

Metode sekvenciranja u preciznoj medicini

izv. prof. Mile Šikić
Sveučilište u Zagrebu,
Fakultet elektrotehnike i računarstva
Laboratorij za bioinformatiku i računalnu biologiju
Bioinformatics Institute, A*STAR, Singapore

Neprecizna medicina

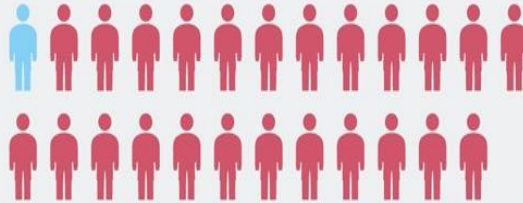
IMPRECISION MEDICINE

For every person they do help (blue), the ten highest-grossing drugs in the United States fail to improve the conditions of between 3 and 24 people (red).

1. ABILIFY (aripiprazole)
Schizophrenia



2. NEXIUM (esomeprazole)
Heartburn



3. HUMIRA (adalimumab)
Arthritis



4. CRESTOR (rosuvastatin)
High cholesterol



5. CYMBALTA (duloxetine)
Depression



6. ADVAIR DISKUS (fluticasone propionate)
Asthma



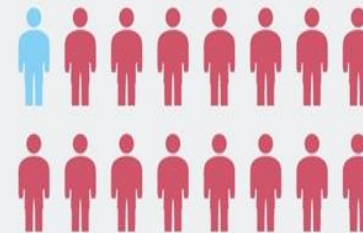
7. ENBREL (etanercept)
Psoriasis



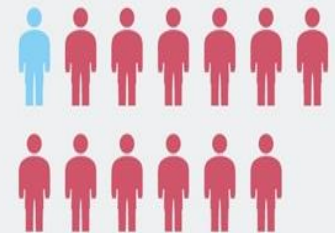
8. REMICADE (infliximab)
Crohn's disease



9. COPAXONE (glatiramer acetate)
Multiple sclerosis



10. NEULASTA (pegfilgrastim)
Neutropenia

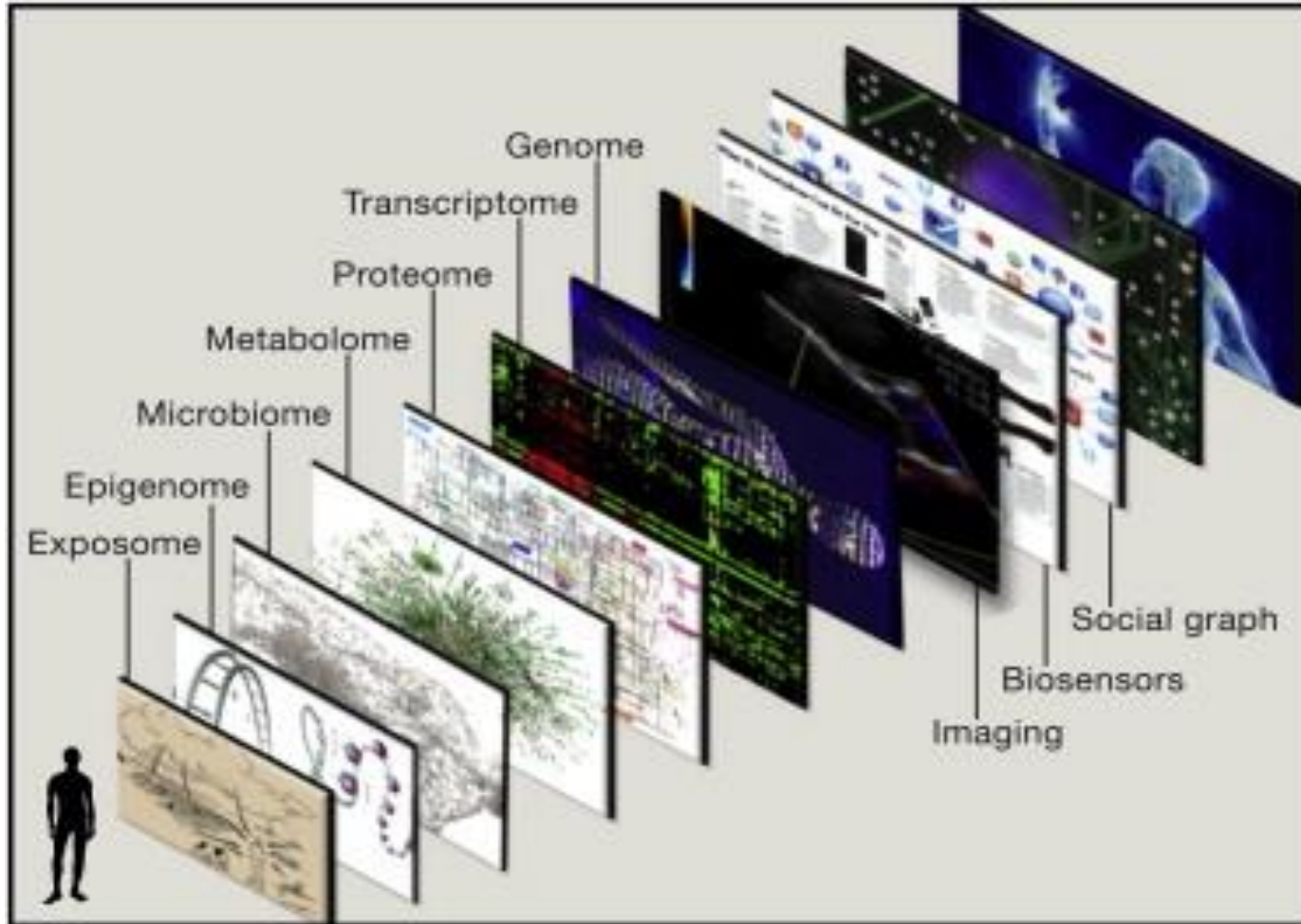


Medicina danas (US)

- >60% lažno pozitivnih nalaza u masovnim pregledima*
- 1 od 4 pacijenta su se dodatno oboljeli u bolnici
- >12 milijuna ozbiljnih dijagnostičkih pogrešaka
- Medicinske pogreške i nuspojave su 4-ti vodeći uzročnik smrti
- U ~80% pacijenata vodeći prepisani lijekove ne djeluju

*Godišnja mamografija u 10 godina

Precizna medicina



Topol, E. Individualized Medicine from Prewomb to Tomb, *Cell*, Volume 157, Issue 1, p241–253, 27 March 2014

Inicijativa za preciznu medicinu (\$215 million)

Mission statement:

To enable a new era of medicine through research, technology, and policies that empower patients, researchers, and providers to work together toward development of individualized care.



Stara nasuprot nove medicine

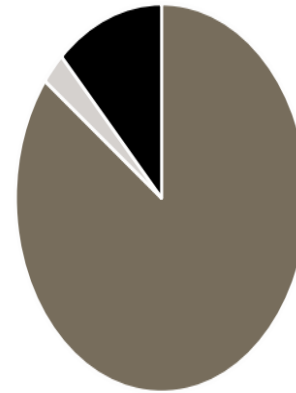
Old Medicine	New Medicine
Population-Based	Individualized
One-Off, Doctor's Office	Real-Time Streaming, Real World
Doctor Ordered Data	Patient Generated Data
Doctor's Notes, Unshared	Our Notes, Patient Edited
Information Owned by Doctors and Hospitals	Information Owned by Rightful Owner
Expensive, Big-Ticket Tech	Cheap Chips, Moore's Law
Data Limited	Panoromic

Author: Eric Topol

Učinkovita dijagnoza neuroleptospiroze koristeći sekvenciranje

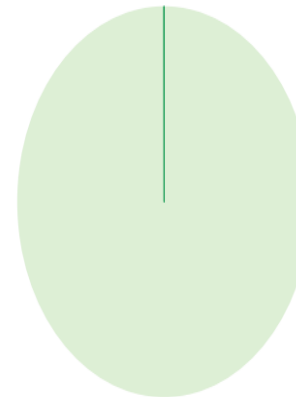


Joshua Osborne, 14 godina



Bacterial (N=589 reads)

- Leptospira (N=475; 80.6%)
- Propionibacterium (N=15; 2.5%)
- Other bacteria (N=99; 16.8%)



Viral (N=107,016 reads)

- Anelloviridae (N=106,988; 99.97%)
- Bacteriophage (N=28; 0.03%)

The New York Times



The NEW ENGLAND
JOURNAL of MEDICINE

Projekt određivanja ljudskog genoma



Određivanje ljudskog genoma u brojkama

- 3 milijarde nukleotida
- 3 milijarde USD
- 13 godina
- 6 uključenih država (US, UK, Francuska, Njemačka, Japan, Kina)
- 20 ustanova

DNA sekvenciranje



DNA sekvenciranje



DNA sekvenciranje



DNA sekvenciranje

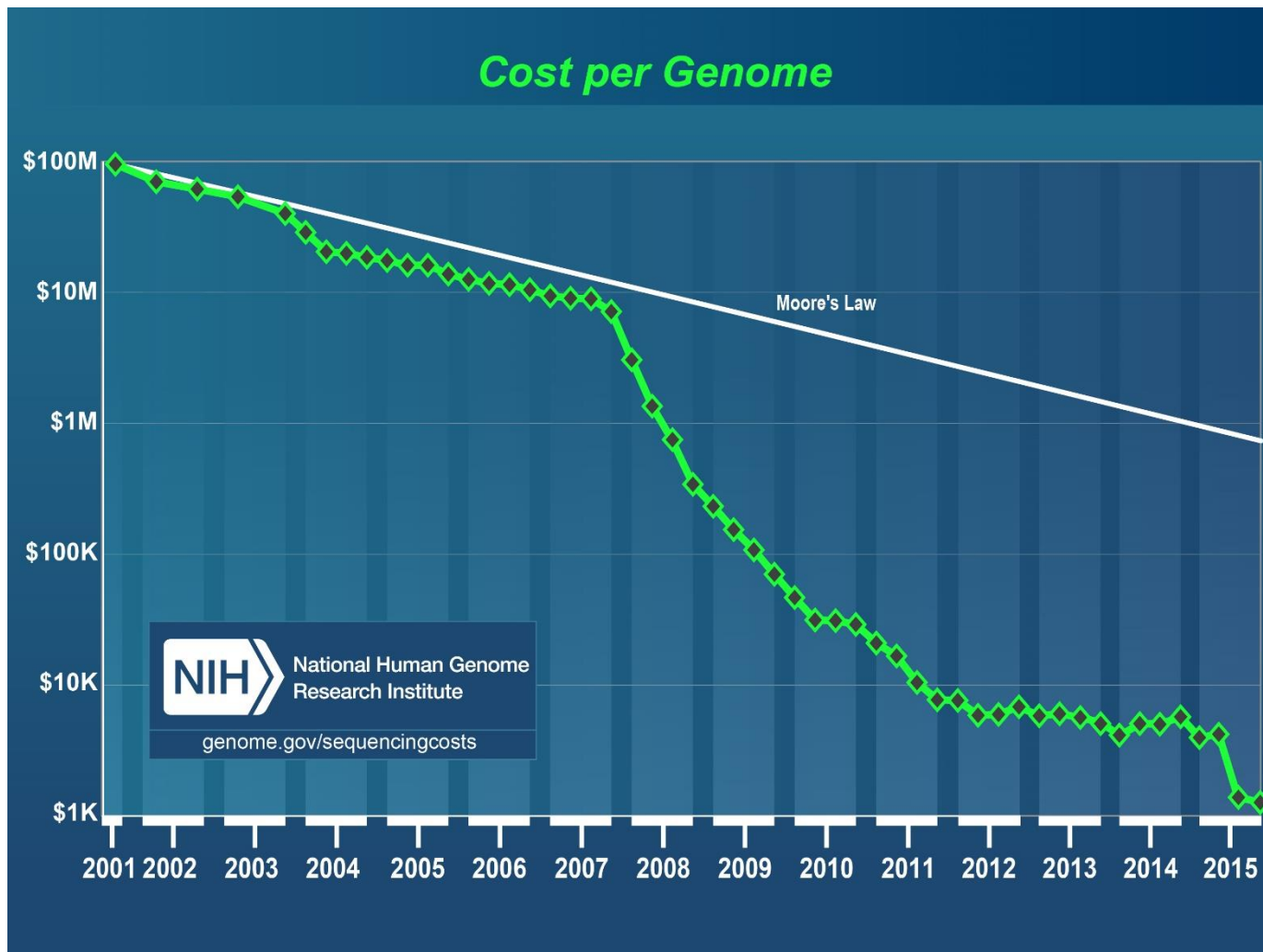


Projekt sekvenciranja ljudskog genoma



Izvor: Wikipedia

Troškovi sekvenciranja

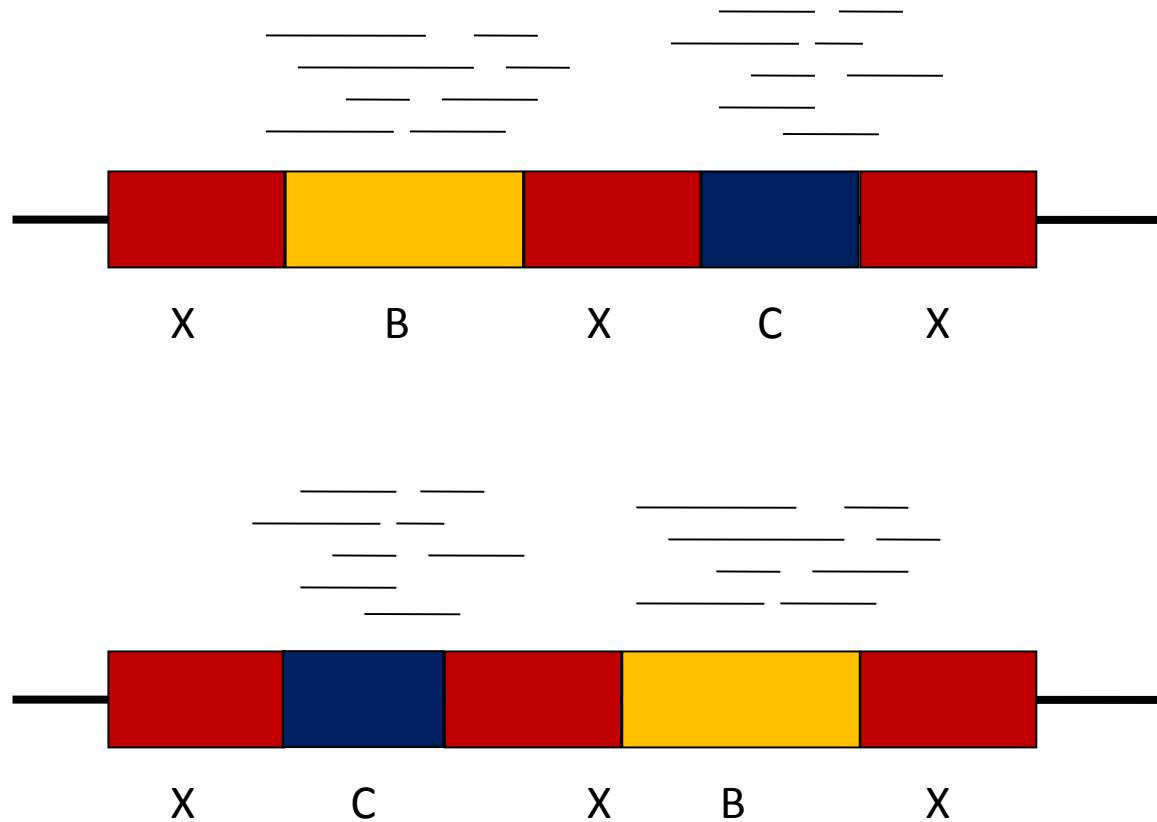


Druga generacija sekvenciranja



Izvor: Illumina

Problemi u sastavljanju genoma (ponavljanja)

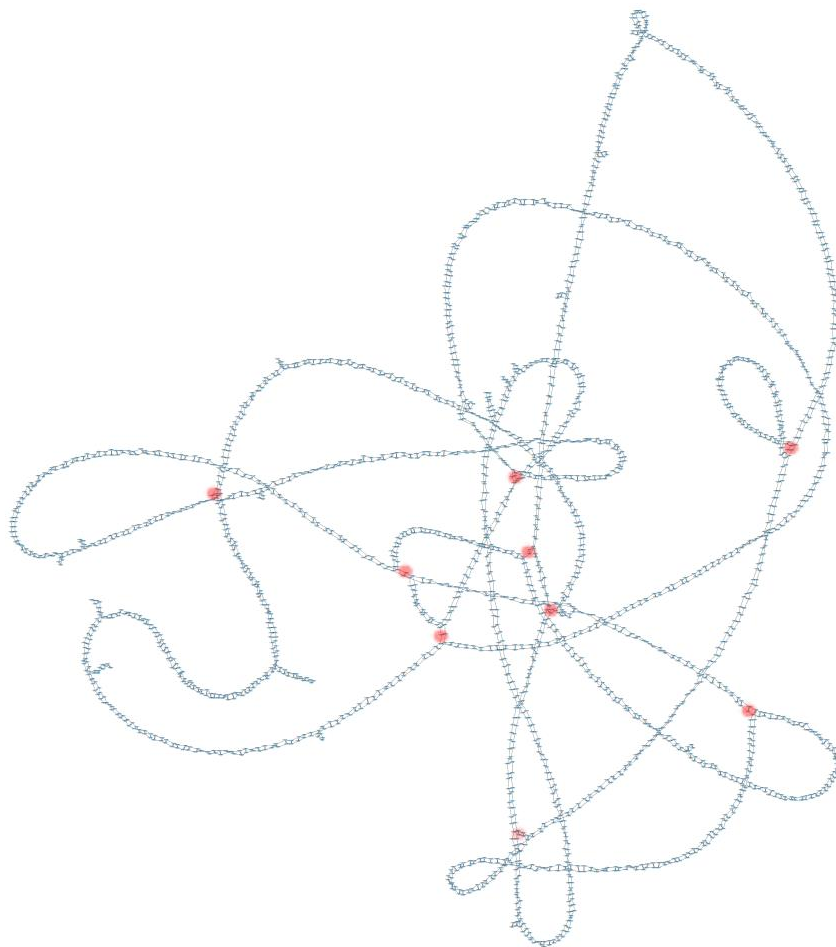


Problemi u sastavljanju (kimerna očitavanja)



AGACGACTTTAGATACTGGGTACTAGAACCCTCAGG

Problemi u sastavljanju



Dugačka očitavanja

Pacific Biosciences Sequel



Oxford nanopore technologies MinION MK I



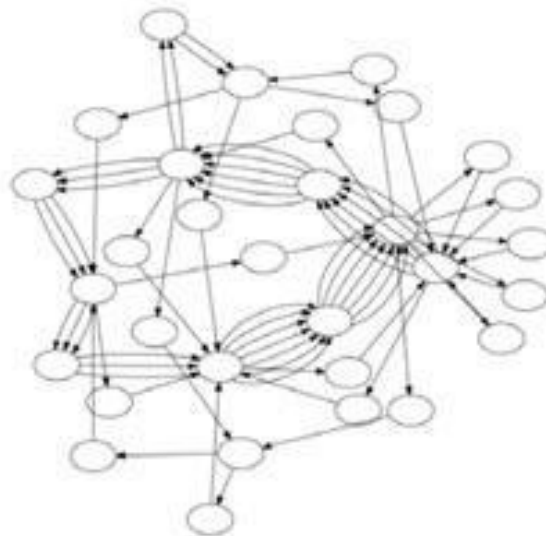
Dugačka očitavanja otpetljavaju graf

(a)



K=100
Contigs=98

(b)



K=1,000
Contigs=31

(c)



K=5,000
Contigs=1

PromethION

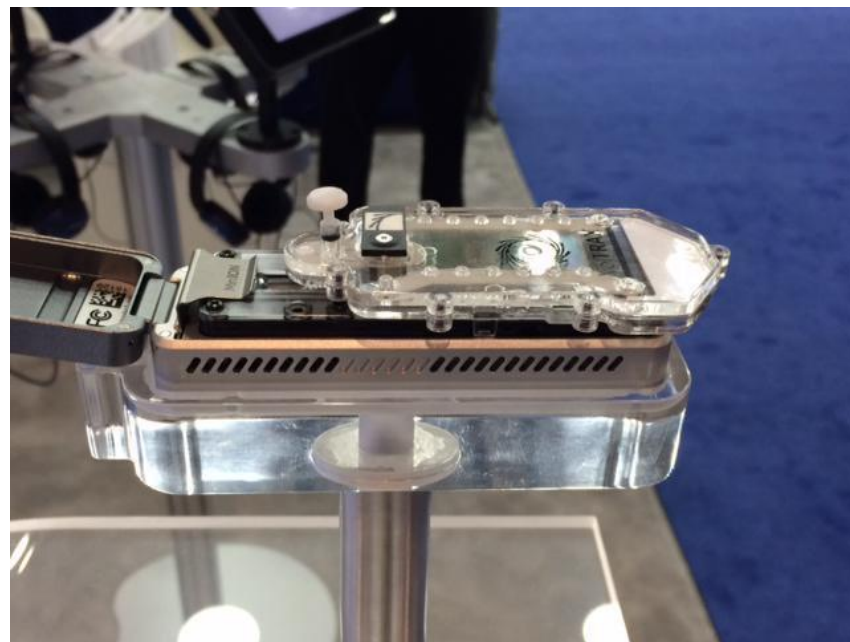
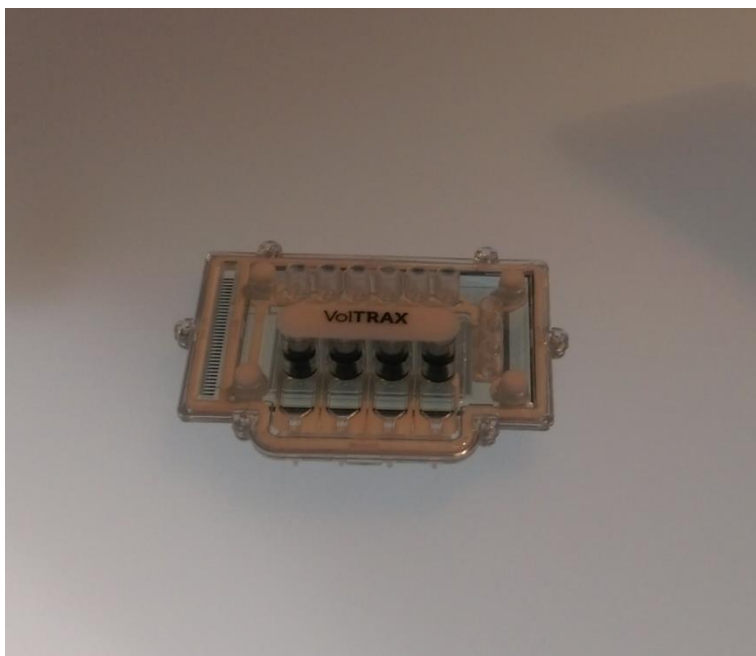


- Novija generacija uređaja za sekvenciranje
- Puno veća propusnost

Pripremanje uzoraka



Pripremanje uzorka - Voltrax



Izvor: Oxford Nanopore Technologies

SmidgION



Detekcija ebole



Photograph: Tommy Trenchard/EMLabs

- Pomoć u identificiranju mutacija u realnom vremenu
- Bez potrebe za prebacivanje uzoraka krvi u visoko kontrolirane laboratorije
- Rezultati unutar 24 sata
- J. Quick et al. Real-time, portable genome sequencing for Ebola surveillance, *Nature*, February 2016

Temeljni računalni problemi

- Sastavljanje genoma – genom za vrstu unaprijed poznat (mapiranje)
- Sastavljanje genome – genom za vrstu nepoznat (*de novo* sastavljanje)
- Sastavljanje kada je u genomu prisutno nekoliko uzoraka
- Određivanje prisutnih RNA molekula

Mapiranje

>1:866511 in NA12878/NA12878.bam

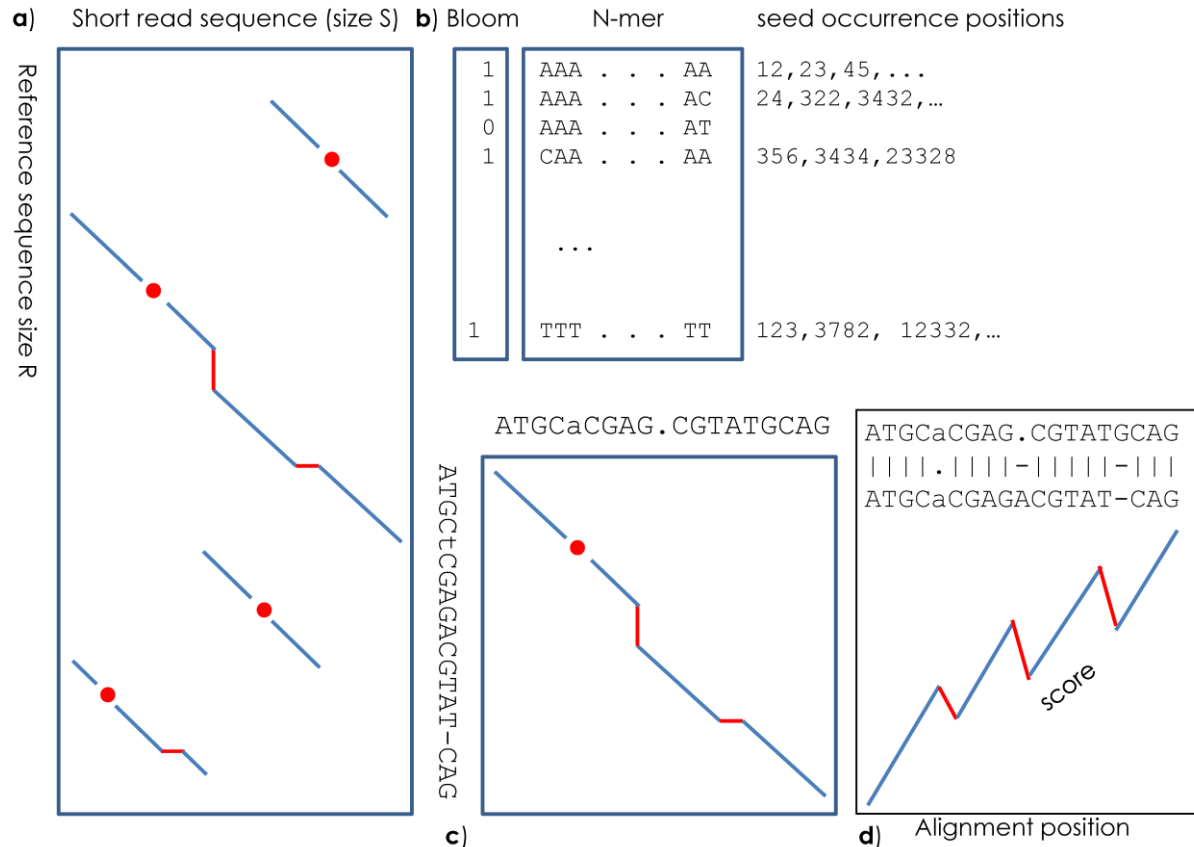
Ref: TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCTGTCTGCTGTCGCGTGTCTCAGCGTGAGC

60> TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----C
60< TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----C
60> TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CC
60< TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCC
60< TcGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CC
60> TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----cCC
60< TCCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCCT----C
60< CCGAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCCT----CC
60> GAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCct----cCC
70< GAGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT
70< gagaGGCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT
60> AGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCctCCCT----CCct----cccc
70< AGAGCCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----C
70> GAGGCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CC
70< GAGGCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CC
60< GGCCTCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCCT----CCC
60> TCCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCct----CCC
60> TCCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCct----CCC
60< TCCTGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCCT----CCC
60< ccTGCAGGTAGGAGCCGTGcTGTGCGTGCATAAAGAGGGGGCCGTGACTcCCCTCCCTCCCT----CCCT----CCCACCCCT
60> TGCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTga
60> GCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTC----CCCTCCCT----CCCT----cCC
60< GCAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGAC
60< GCAGgtAGGAGCCGTGCTgCCTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGAC
60< CAGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACC
60< AGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCG
60< AGGTAGGAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCG
60< GAGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCT
60< gaGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCT
60< gaGCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCT
60> GCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCctCCCTCCCT----CCct----CCCACCCCTGACCGtgccctgc
60< GCCGTGCTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGC
60< CCGTGTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCT
60< tgCTGTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCTGTG
60< tgctgtGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCTGTG
60< GCTGTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCaCCCTGACCGTGCCCTGCTGTCT
60< TGTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCTGTGCG
60< GTGCGTGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCTGTCTGT
60< tgcgtgcATAAAGAGGGGGcCgtgaCTCCCCTCCCTCCCT----cCCT----CCCaCCCTGACCGTGCCCTGCTGTCTGTG
60< TGCATAAAGAGGGGGCCGTGACTCCCCTCCCTCCCT----CCCT----CCCACCCCTGACCGTGCCCTGCTGTCTGTGCGC
60> AGAGGGGGCCGTGACTC----CCCTCCCT----CCCTCCCccccccccctgacctgacctgctgtctggtgtccccctggct
60< gaGGGGGcGTGACTCCCctCCCTCCct----CCct----CCCACCCCTGACCGTGCCCTGCTGTCTGCTGCGCTGCTCT

Mapiranje – optimalno (ljudski genom 3 milijarde nukleotida)

Algoritam	Procesor	Vrijeme
SW or NW	CPU i7	2000 godina
SW or NW	1000 * CPU i7	2 godine
SW or NW	GPU	60 godina
SW or NW	100 * GPU	7 mjeseci

Mapiranje postupak (heuristika)



Complexity and memory requirements for a single Dynamic matrix computation are $O(S \cdot R)$

LCSk



Cornell University
Library

arXiv.org > cs > arXiv:1407.2407

Search or Article ID inside arXiv

All papers



Broaden your search

[Help](#) | [Advanced search](#)

Computer Science > Data Structures and Algorithms

LCSk++: Practical similarity metric for long strings

Filip Pavetić, Goran Žužić, Mile Šikić

(Submitted on 9 Jul 2014)

In this paper we present *LCSk++*: a new metric for measuring the similarity of long strings, and provide an algorithm for its efficient computation. With ever increasing size of strings occurring in practice, e.g. large genomes of plants and animals, classic algorithms such as Longest Common Subsequence (LCS) fail due to demanding computational complexity. Recently, Benson et al. defined a similarity metric named *LCSk*. By relaxing the requirement that the k -length substrings should not overlap, we extend their definition into a new metric. An efficient algorithm is presented which computes *LCSk++* with complexity of $O((|X| + |Y|) \log(|X| + |Y|))$ for strings X and Y under a realistic random model. The algorithm has been designed with implementation simplicity in mind. Additionally, we describe how it can be adjusted to compute *LCSk* as well, which gives an improvement of the $O(|X| |Y|)$ algorithm presented in the original *LCSk* paper.

Edlib

OXFORD
ACADEMIC

Bioinformatics

Issues Advance Articles Publish ▾ Purchase Alerts About ▾

Article Contents

- 1 Introduction
- 2 Methods
- 3 Results
- Acknowledgements
- References
- Author notes
- Supplementary data
- Comments (0)

CORRECTED PROOF

Edlib: a C/C++ library for fast, exact sequence alignment using edit distance

Martin Šošić; Mile Šikić 

Bioinformatics btw753. DOI: <https://doi.org/10.1093/bioinformatics/btw753>

Published: 31 January 2017 Article history ▾

 Views ▾  PDF  Cite  Share ▾  Tools ▾

Summary: We present Edlib, an open-source C/C++ library for exact pairwise sequence alignment using edit distance. We compare Edlib to other libraries and show that it is the fastest while not lacking in functionality and can also easily handle very large sequences. Being easy to use, flexible, fast and low on memory usage, we expect it to be easily adopted as a building block for future bioinformatics tools.

Availability and Implementation: Source code, installation instructions and test data are freely available for download at <https://github.com/Martinos/edlib>, under the MIT licence. Edlib is implemented in C/C++ and supported on Linux, MS Windows, and MacOS.

Contact: mile.sikic@fer.hr

Supplementary information:
Supplementary data
are available at *Bioinformatics* online.

Issue Section: APPLICATIONS NOTE

GraphMap

nature.com > nature communications > articles > article

MENU ▾



Altmetric: 68 Views: 8,104 Citations: 11

[More detail >>](#)

Article | [OPEN](#)

Fast and sensitive mapping of nanopore sequencing reads with GraphMap

Ivan Sović, Mile Šikić, Andreas Wilm, Shannon Nicole Fenlon, Swaine Chen & Niranjan Nagarajan 

Nature Communications **7**,
Article number: 11307 (2016)
doi:10.1038/ncomms11307

[Download Citation](#)

[Bioinformatics](#) [DNA sequencing](#)

[Molecular biology](#)

Received: 30 December 2015

Accepted: 11 March 2016

Published online: 15 April 2016



De novo sastavljanje (1000 CPU sati)

nature **methods**
Techniques for life scientists and chemists

[Home](#) | [Current issue](#) | [Comment](#) | [Research](#) | [Archive](#) ▾ | [Authors & referees](#) ▾ | [About the journal](#) ▾

[home](#) ▶ [archive](#) ▶ [issue](#) ▶ [brief communication](#) ▶ [abstract](#)

ARTICLE PREVIEW
[view full access options](#) ▶





NATURE METHODS | **BRIEF COMMUNICATION**  

A complete bacterial genome assembled *de novo* using only nanopore sequencing data

Nicholas J Loman, Joshua Quick & Jared T Simpson

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)

Nature Methods **12**, 733–735 (2015) | doi:10.1038/nmeth.3444
Received 11 March 2015 | Accepted 22 May 2015 | Published online 15 June 2015

 [Citation](#)  [Reprints](#)  [Rights & permissions](#)  [Article metrics](#)

We have assembled *de novo* the *Escherichia coli* K-12 MG1655 chromosome in a single 4.6-Mb contig using only nanopore data. Our method has three stages: (i) overlaps are detected between reads and then corrected by a multiple-alignment process; (ii) corrected reads are assembled using the Celera Assembler; and (iii) the assembly is polished using a probabilistic model of the signal-level data. The assembly reconstructs gene order and has 99.5% nucleotide identity.

Može brže

OXFORD
ACADEMIC

Bioinformatics

Issues

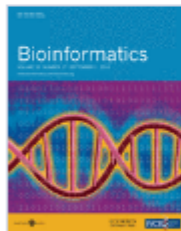
Advance Articles

Publish ▼

Purchase

Alerts

About ▼



Volume 32, Issue 17

1 September 2016

Article Contents

Evaluation of hybrid and non-hybrid methods for de novo assembly of nanopore reads

Ivan Sović; Krešimir Križanović; Karolj Skala; Mile Šikić ✉

Bioinformatics (2016) 32 (17): 2582-2589.

DOI: <https://doi.org/10.1093/bioinformatics/btw237>

Published: 09 May 2016 **Article history** ▼

“ Cite

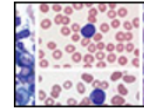
↻ Share ▼

🔧 Tools ▼

Racon



GENOME
RESEARCH



Second AACR Conference on
Hematologic Malignancies
May 6-9 | Boston, MA [REGISTER TODAY!](#)

[HOME](#) | [ABOUT](#) | [ARCHIVE](#) | [SUBMIT](#) | [SUBSCRIBE](#) | [ADVERTISE](#) | [AUTHOR INFO](#) | [CONTACT](#) | [HELP](#)

Fast and accurate de novo genome assembly from long uncorrected reads

Robert Vaser¹, Ivan Sovic², Niranjan Nagarajan³ and Mile Sikic^{1,4}

[+](#) Author Affiliations

[↵](#)* Corresponding author; email: mile.sikic@fer.hr

Abstract

The assembly of long reads from Pacific Biosciences and Oxford Nanopore Technologies typically requires resource intensive error correction and consensus generation steps to obtain high quality assemblies. We show that the error correction step can be omitted and high quality consensus sequences can be generated efficiently with a SIMD accelerated, partial order alignment based stand-alone consensus module called Racon. Based on tests with PacBio and Oxford Nanopore datasets we show that Racon coupled with Miniasm enables consensus genomes with similar or better quality than state-of-the-art methods while being an order of magnitude faster.

ACCEPTED MANUSCRIPT

This Article

Published in Advance January 18, 2017, doi: 10.1101/gr.214270.116

Genome Res. 2017.

Published by Cold Spring Harbor Laboratory Press

- » Abstract *Free*
- » Full Text (PDF)

- Article Category

Method

+ Services

[+](#) Google Scholar

[+](#) PubMed/NCBI

[+](#) ORCID

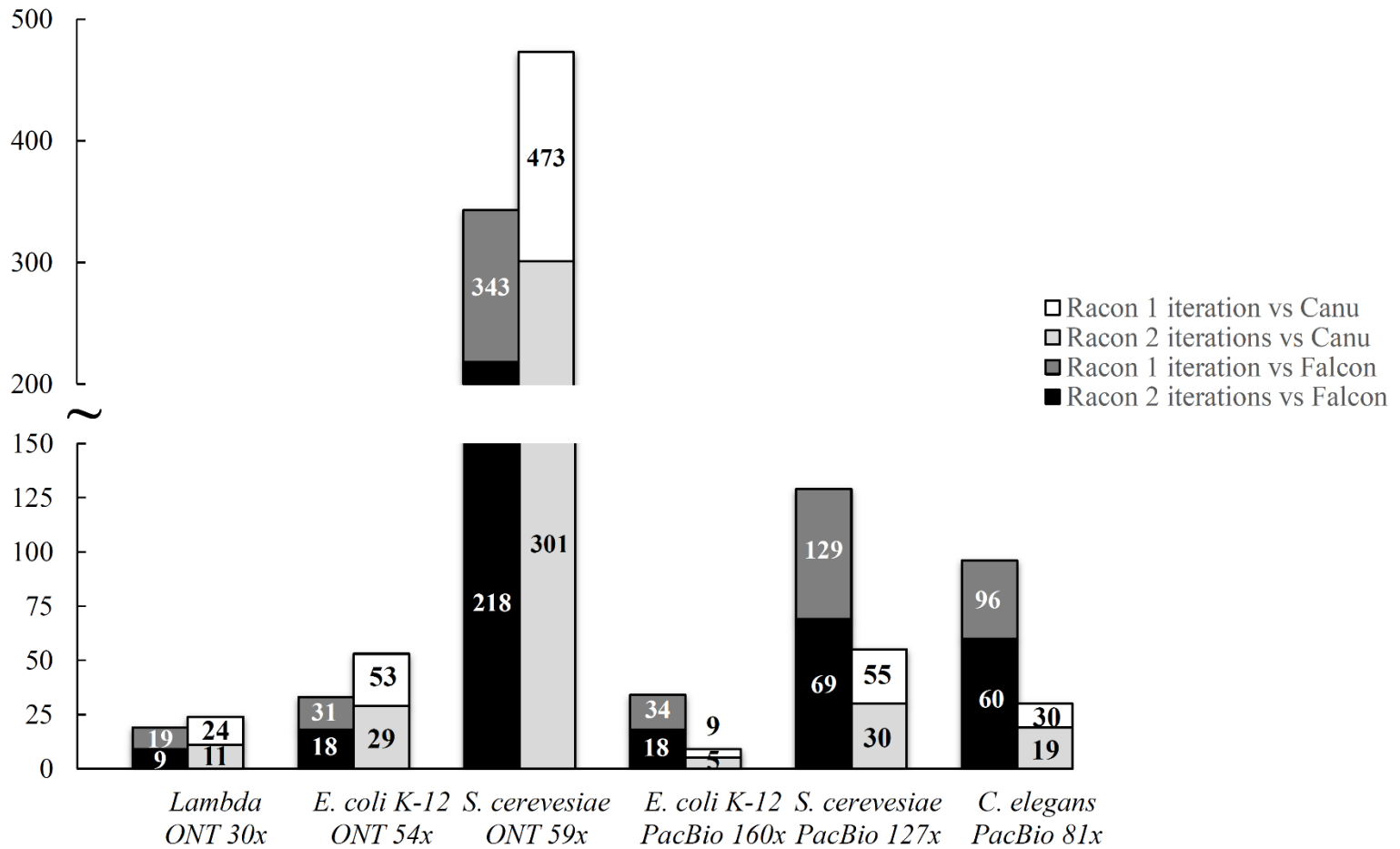
[+](#) Share

[-](#) Metrics

Total Downloads



Racon



Primjena - dijagnoza



“Naši” mikrobi



Genomi biljaka



Gdje smo sada?

- Genomi veličine bakterija – uspješno
- Eukarioti – fragmentirani genomi
- Ploidnost
- Metagenomi (de novo)

Laboratorij za bioinformatiku i računalnu biologiju



<https://www.facebook.com/ferlbcb/>

Laboratorij za bioinformatiku i računalnu biologiju

- 1 poslijedoktorand
- 5 doktoranda
- 15 diplomskih studenata u području računarske znanosti
- Alumni
 - Carnegie Mellon, EPFL, ETH
 - Google, Facebook, Microsoft, Amazon,...

Projekti

- Algorithms for genome sequence analysis
- Utvrđivanje patogena iz niza sekvenciranih RNA podataka
- Utvrđivanje patogena iz niza sekvenciranih DNA podataka



ZAKLADADRIS

Hvala vam!
Mile.Sikic@FER.hr