

ALGORITMISKĀS METODES GENOMA VIRKŅU ANALĪZEI

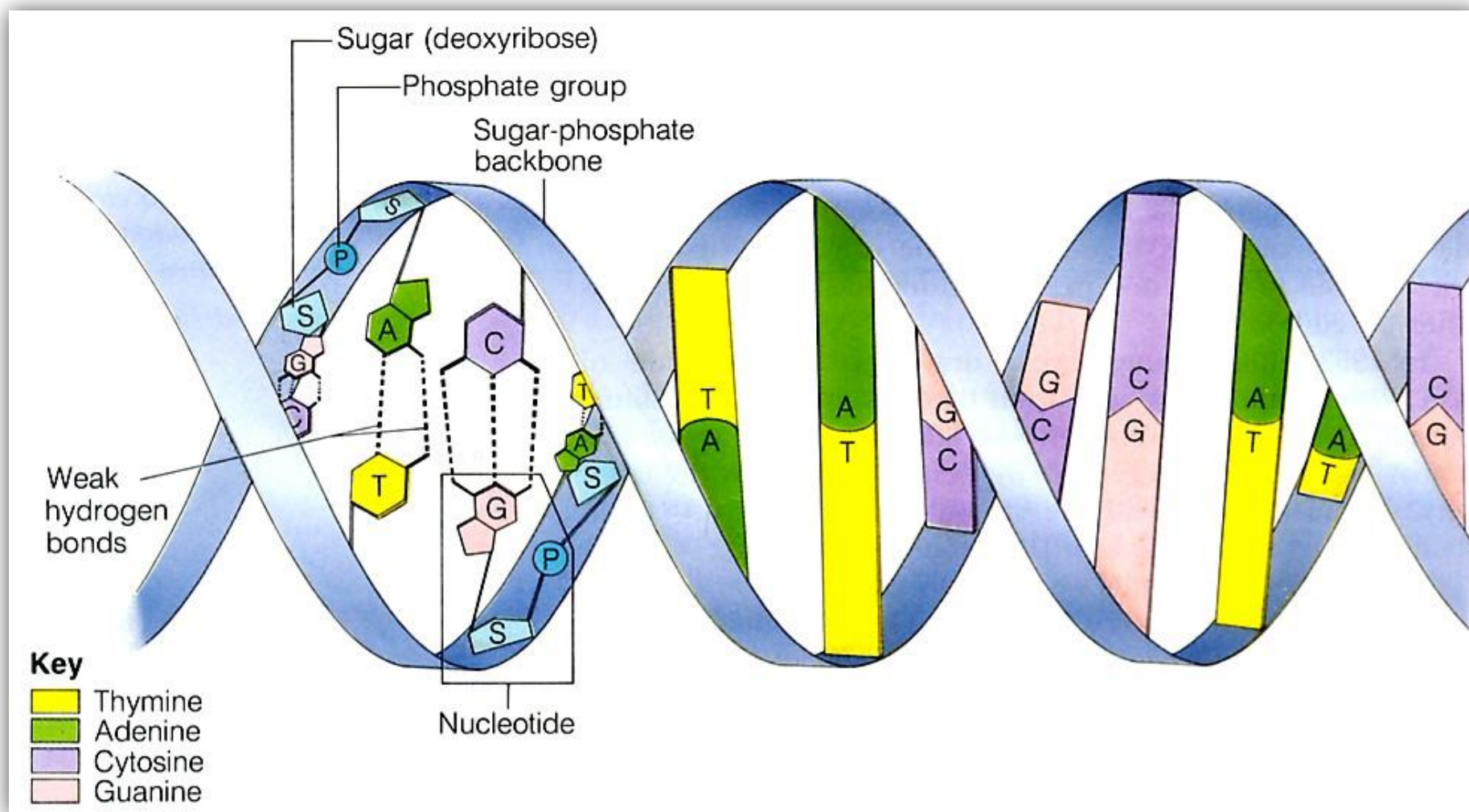
Kas, kāpēc, cik daudz?

2

- Genoms – gara, 4 dažādu nukleotīdu (A-T, C-G) virkne
- Satur **visu** informāciju par organisma uzbūvi
- Sadalīts vairākās hromosomās (cilvēkam – 2 x 23)
- Cilvēka DNS garums ~3 000 000 000 bp
- Nekompresētā veidā ~3GB, kompresētā ~800MB

Kas, kāpēc, cik daudz?

3



Genoma datu iegūšana

4

- Pirmo reizi pilns genoms iegūts: Human Genome Project (1990. – 2003.), izmaksas – \$ 2 700 000 000
- Šobrīd – iegūšanas laiks dažas dienas, izmaksas zem \$ 1 000 (pie lieliem apjomiem)
- Pēc dažiem gadiem, iespējams, \$ 100

- Plašākai auditorijai: <https://www.23andme.com/>, tikai vispārīgi dati, izmaksas \$ 99

Tehnoloģijas

5

Method	Sequencing by synthesis (Illumina)	Sequencing by ligation (SOLiD sequencing)
Read length	50 to 300 bp	50+35 or 50+50 bp
Accuracy	98%	99.9%
Reads per run	up to 3 billion	1.2 to 1.4 billion
Time per run	1 to 10 days, depending upon sequencer and specified read length ^[48]	1 to 2 weeks

Motivācija

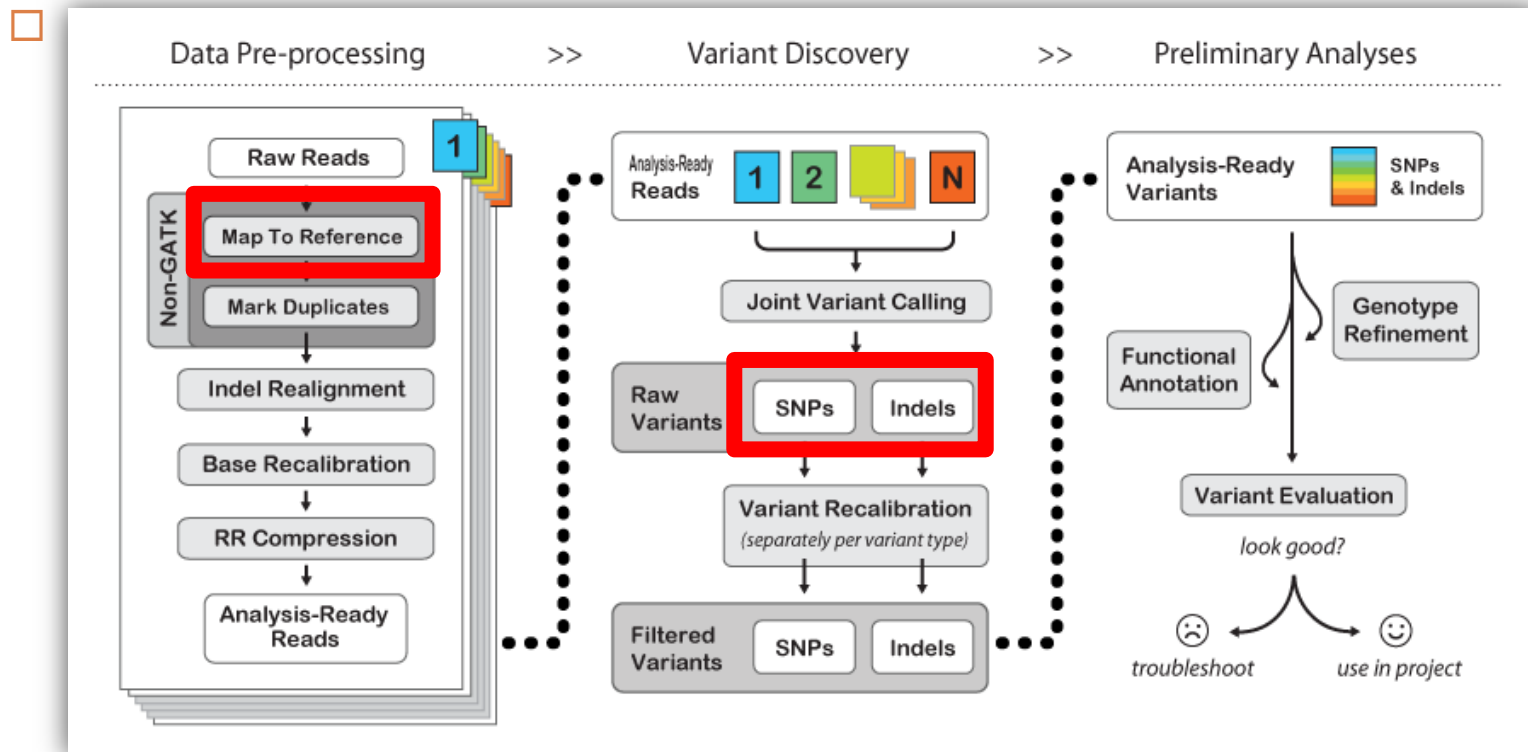
6

- Absolūti primārais - cilvēka DNS
- Atšķirība cilvēku DNS ir ar kārtu 0.1%
- Pielietojums genoma analīzei – iedzimtu slimību izpausmes prognozēšana (gan individuāli, gan cilvēku kopai) kā arī gēnu terapijā.
- Ņemot vērā sekvencēšanas izmaksu samazināšanos, jāspēj pārskatāmā laikā apstrādāt datus un sniegt pacientam ieteikumus, etc.
- Šobrīd visaktīvāk tiek pētīti vēža pacientu DNS, kodējošos reģionos (1.5% no kopējā cilvēka DNS)

Tipiskā darba plūsma DNS datu apstrādē

7

□ **Genoma de-novo salikšana** +



□ + **haplotipu atšifrēšana**

De-novo salikšana

8

- Daudz miljonu nolasījumu (atkarībā no tehnoloģijas, 30-50 bp), kuri jāsavieno vienā garā virknē
- Salikšana bez 'references' genoma
- Metodes – garākiem lasījumiem OLC (overlap/layout/consensus) algoritmi (parasti 'greedy'), īsākiem uz de Bruijn grafu balstītas metodes
- Nav vairs tik aktuāli, jo cilvēka genoms sekvencēts un salikts ļoti daudz reižu.

Meklēšana pēc 'references'

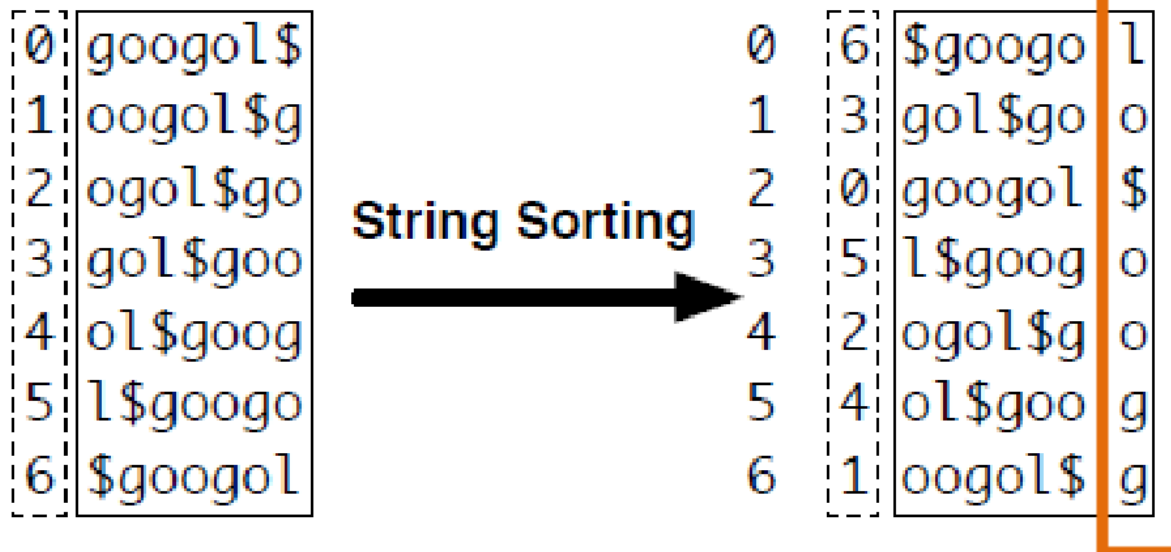
10

- Nepieciešams konkrēta gadījuma analīzei (rase, tautība, ķermeņa uzbūve, slimības uc.)
- Nepieciešamība pēc algoritmiem/programmatūras parastai darba stacijai. Šobrīd tāda ir – Bowtie2, BWA, GMAP
- Visi izmanto Burrows-Wheeler transformāciju un tās pielāgotu indeksāciju, kas nodrošina mazu atmiņas nepieciešamību meklējot, un iespēju saspiest uzglabājot

Meklēšana pēc 'references' - BWT

11

1. Original text = "googol"
2. Append '\$' to mark the end = X = "googol\$"
3. Sort all rotations of the text in lexicographic order
4. Take the last column.



Meklēšana pēc 'references' - BWT

12

- No transformācijas viennozīmīgi var atgūt oriģināltekstu
- Meklēšana notiek kāpjoties atpakaļ
- Precīzai meklēšanai sarežģ. $O(|W|)$
- Neprecīzai meklēšanai izmanto 'backtracking', sarežģītība parametriska $O(|W|) * f(|hg|, d)$

- All occurrences of substrings with a common suffix, W appear next to each other

$W = \text{"go"}$

0	6	\$googol
1	3	gol\$go o
2	0	gogol \$
3	5	l\$goog o
4	2	ogol\$g o
5	4	ol\$goo g
6	1	oogol\$ g

$X = \underline{g}o\underline{o}l\$$
0 3

Meklēšana pēc 'references'

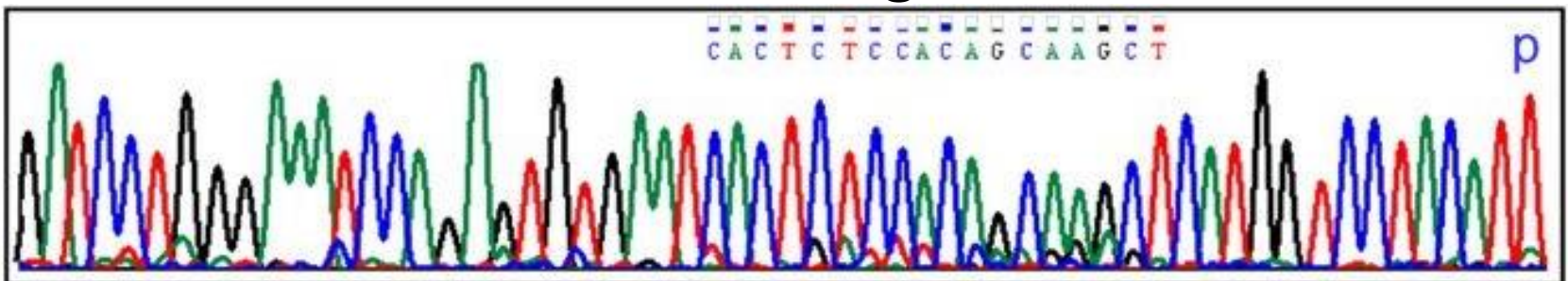
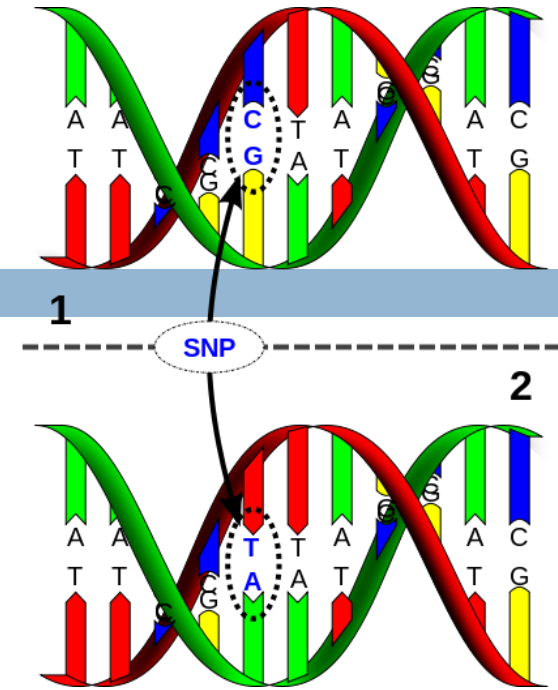
13

- Nozīmīgākās problēmas, ar ko jātiek galā: SNP/indels, atkārtotošies reģioni.
 - ▣ SNP/indels – dažādi sodi par dažādām neatbilstībām + max. robeža vai x pieļaujamās neatbilstības, resp. rēķina labošanas distanci.
 - ▣ Atkārtotošies reģ. - vairāku atbilstību gadījumā nolasījumu piekārto varbūtiski vienai no atbilstošajām pozīcijām
- Visi ātrākie un populārākie ir 'short read aligners' un radīti Illumina datiem, tātad problēmas palielinoties raksturīgajiem eksperimentālo datu garumiem (>500bp).
- Atrākai un precīzākai meklēšanai var izmantot SNP datu bāzes. hg ~3Gbp, dbSNP ~30Mbp. Vai var iekļaut BWT indeksa veidošanas laikā?

SNP

14

- Single nucleotide polymorphism (parasti 2 varianti, nav vienmērīgi)
- Nozīmīgi meklējot noteiktas cilvēku grupas un slimības korelāciju
- Meklēšanā plaši lieto bezmaksas programmatūru: GATK, SAMtools.
- Abas darbojas līdzīgi - izmantojot nolasījumu varbūtības un 'references' genomu



Haplotipu atšifrēšana

15

- Diploīdos organismos (katrā šūnā 2x no katras hromosomas) – katrs SNP dod 3 dažādus variantus. Piem SNP = A/T → 3 haplotipi (AA, AT, TT)
- Vienīgā viennozīmīgā metode - sekvencēšana
- Izvēlētai cilvēku grupai var noteikt varbūtību.
- Izmantotās metodes šobrīd – novērojumi par noteiktu haplotipu biežumu noteiktās kopās
- Ja pieder dažādām kopām – pēc iespējas mazāks skaits haplotipu, kas izmantots kombinācijā

Apzinātie resursi

16

- <http://www.ncbi.nlm.nih.gov/>
- <http://www.ebi.ac.uk/>
- EBI (A.Brāzma, N. Kurbatova)
- LV Biomedicīnas centrs (D. Frīdmanis)
- Bioloģijas Mg/PhD studenti (A. Kalviša)

Tuvākie mērķi

17

- Esošās programmatūras/metozu salīdzinājums uz reāliem (ne simulētiem) datiem
- Algoritms precīzākai/lielāka garuma ($200 \text{ bp} < l < 1 \text{ kbp}$) eksperimentālo datu nolasījumu meklēšanai pēc 'references' (izmantojot SNP datu bāzes)
- Metodes implementācija (iespējams, balstīta uz esošo Bowtie2, BWA indeksēšanu, jo šobrīd BW transformācija piedāvā labāko atmiņas optimizāciju)
- NB! Ļoti nepieciešama cieša sadarbība ar attiecīgās nozares (DNS sekvenčēšana) bioloģiem (reāli dati, konsultācijas, nozares standarti)

Paldies par uzmanību!