

Comparison of short read aligners with Galaxy

Subazini Thankaswamy Kosalai, Jens Nielsen, Intawat Nookaew*

Systems and Synthetic Biology,

Department of Chemical and Biological Engineering,

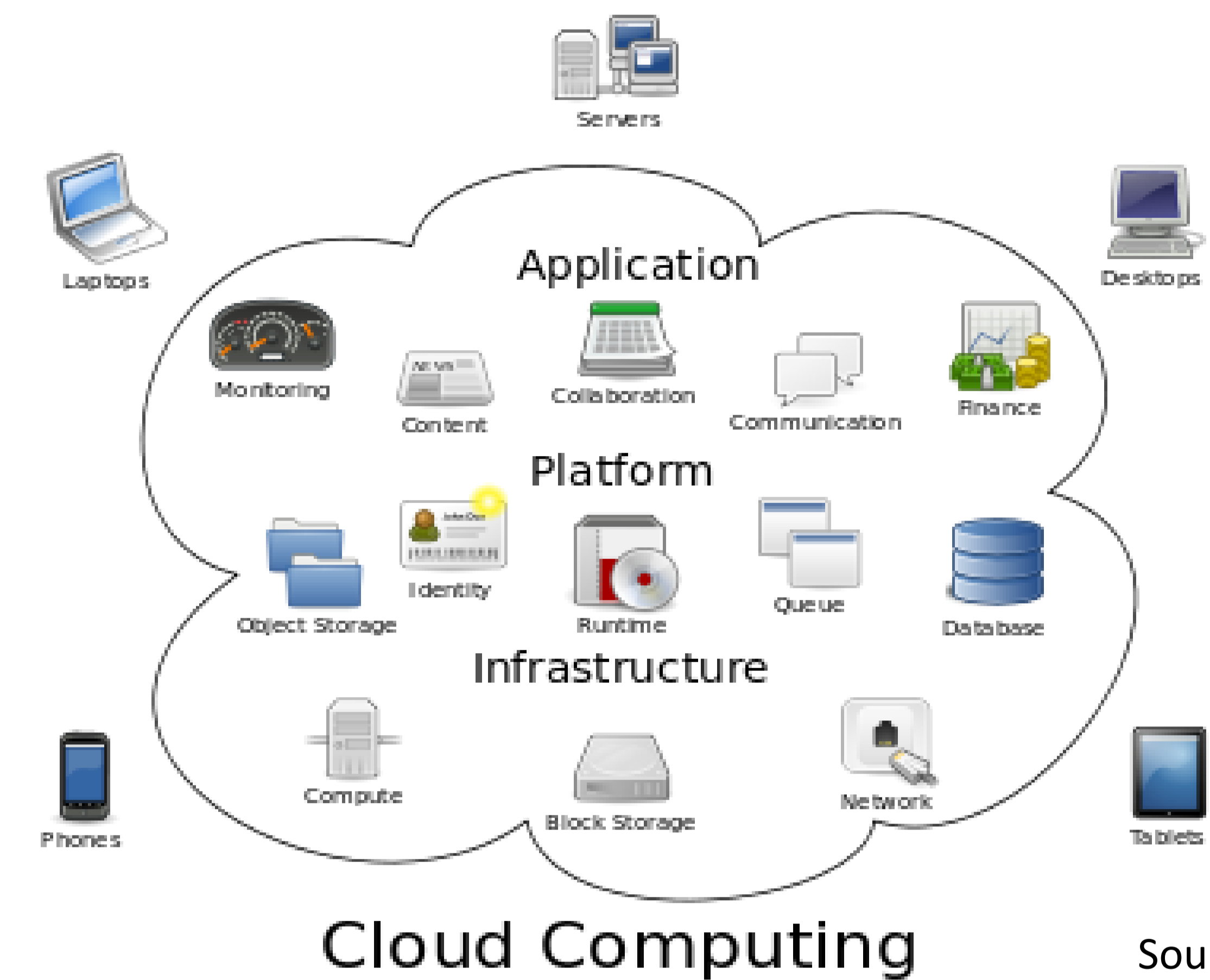
Chalmers University of Technology, Gothenburg, Sweden 41296.



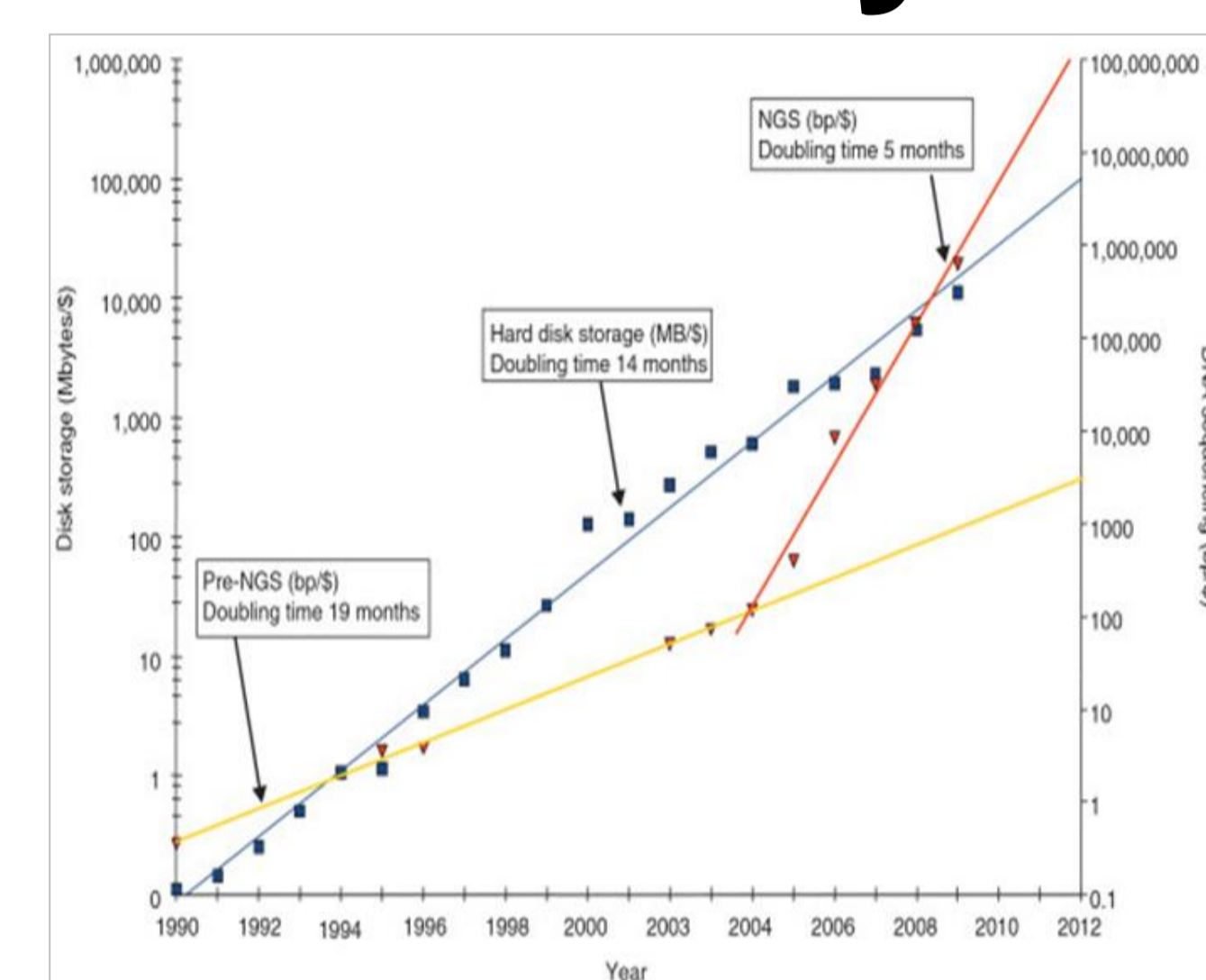
Abstract

The emergence of Next generation sequencing (NGS) technology ensued production of large-scale data in fast pace demanding increased storage resource and computational power. The essential step in NGS analysis is read alignment or mapping with reference genome to determine the desired DNA sequence. The genetic difference between strains attained on mapping can also be used in variant detection and annotation. It is difficult to determine the position of short reads by mapping, mostly in the case of repetitive regions. Many tools developed for short read sequence alignment are available public and mostly command-line. On the other hand end-users find it more convenient when the tools are with user-interface. Galaxy is an integrated frame, which can be used in resolving computational issues, by allowing the tools to be deployed in cloud called Galaxy CloudMan. It also allows user to create a well-defined user-interface for command-line tools in XML. In this work, we have deployed different mappers or aligners based on different algorithms in Galaxy CloudMan and compared them for sensitivity and speed with allowed mismatch. XML Wrapper files are generated to create user-defined interface for the command-line mappers and deployed in galaxy so that it can be utilized for constructing workflows. The challenge is to select a mapping tool with fundamental priorities of speed, sensitivity and minimal memory usage. We made criteria for setting different parameters suitable for researchers' project and evaluating the aligners using mapping speed, RAM occupancy, sensitivity and accuracy using short read simulators and some real data.

Service over network



NGS analysis



Galaxy

Galaxy provides a web based interface for deploying tools in cloud and also for analyzing and manipulating NGS data

Advantages

- Provides computational and storage resources
- Independent queries on genomic data from different sources (UCSC, Yeast mine ...)
- Can share history with other Galaxy users.
- Workflow, can be developed in combination, refinement, calculation, extraction and visualization of queries
- Multiple analysis by Query intersection, subtractions and proximity searches

Stein *Genome Biology* 2010 11:207

Challenges

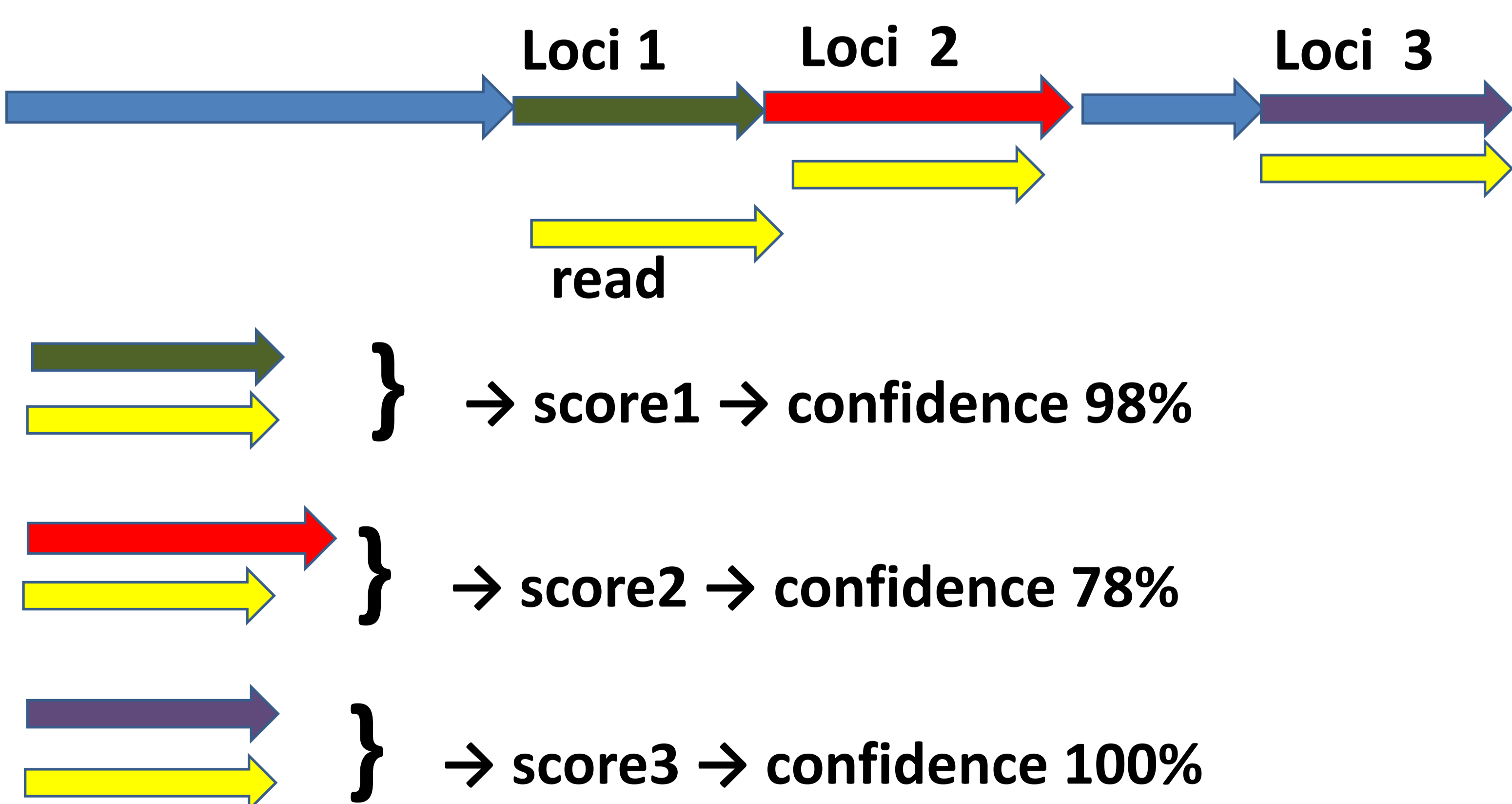
- Big data size of NGS sequences
- Extensive memory required for data storage
- Servicing of certain softwares used in NGS analyses imposes additional cost
- Moving data is non trivial
- Requires extensive computing power for data management and analysis

Galaxy CloudMan

Parameters selection

	A	B	C	D	E	F	G	H	I	J
1	Blat	Rowle2	BWA	GEM	Novocraft	RAZER	smallt	stampy	scap	
2	maxGap	-<fup>	-c INT							
3		Size of max gap between tiles in a clump. Usually between 0 and 3	Maximum number of gap opens [1]							maximum gap size allowed on a read, default=0bp
4	extendthroughN									
5		Allows extension of alignment through large blocks of N's								
6	Extension									
7										
8	Max insert									
9		The maximum fragment length for valid paired-end alignments								
10										
11	Min insert									
12		The minimum fragment length for valid paired-end alignments								
13										
14	Score									
15		match minus mismatch minus gap penalty. Default 90								
16										
17	Gap open penalty									
18		sets the read gap open (<int>) and extend (<int>) penalties.	Gap open penalty [11]							
19										
20	Gap extend penalty									
21		sets the reference gap open (<int>) and extend (<int>) penalties. A reference gap of length N gets a penalty of <int> * N * <int>. Default: 5, 3.	Gap extension penalty [4]							
22										
23	Mismatch penalty									
24		Allows one mismatch in tile penalties, both integers.								
25										
26										
27										
28										

Read mapping



The screenshot shows the Galaxy web interface. On the left, there's a 'Tools' panel with 'Blat' selected. The main area shows the command line for 'novoalign' and its output, which includes alignment scores and sequence alignments. On the right, there's a 'History' panel showing the workflow history.

References:

- Thakur RS, Bandopadhyay R, Chaudhary B, Chatterjee S. Now and next-generation sequencing techniques: future of sequence analysis using cloud computing. *Front Genet.* 2012;3:280
- Afgan E, Baker D, Coraor N, Chapman B, Nekrutenko A, Taylor J. Galaxy CloudMan: delivering cloud compute clusters. *BMC Bioinformatics.* 2010 Dec 21;11 Suppl 12:S4.

