

Dottorato di Ricerca in Informatica - Ciclo XXV Dipartimento di Informatica, Sistemistica e Comunicazione Facoltà di Scienze Matematiche, Fisiche e Naturali Università degli Studi di Milano - Bicocca



#### Algorithms for detecting variations from Next-Generation Sequencing Data

Presentazione Dottorato

30 Settembre 2010

Candidato: Stefano Beretta

Tutor: Prof.ssa Lucia Pomello









Stefano Beretta Algorithms for detecting variations from NGS data

∃ → < ∃</p>

Sequencing NGS Data Goals

### Outline



State of the Art & Ongoing Works

#### 3 Conclusions

Stefano Beretta Algorithms for detecting variations from NGS data

・ 同 ト ・ ヨ ト ・ ヨ

э

Sequencing NGS Data Goals

## Motivations

- Revolution in genome sequencing and analysis: from traditional methods to NGS (Next-Generation Sequencing)<sup>1</sup>
- Need to develop novel computational frameworks to analyze NGS data
- Goal: design algorithms to analyze NGS data for detecting sequence variations

<sup>&</sup>lt;sup>1</sup>Venter J.Craig, Nature (2010), *Multiple personal genomes await* 

Sequencing NGS Data Goals

# Genome Sequencing

Determination of the primary structure of a molecule  ${\sf DNA/RNA} \rightarrow$  sequence of nucleotides

- Traditional Methods (Sanger, 1977)
- Next-Generation Sequencing Methods (2005)



MotivationsSequencingState of the Art & Ongoing Works<br/>ConclusionsNGS Data<br/>Goals

## Sanger Vs. Next-Generation Sequencing

#### • Sanger (1977)



#### • NGS (2005)



 Motivations
 Sequencing

 State of the Art & Ongoing Works
 NGS Data

 Conclusions
 Goals

## Sanger Vs. Next-Generation Sequencing

#### Sanger (1977)

- Long Reads ( $\sim$ 1000 bp)
- 2 Low Throughput  $(\sim 10^6 \text{ bp/day})$
- 3 Low Coverage ( $\sim$  1x)
- Expensive (10<sup>3</sup> bp/\$)

#### NGS (2005)

- Short Reads (25-300 bp)
- I High Throughput  $(\sim 10^9 \text{ bp/day})$
- High Coverage (>10x)
- Low Costs (> 10<sup>5</sup> bp/\$)

くほし くほし くほし

 Motivations
 Sequencing

 State of the Art & Ongoing Works
 NGS Data

 Conclusions
 Goals

## Sanger Vs. Next-Generation Sequencing

#### Sanger (1977)

- Long Reads (~1000 bp)
- 2 Low Throughput  $(\sim 10^6 \text{ bp/day})$
- 3 Low Coverage ( $\sim$  1x)
- Expensive (10<sup>3</sup> bp/\$)

#### NGS (2005)

- Short Reads (25-300 bp)
- I High Throughput  $(\sim 10^9 \text{ bp/day})$
- High Coverage (>10x)
- Low Costs (> 10<sup>5</sup> bp/\$)

< ロ > < 同 > < 三 > <

Traditional Algorithms Models and Tools

Sequencing NGS Data Goals

## NGS Data



э

Sequencing NGS Data Goals

## **Biological Problems**

- Compare a known genome reference with an unknown genome (but sequenced by NGS)
- Infer transcripts data using short reads sampled by NGS

## **Detect Variations**

Sequencing NGS Data Goals

# Algorithmic Challenges

- More than  $10^9$  short sequences  $\Rightarrow$  Linear time algorithms
- Need for data compression / succint data structures
- New computational model and data structure for pattern matching (es. hashing, Burrows-Wheeler transf., suffix array)<sup>2 3 4</sup>

 $<sup>^2</sup>$ Langmead B, et al., Genome Biology (2009), Ultrafast and memory-efficient alignment of short DNA sequences to the human genome

 $<sup>^{3}</sup>$ Dalca V.A. and Brudno M., Briefings in Bioinformatics (2010), Genome variation discovery with high-throughput sequencing data

<sup>&</sup>lt;sup>4</sup>Li H. and Homer N., Briefings in Bioinformatics (2010), A survey of sequence alignment algorithms for next-generation sequencing  $\langle \Box \rangle \land \langle \Box \rangle \land \langle \Xi \rangle \land \langle \Xi \rangle$ 

Computational Problems dentification of Structural Variations Characterization of Alternative Splicing Events

### Outline





#### 3 Conclusions

Stefano Beretta Algorithms for detecting variations from NGS data

・ 同 ト ・ ヨ ト ・ ヨ

Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

## **Computational Problems**

- Identification of differences (Structural Variations) between a known genome (*reference*) and an unknown genome sequenced by NGS (*donor*).
  - Biological Motivations:
    - SV are common in human individuals and are related to diversity and disease susceptibility  $^{5\ 6}$
    - Detecting SV is crucial in medical and biological studies of several diseases

<sup>5</sup>Korbel, et al., Science (2007), Paired-End Mapping Reveals Extensive Structural Variation in the Human Genome <sup>6</sup>Tuzun, et al., Nature Genetics (2005), Fine-scale structural variation of the human genome < ≧ > < ≧ > < ≥ <

Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

### **Computational Problems**

- Characterization of variations (i.e alternative splicing events) among different transcripts sequences (sequenced by NGS) of the same gene.
  - Biological Motivations:
    - Human genes undergo AS (alternative splicing)
    - AS is the key process in determining transcriptomes diversity

- - E - - E

Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

# Identification of Structural Variations (SVs)

#### • Structural Variations (SVs)

- Insertions
- Deletions
- Inversions (>5 Kb)



Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

# Identification of Structural Variations (SVs)

- $\bullet\,$  Structural variation discovery using maximum parsimony is NP-hard  $^7$
- $\bullet\,$  Actual tools consider only one alignment (for each short read) and discard all the others  $^8\,$
- Probabilistic frameworks have been designed for the identification of specific SVs

<sup>&</sup>lt;sup>4</sup> Hormozdiari F., Alkan C., Eichler E., Sahinalp C., Genome research (2009), *Combinatorial algorithms for structural variation detection in high-throughput sequenced genomes* 

<sup>&</sup>lt;sup>8</sup>Medvedev P., Stanciu M., Brudno M., Nature methods (2009) *Computational methods for discovering structural* variation with next-generation sequencing

Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

# Identification of Structural Variations (SVs)

- Problem: predicting structural variations
  - **Input**: a set *S* of paired-ends (PEs) from a donor genome *D* and a reference genome *R*.
  - **Goal**: compute the set of structural variations that explains how *D* differs from *R*
- Previous approaches:
  - consider each SV separately
  - adopt probability based formulation
- Our Approach:
  - design a specific tool for PEs
  - develop an integrated approach for all SVs

4 3 6 4 3

Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

# Identification of Structural Variations (SVs)

#### • Algorithmic Solution

- PEs are aligned to the reference genome *R* (> 1 locations and > 1 orientations) and clustered into
  - Concordant mapped  $PEs \Rightarrow Donor = Reference$
  - $\bullet \ \ \mathsf{Discordant} \ \ \mathsf{mapped} \ \mathsf{PEs} \Rightarrow \mathsf{Structural} \ \mathsf{Variations}$
- Discordant PEs are analyzed to detect
  - Insertion / Deletion, Inversion, Other complex cases

#### Issues

- Large SVs are hard to detect
- Detecting combinations of different SVs

伺 ト イ ヨ ト イ ヨ ト

Computational Problems Identification of Structural Variations Characterization of Alternative Splicing Events

### Characterization of Alternative Splicing Events

#### Alternative Splicing



(日) (同) (三) (三)

э

## Characterization of Alternative Splicing Events

- No techniques based on short reads comparison
- No characterization of differences of transcripts
- Developed algorithms map the NGS data into the given reference genome to infer splice junctions<sup>9</sup> <sup>10</sup>



 $<sup>^9</sup>$ Bryant, et. al., Bioinformatics (2010) Supersplatspliced RNA-seq alignment

## Characterization of Alternative Splicing Events

- Problem: inference of alternative splicing (AS) events
  - Input: a set of short reads from transcripts of a gene
  - Goal: graph representation of AS events (genome scale)
- Previous approaches:
  - detect splice junctions
  - validate transcripts
- Our Approach:
  - detect differences (which are few) and discard similarities (too many)
  - no alignment to the reference genome

3 1 4 3

Motivations Computational Problems State of the Art & Ongoing Works Identification of Structural Variations Conclusions Characterization of Alternative Splicing Events

### Characterization of Alternative Splicing Events

#### • Algorithmic Solution

- We index short reads with a hash table in order to:
  - assembly chains of short reads to compose Exons
  - identify junction points of Exons



#### Issues

- Not unique identification of splicing junction
- Some AS events are hard to characterize
- Needs for a topological sort of Exons

#### Outline



#### 2 State of the Art & Ongoing Works



Stefano Beretta Algorithms for detecting variations from NGS data

Image: Image:

## Conclusions



Stefano Beretta Algorithms for detecting variations from NGS data