

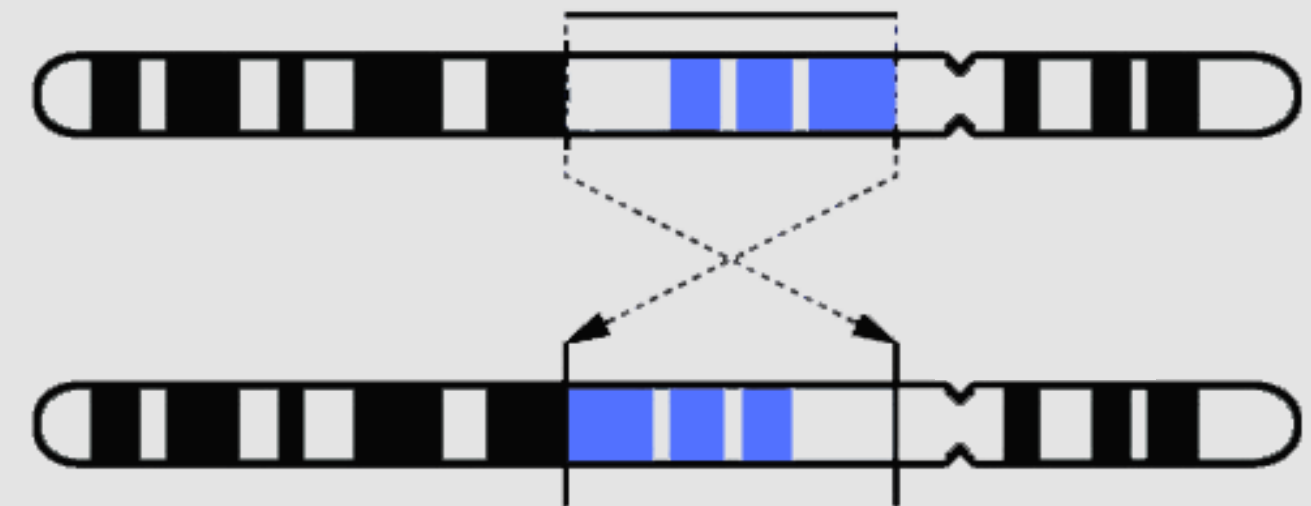
Detection of Genomic Inversions Using Sample-Specific Strings

Silvia Cambiago
Matr. 879382

Relatore: Prof.ssa Paola Bonizzoni
Correlatore: Dott. Luca Denti

Inversioni

- Variazioni strutturali
- Inversione e complemento basi
 - $A \Leftrightarrow T, C \Leftrightarrow G$
- Generate da Non-Allelic Homologous Recombination



R: AAC TTT TGG CTA AGA CCG AAA ACCA
T: AAC TTT CGG GTCT TAG CCA AAC CA

Motivazione



PROBLEMA

Identificare le inversioni



OBIETTIVO

Proporre un algoritmo di pattern matching che le rilevi con precisione



IMPORTANZA

Sono associate a diverse patologie e coinvolte nell'evoluzione

Idea di base

Article | Published: 22 December 2022

SVDSS: structural variation discovery in hard-to-call genomic regions using sample-specific strings from accurate long reads

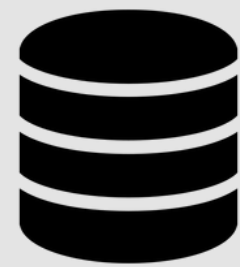
[Luca Denti](#), [Parsoa Khorsand](#), [Paola Bonizzoni](#) , [Fereydoun Hormozdiari](#)  & [Rayan Chikhi](#) 

[Nature Methods](#) **20**, 550–558 (2023) | [Cite this article](#)

4736 Accesses | **11** Citations | **49** Altmetric | [Metrics](#)

Problema

IDENTIFICAZIONE DELLE INVERSIONI

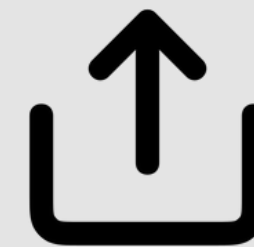


INPUT

Genoma R

Target T

Collezione di stringhe
specifiche (SFS) per T



OUTPUT

Collezione di segmenti
 $T[i, j]$ invertiti rispetto al
genoma di riferimento R

Approccio

R: CCATGATTGCGCATGCACTTG

T: CCATTGCATGCGAATCCTTG

SFSs: [CATT, CCT]

R: CCATGATTGCGCATGCACTTG

T: CCATTGCATGCGAATCCTTG



GCATGCGAAT



ATTCGCGATGC

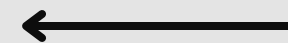


ATTCGCGATGC



R: CCATGATTGCGCATGCACTTG

Utilizzando KMP (Knuth-Morris-Pratt)



Complessità

n: numero SFSs

r: lunghezza genoma

\bar{m} : lunghezza media segmenti invertiti



TEMPO

$$\Theta(r + n \cdot \bar{m})$$



SPAZIO

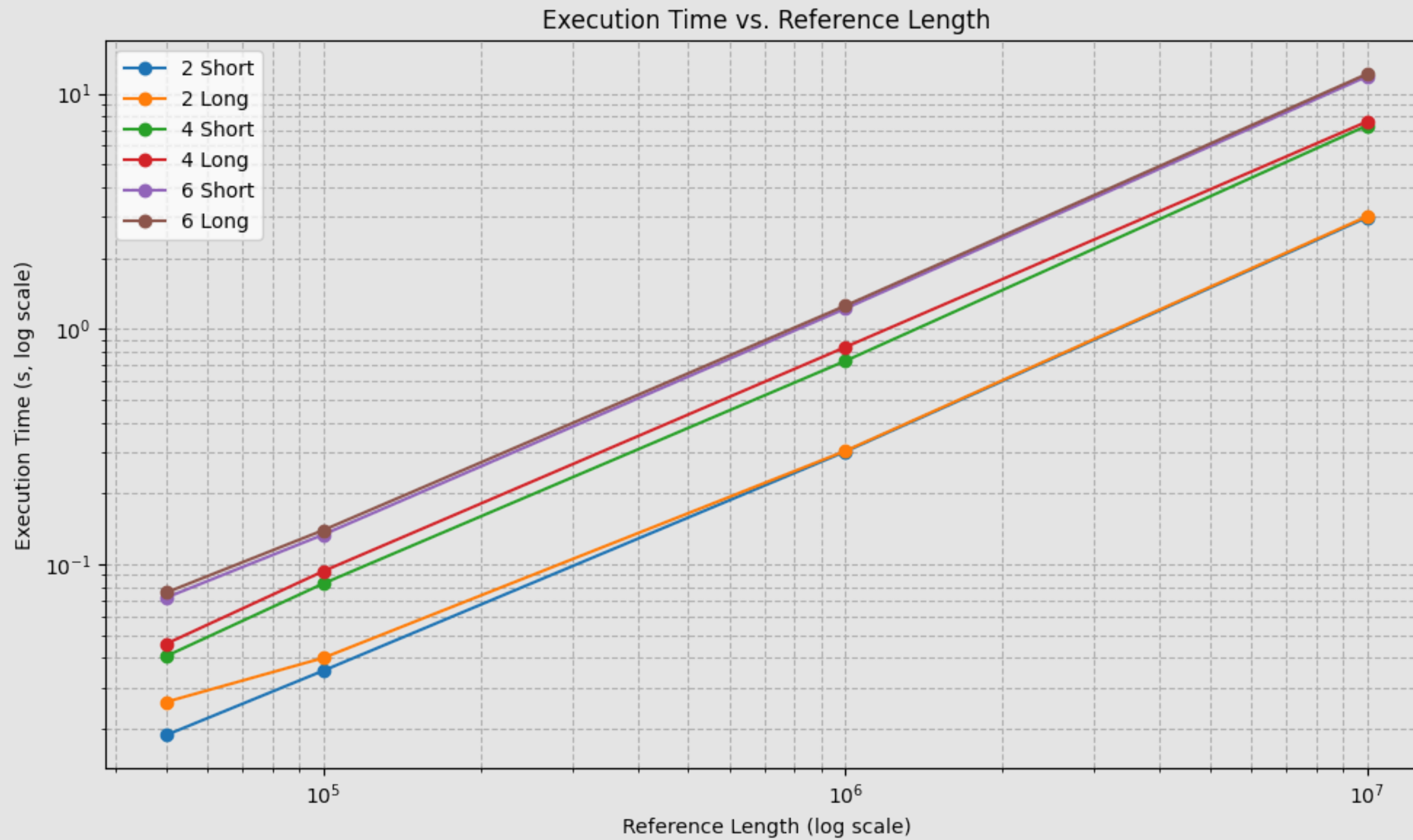
$$\Theta(n \cdot \bar{m})$$

Risultati

- 2 lunghezze
 - Corte: ~ 250 bp
 - Lunghe: ~ 6500 bp

REFERENCE	TARGET	2 INVERSIONI		4 INVERSIONI		6 INVERSIONI	
		CORTE	LUNGHE	CORTE	LUNGHE	CORTE	LUNGHE
50000	50000	0.0188	0.0260	0.0407	0.0459	0.0721	0.0758
100000	50000	0.0353	0.0400	0.0828	0.0934	0.1336	0.1396
1000000	50000	0.03009	0.03029	0.7305	0.8364	1.2190	1.2535
10000000	50000	2.9765	3.0082	7.3140	7.6067	11.8296	12.1238

Risultati



Duplicazioni Invertite

- Riconosce l'inversione ma non la presenza di una duplicazione invertita

ACTTCAGGGGAAAACTA

ACTTCAGGGGAAAACTCAGGGGAAAACTA

ACTTCAGGGGAAAACTTTCCCTGACTA

R: ACTTCAGGGGAAAACTA

T: ACTTCAGGGGAAAACTTTCCCTGACTA

Conclusioni



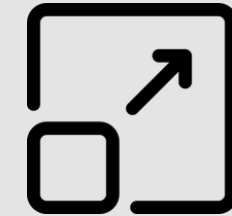
FATTORE DOMINANTE

La lunghezza del
genoma di riferimento r



BUONA EFFICIENZA

Per genomi di media
lunghezza



SCALABILITÀ LIMITATA

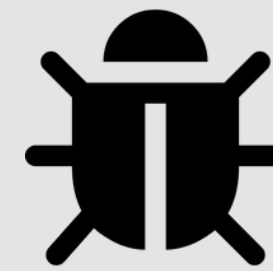
Per l'influenza di r

Prospettive Future



OTTIMIZZAZIONE

Migliorare la gestione di
sequenze di riferimento
lunghe



GESTIONE ERRORI

Per poter gestire dati
reali con mutazioni