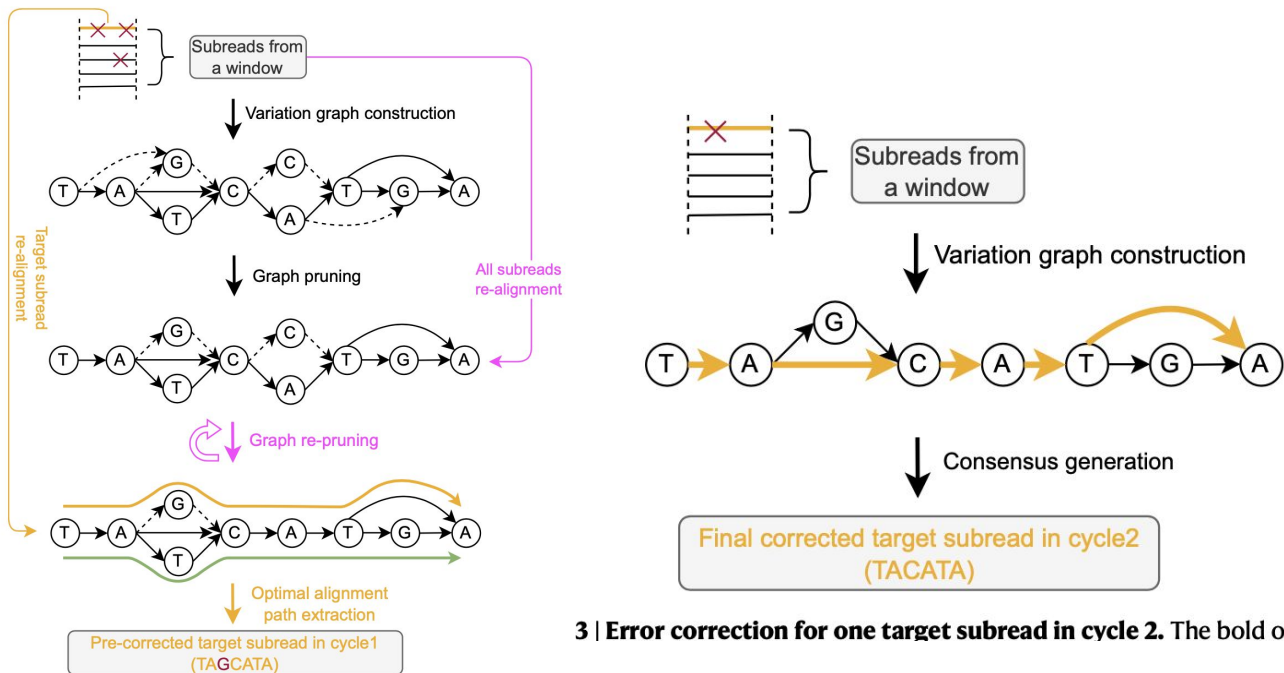
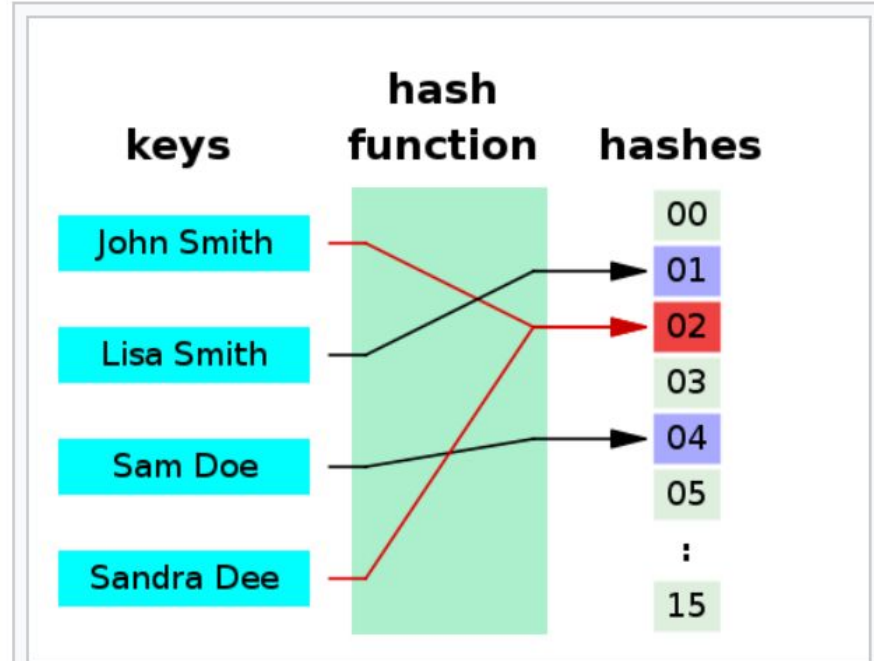


Fig. 2 | Error correction for one target subread in cycle 1. The error correction

3 | Error correction for one target subread in cycle 2. The bold orange

Hash function

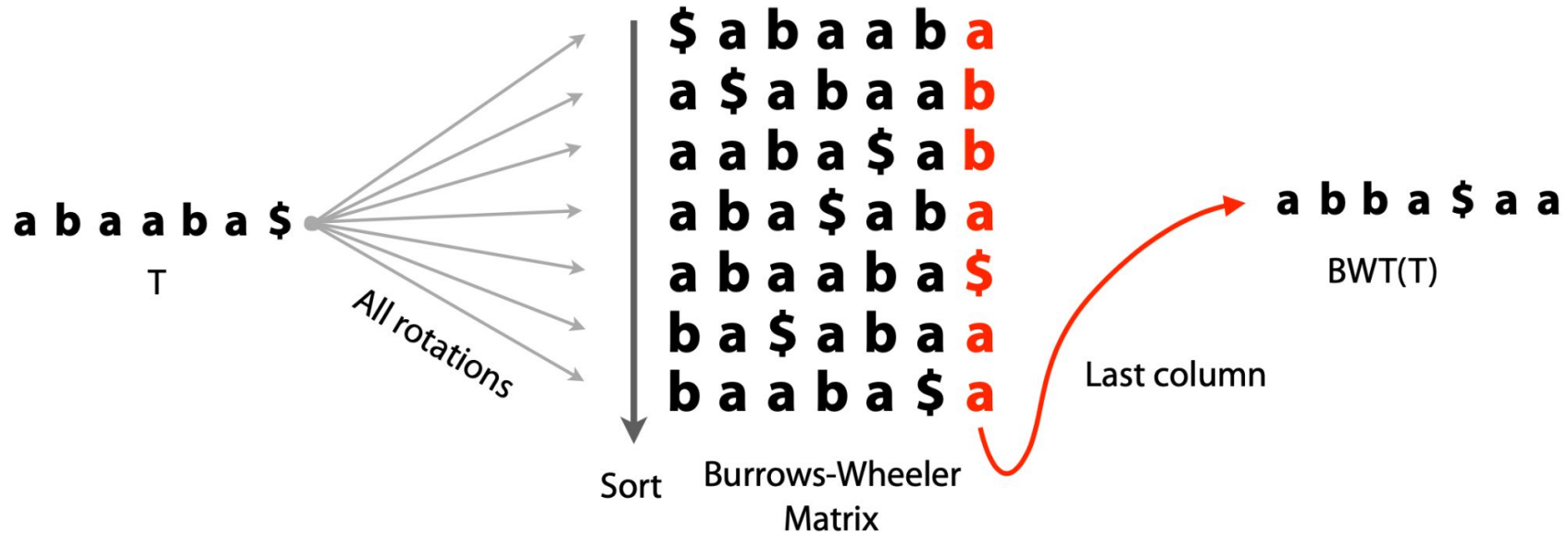
A hash function is any function that can be used to map data of arbitrary size to fixed-size values.



A hash function that maps names to integers from 0 to 15. There is a **collision** between keys "John Smith" and "Sandra Dee".

wikipedia

Burrows–Wheeler transform



Suffix array

\$ a b a a b a
a \$ a b a a b
a a b a \$ a b
a b a \$ a b a
a b a a b a \$
b a \$ a b a a
b a a b a \$ a

BWT

6	\$						
5	a	\$					
2	a	a	b	a	\$		
3	a	b	a	\$			
0	a	b	a	a	b	a	\$
4	b	a	\$				
1	b	a	a	b	a	\$	

FM index

Using the BWT matrix

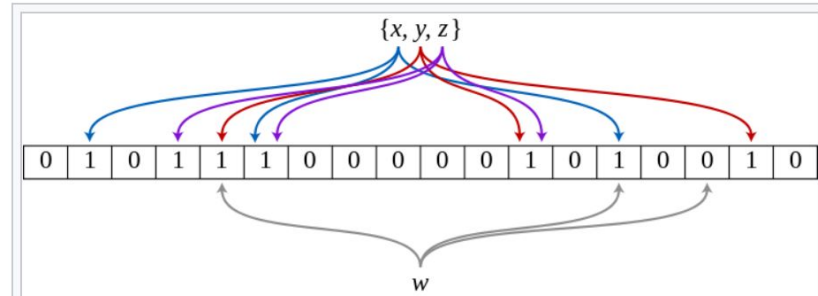
a last-to-first column mapping $LF(i)$ from an index i to $L[i]$

<i>F</i>							<i>L</i>
\$	a	b	a	a	b	a	a
a	\$	a	b	a	a	b	b
a	a	b	a	\$	a	b	b
a	b	a	\$	a	b	a	a
a	b	a	a	b	a	\$	\$
b	a	\$	a	b	a	a	a
b	a	a	b	a	\$	a	a

┌──────────────────┐
Not stored in index

Bloom filter

- a space-efficient probabilistic data structure
- test whether an element is a member of a set
- False positive matches are possible, but false negatives are not
- An empty Bloom filter is a bit array of m bits, all set to 0.
- Using k different hash functions
- each function maps (hashes) an element to one of the m array positions.



An example of a Bloom filter, representing the set $\{x, y, z\}$. \square
The colored arrows show the positions in the bit array that each set element is mapped to. The element w is not in the set $\{x, y, z\}$, because it hashes to one bit-array position containing 0. For this figure, $m = 18$ and $k = 3$.

Bloom filters

independent

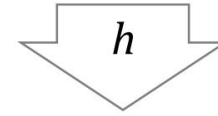
hashfn \emptyset
hashfn1
hashfn2
hashfn3

	\emptyset	1	2	3	
kmer 1	5	12	7	3	
kmer 2	18	9	11	5	
kmer 1	5	12	7	3	TP ✓
kmer 3	18	10	21	16	TN X
kmer 4	11	5	3	18	FP ✓



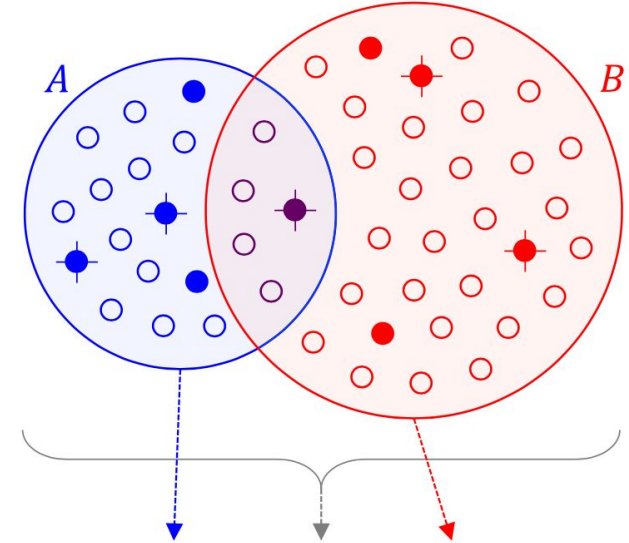
MinHASH

.....
 GGATT AATCG
 TGACG AAGCT
 GTACT GGCAT



$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \approx \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$

Fig. 1 Overview of the MinHash bottom sketch strategy for estimating the Jaccard index. First, the sequences of two datasets are decomposed into their constituent k-mers (*top, blue and red*) and each k-mer is passed through a hash function h to obtain a 32- or 64-bit hash, depending on the input k-mer size. The resulting hash sets, A and B , contain $|A|$ and $|B|$ distinct hashes each (*small circles*). The Jaccard index is simply the fraction of shared hashes (*purple*) out of all distinct hashes in A and B . This can be approximated by considering a much smaller random sample from the union of A and B . MinHash sketches $S(A)$ and $S(B)$ of size $s = 5$ are shown for A and B , comprising the five smallest hash values for each (*filled circles*). Merging $S(A)$ and $S(B)$ to recover the five smallest hash values overall for $A \cup B$ (*crossed circles*) yields $S(A \cup B)$. Because $S(A \cup B)$ is a random sample of $A \cup B$, the fraction of elements in $S(A \cup B)$ that are shared by both $S(A)$ and $S(B)$ is an unbiased estimate of $J(A, B)$.



$S(A)$	$S(A \cup B)$	$S(B)$
42	42	66
64	64	82
82	66	87
128	82	104
139	87	127

StrobeMer

$k=3$, $w=5$, $n=4$ 

strobemer 1:

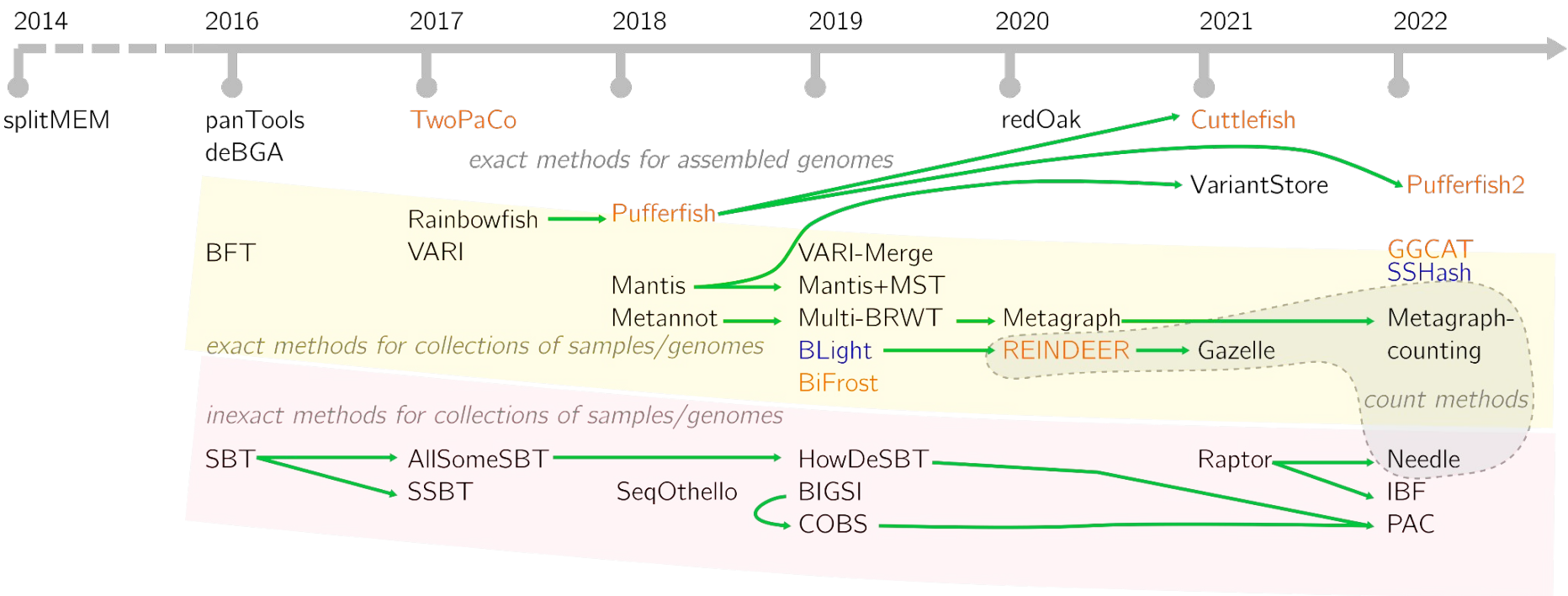
k_1 
ACTCTAAACATTGGCGTCT...

strokes: $k_1=ACT$
 $k_2=AAA$
 $k_3=ATT$
 $k_4=CGT$

strobemer 2:

k_1 
ACTCTAAACATTGGCGTCT...

strokes: $k_1=CTC$
 $k_2=AAA$
 $k_3=ATT$
 $k_4=CGT$



acgt compacted colored de Bruijn graphs
acgt general purpose k -mer hash tables with possible color/feature coding
→ some level of relatedness, either from shared data-structures or inherited ideas,

https://kamimrcht.github.io/webpage/sets_kmer_sets.html

Time and space

count information with each k -mer. Rather than present how this is done, we present a data structure that supports it, we have a separate section (Section 7.4) describing how other data structures can be adapted to store count information.

3 BASIC APPROACHES

Perhaps the most basic static representation that is used in practice is a sorted list of k -mers. The construction time is $O(nk)$ using any linear time algorithm, and the space needed to store the list is $\Theta(nk)$. A membership query is expensive, taking time $O(k \log n)$. This representation is both space- and time-inefficient. Other approaches we will discuss (e.g., unitig-based approaches or BOSS). But with very limited computer science background, making it still relevant.

Sorted lists can be partitioned to speed up queries. In this approach, [2014], the k -mers are partitioned according to a minimizer function. The minimizer of a k -mer x is the smallest (according to some given p

Swap and compare kmer counting

Example:

workers **W1** & **W2**.

W1: read count of AAT (0) W1: increment of AAT 0->1

W2: read count of ATT (0) W2: increment of ATT 0->1

AAT	1
ATT	1

W1:

W2: read count AAT (1)

W1: read count AAT (1)

W2: increment AAT 1->2

AAT	2
ATT	1

W1: increment AAT 1->2

W2:

AAT	2
ATT	1

W1, No, I won't increment AAT as 1->2 does not hold

to ensure low random access memory (RAM) usage, as has been argued by Li et al. (2013).

In this study, we introduce a new approach for adapting the minimizer order to the target sequence data. The method, called Adaptive Order (AdaOrder), iteratively updates the minimizer order based on an estimate of the minimizer loads in the data to reduce the maximum load. We demonstrated its ability to lower the maximum load compared with all predefined orders, except the frequency-based order.

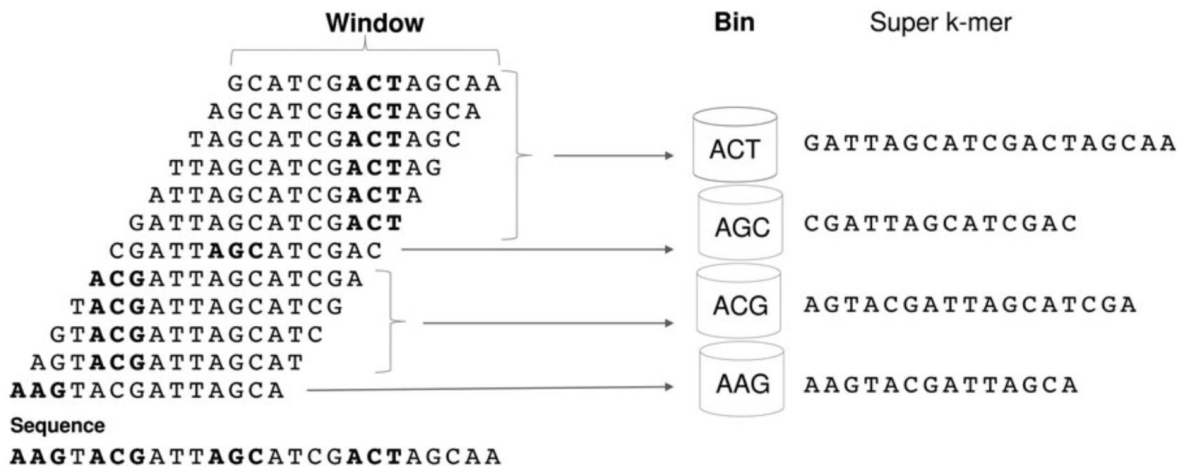


FIG. 1. Illustration of a minimizer scheme and binning. Here $k=12$ and $m=3$. The input sequence is broken into windows of length k , and in each window the m -long minimizer according to the lexicographic order (shown in bold) is selected. The k -mer is assigned to the partition whose label is the minimizer. Consecutive windows tend to select the same minimizer, and the concatenation of the consecutive windows forms the super- k -mer that is stored in the bin.

Sequence comparison with minimizer

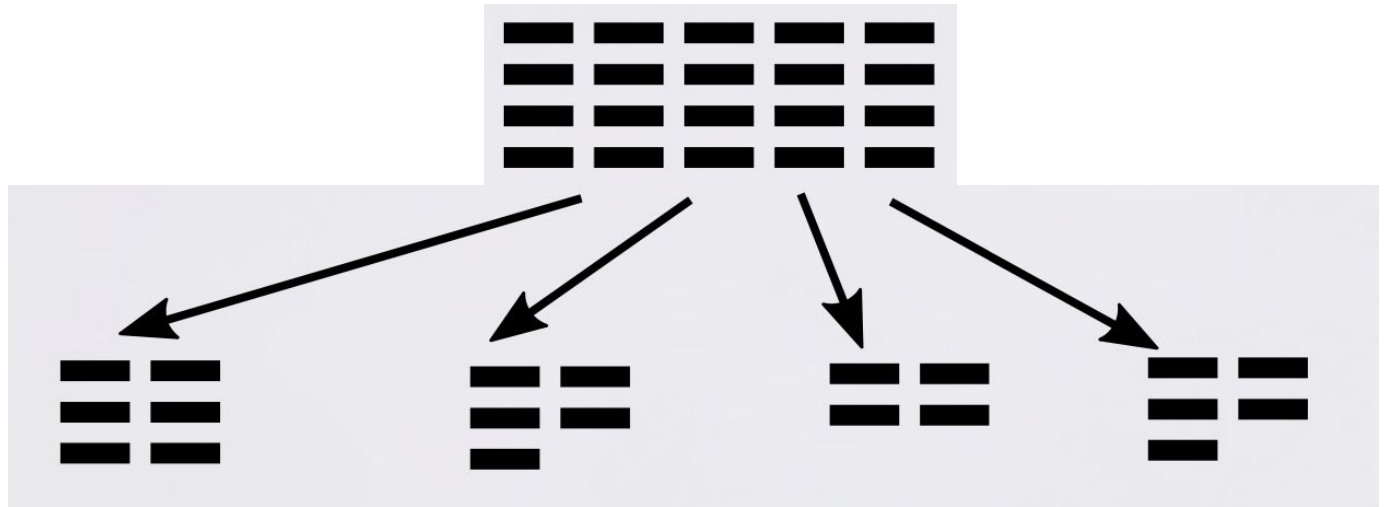
- Number of minimizer at most = $L-w+1$ but
- Similar sequences have high probability of having the same minimizers.

- Sequence of minimizers is also called a sketch of the original seq.

- With non-overlapping windows $\rightarrow L/w$ minimizers for sequence of length L .
- For a random order (w not too large), a factor of $2/(w + 1)$.

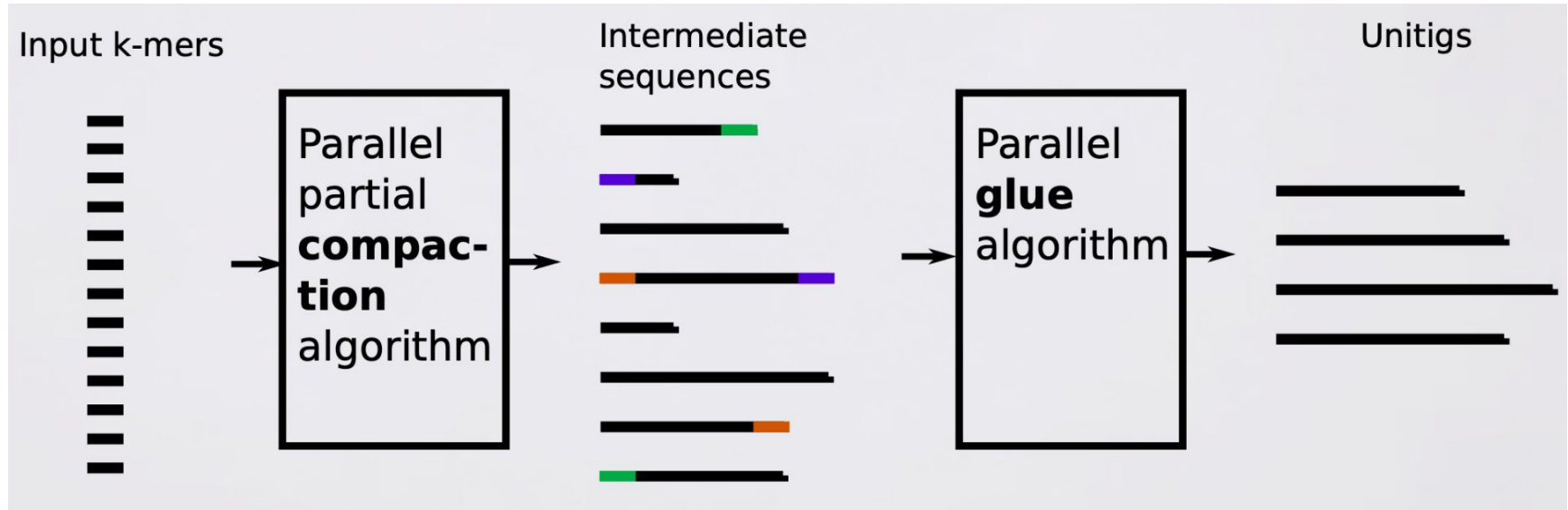
Memory and speed; minimizer

- partitioning input k-mers on disk
- based on minimizers

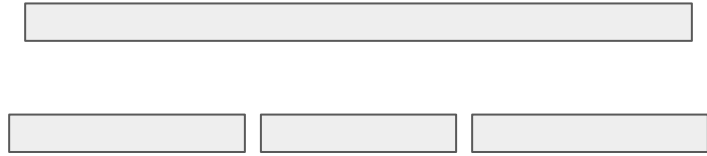


BCALM 2

A tool for compacting the de Bruijn graph



Parallel counting: splitting the sequence



Issues: missing kmers, combining stats of each part.

Gerbil: with minimizer

- 4-mers
- 3-mer minimizer (bold)
- Finding (overlapping) part of sequence with the same minimizer on temp files



CAAGAACAGTG

CAAG

1. CAAGA

AAGA

AGAA

2. AGAA

GAAC

3. GAACA

AACA

ACAG

4. ACAG

CAGT

5. CAGTG

AGTG