



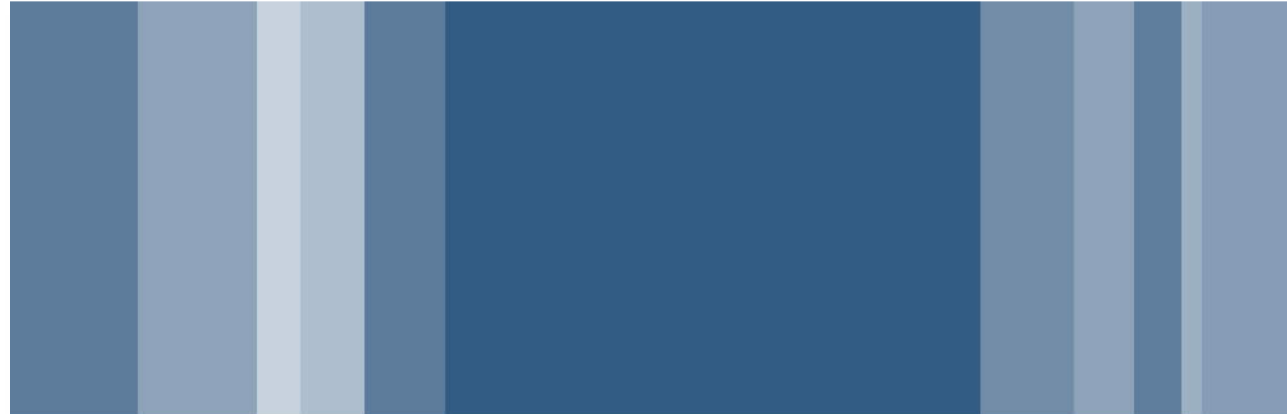
Genomic Computing

Politecnico di Milano



**POLITECNICO
DI MILANO**

Dipartimento
di Elettronica, Informazione e
Bioingegneria



Discovering similar (epi)genomics feature patterns in multiple genome browser tracks

Piero Montanari¹, Arnaud Ceol², Ilaria Bartolini¹, Paolo Ciaccia¹, Marco Patella¹, Stefano Ceri³, Marco Masseroli³

¹ DISI - Università di Bologna, ² IIT@SEMM – IIT,

³ DEIB - Politecnico di Milano



Background

- **Next Generation Sequencing (NGS)** is opening many interesting practical and theoretical computational problems
- **Genome browsers** (UCSC Genome Browser, Integrated Genome Browser (IGB)) allow:
 - Visual inspection and identification of interesting **patterns** on multiple genome browser tracks
- **Pattern**: a set of (epi)genomic regions / peaks at given distances from each other in different tracks



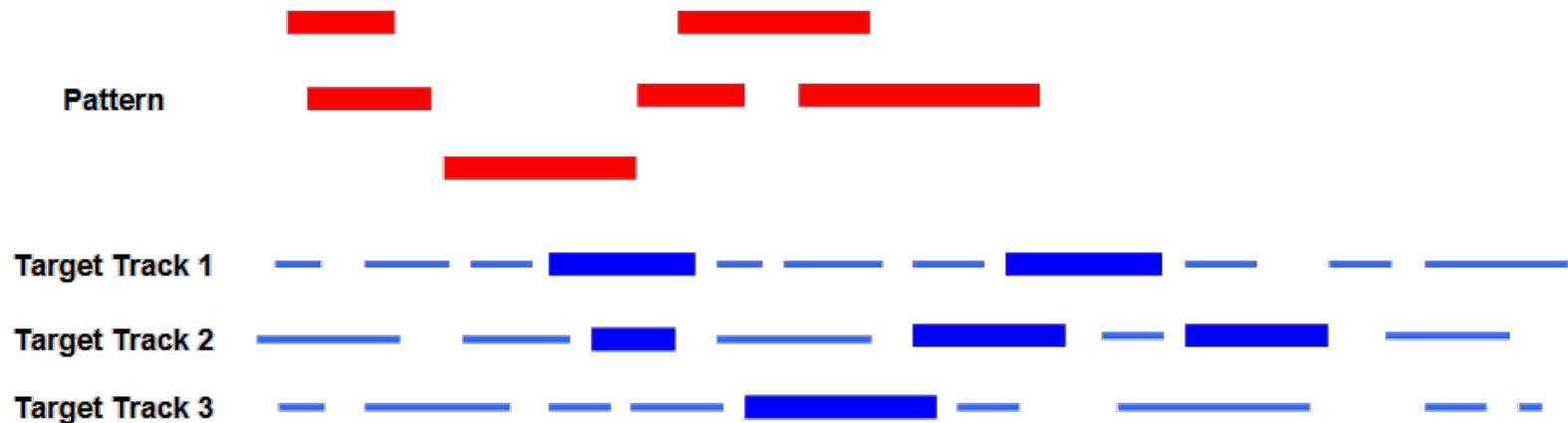
- e.g.: gene regulatory DNA areas, which include heterogeneous (epi)genomic features (different histone modifications, TFBSs, , ...)



Motivation

Once such patterns are visually identified in a genome section:

- **Search for pattern occurrences in whole genome:**
 - Complex computational task
 - Currently not supported
- **Their discovery in whole genome** very important for:
 - Biological interpretation of NGS experimental results
 - Comprehension of biomolecular phenomena





- We defined an **optimized pattern-search algorithm** to find **genomic region sets** that are **similar to a given pattern**
 - Efficiently
 - In large (epi)genomic / transcriptomic data sets
- We implemented the algorithm within an **IGB plugin**, named **SimSearch**, which allows intuitive user interaction in both:
 - **Visual selection** of an interesting **pattern** on loaded IGB tracks
 - **Visualization** of occurrences of **similar patterns** identified in the whole genome



Method

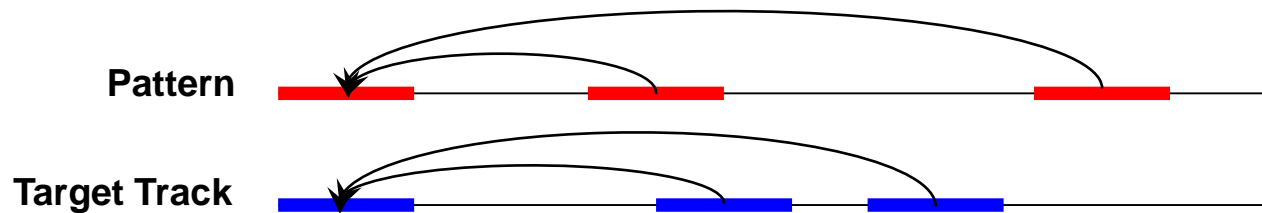
- **Best-Matching Problem (BMP)** is suspected to be NP-hard
- For the considered application:
 - **Aligned genomic data** -> strictly increasing region sequences
 - **$M \ll N$** (M and N: number of elements in a pattern and in the target tracks to be compared, respectively)
- **Our proposal:** a **Root-element approach** (R-BMP) and **Windowed Dynamic Programming algorithm** (WDP-BMP) to lower complexity:
 - Best matching for each element of the pattern in the target track through a **binary search**
 - **Complexity** only of $O(N \cdot \log(N))$
 - **Applicable** also to (**very**) **large data** problems



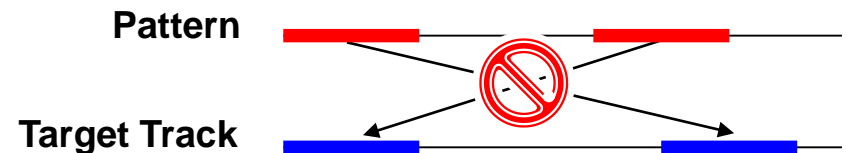
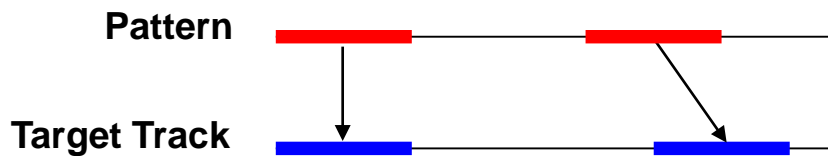
Discovering similar patterns in genomics tracks

Method – *Base problem / model*

- **Pattern** = single track
- **Regions** = points, identified by their linear genomic coordinate
 - Patterns and tracks = sequences of integers (region coordinates)
- Only relative distances between regions/points are relevant



- A **matching** is a strictly increasing function f that assigns to each pattern (query) element a (target) track element
 - **Preserving element order**





- **Pattern matching** typically solved by a cost based approach, where **lower cost** implies **higher similarity**
- **Root-element approach** (R-BMP): The cost C of a **matching** f is the sum of squared distances relative to the first (root) matching pair

$$C_f(Q, T) = \sum_i (T_{f(i)} - Q_i - (T_{f(1)} - Q_1))^2$$

- Q and T : sequences of query (pattern) and target elements
- $T_{f(1)} - Q_1$: root-distance (relative distance between root elements)
- **Goal**: find the matching with minimum cost



Example

- Pattern (query) $Q = \langle 1, 7, 10 \rangle$
- Target track $T = \langle 3, 5, 9, 11, 13, 14, 18, 21 \rangle$
- A possible matching: $\{(1 \rightarrow 3), (7 \rightarrow 9), (10 \rightarrow 21)\}$
 - Root-distance = $3 - 1 = 2$
 - Cost = $(3 - 1 - 2)^2 + (9 - 7 - 2)^2 + (21 - 10 - 2)^2 = 81$
- **Best matching** $\{(1 \rightarrow 5), (7 \rightarrow 11), (10 \rightarrow 14)\}$
 - Root-distance = $5 - 1 = 4$
 - Cost = $(5 - 1 - 4)^2 + (11 - 7 - 4)^2 + (14 - 10 - 4)^2 = 0$



Windowed Dynamic Programming algorithm (WDP-BMP)

- **Main result:** A matching is **optimal** if and only if all its **partial** matchings have **minimum cost**
- For all possible root positions in T:
 - They are $|T|-|Q|+1$
- For all regions Q_i in Q:
 - **Find** (with **binary search**) the best match for Q_i in T
 - Possibly, widen the “window” to avoid conflicts (regions in T that are best match for multiple regions in Q)
- For all possible matches:
 - **Compute** the cost using the partial cost obtained till then
 - **Abandon** the current solution if partial cost \geq best cost till then
- Overall **complexity:** $O(|T||Q|(\log|T|+|Q|))$
 - But **$O(|T|\log|T|)$** if $|Q| \ll |T|$ (usually $|Q|=1 \div 10$, $|T| \sim 10^5 \div 10^6$)



- **Multi-track patterns**
 - Same approach is repeated for each pattern track
- **Negative matching tracks**
 - Pattern tracks with regions that should not appear in the solution
 - Removed from the solution search space (before search start)
- **Partial matching tracks**
 - Pattern tracks with regions that might be missing in the results
 - A cost is considered for those regions that remain unmatched
- **Region features**
 - The cost function includes the (dis-)similarity between features (attributes) of query and target track regions
- **Regions as intervals**
 - Region size modeled as a region feature
- **Top-K distinct matchings**
 - To consider diverse results, we require that the best K results have no matched region in common



Finding enhancer regions

- Complex pattern
 - 2 single region **positive** tracks (H3K4me1 and H3K27ac)
 - 2 single region **negative** tracks (TSS and H3K4me3)
 - 4 single region **partial** tracks (CTCF, DHS, P300 and Pol2)
- Target data tracks about **K562** cell line (AML) from **ENCODE**
- **Top-100** results **visually inspected** by an expert
 - **100% precision** (of the resulting matchings to the pattern)
- **All** (1,651) results **compared** with “*enhancer*” chromatin state regions from **ChromHMM** tool (Broad Institute), data available in **ENCODE**
 - **85.46% precision** (but ChromHMM uses more histone marks)



Discovering similar patterns in genomics tracks

SimSearch: a plugin for IGB

Available at: <https://arnaudceol.github.io/simsearch/>

FG	BG	Load Mode	Track Name
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	Assembly
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	GENCODE
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	Common SNPs(142)
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	CpG Islands
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	Hq19 Diff
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	RefGene
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Manual	Genome

FG	BG
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Bed/BigBed, BigWig, GTF ...

Galaxy

HTS-flow (<http://arnaudceol.github.io/htsflow/>)

The Integrated Genome Browser (<http://bioviz.org/>)



Discovering similar patterns in genomics tracks



SimSearch: *define the pattern to search for*

Load pattern from file **Create** pattern by selecting IGB tracks / regions

Search Options

Load pre-defined pattern:

Dataset ID	Type	Distance/Range	Score	Target Dataset
<input type="checkbox"/> <input type="checkbox"/> H3K4me1	perfect	0		H3K4me1_DATA_ALL.gtf
<input type="checkbox"/> <input type="checkbox"/> H3K27ac	perfect	0		H3K27ac_DATA_ALL.gtf
<input type="checkbox"/> <input type="checkbox"/> H3K4me3	negative			H3K4me3_DATA_ALL.gtf

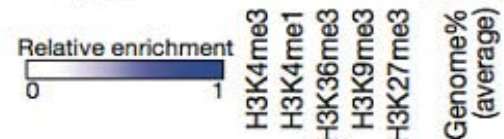
Where to search

enhancer_discovery_pattern

- 01_Active_TSS
- 02_Flanking_active_TSS
- 03_Transcription_at_gene_5prime_and_3prime
- 04_Strong_Transcription
- 05_Weak_Transcription
- 06_Genic_enhancers

Load pre-defined pattern

Chromatin state	Abbreviation	emissions	Cov.
1 Active TSS	TssA		0.7%
2 Flanking active TSS	TssAFlnk		0.5%
3 Transcr. at gene 5' and 3'	TxFlnk		0.1%
4 Strong transcription	Tx		3.6%
5 Weak transcription	TxWk		11.6%
6 Genic enhancers	EnhG		0.4%
7 Enhancers	Enh		2.8%
8 ZNF genes + repeats	ZNF/Rpts		0.2%
9 Heterochromatin	Het		2.6%
10 Bivalent/poised TSS	TssBiv		0.1%
11 Flanking bivalent TSS/Enh	BivFlnk		0.1%
12 Bivalent enhancer	EnhBiv		0.1%
13 Repressed Polycomb	ReprPC		1.2%
14 Weak repressed Polycomb	ReprPCWk		8.3%
15 Quiescent/low	Quies		67.8%



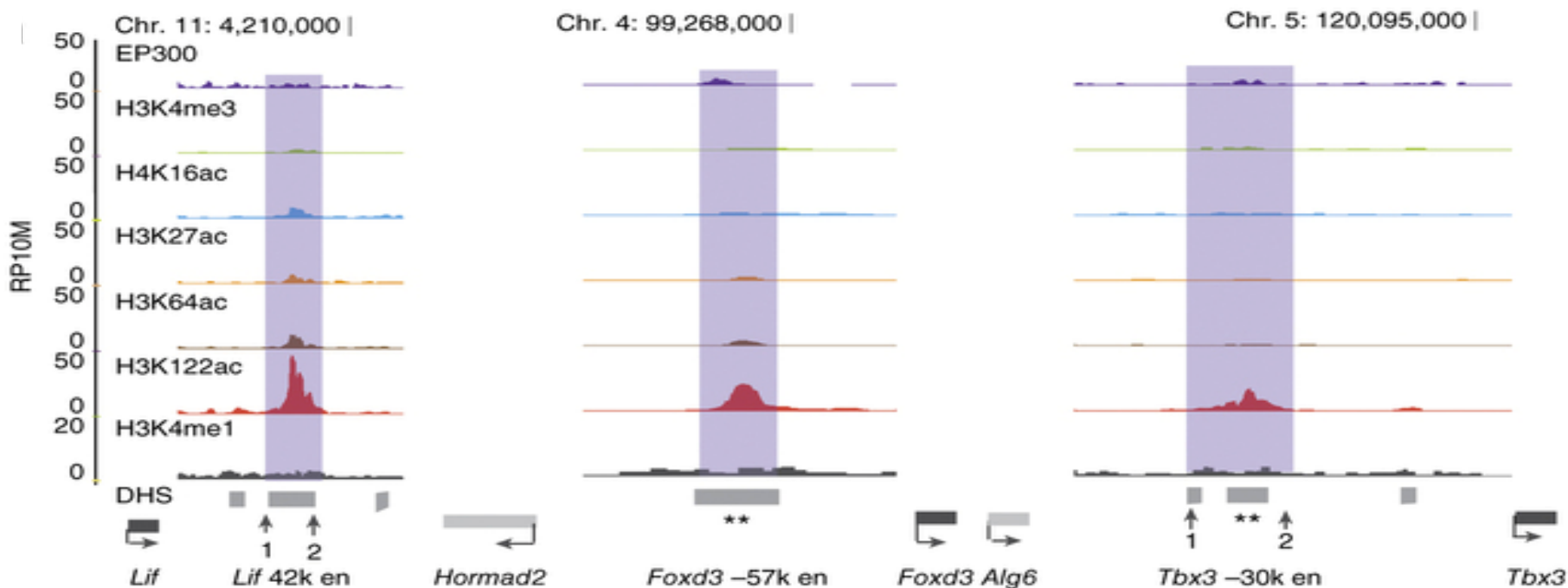
(e.g. CHromHMM chromatin states from Roadmap Epigenomics Consortium, *Nature* 2015)



Discovering similar patterns in genomics tracks

SimSearch: *searching for new enhancers*

Pradeepa et al. (*Nature Genetics* 2016) identified a **new class of enhancer** (with H3K122ac and absence of H3K27ac) in mouse



Several examples validated in mouse

- Homozygous deletion of group 2 putative enhancer 42 kb **downstream** of *Leukemia inhibitory factor* (*Lif*) gene led to **reduced expression** of *Lif*, but not of flanking *Hormad2* gene
- Deletion of one allele of putative enhancer 30 kb **upstream** *Tbx3* led to **downregulation** of *Tbx3*



Discovering similar patterns in genomics tracks

SimSearch: searching for new enhancers

Look for this new pattern in human samples:

1. Load tracks of human data

2. Add tracks to the pattern

3. Add TSS (negative)

The screenshot shows the SimSearch web interface. At the top, the genomic region is set to chr1:0-249,250,621. On the left, a track list includes H3K4me1_broad_peaks (1).bed (+/-), H3K122ac_broad_peaks.bed (+/-), H3K27ac_broad_peaks.bed (+/-), RefGene (+), Coordinates, and RefGene (-). The main visualization area shows a genomic track with a scale from 0 to 200,000,000. Below the tracks, a table lists the loaded tracks and their configuration in the search pattern.

Dataset ID	+/-	Type	Distance/Range	Score	Target Dataset
<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> H3K4me1_broad_peaks (1).bed	<input type="checkbox"/>	perfect	0		H3K4me1_broad_peak...
<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> H3K122ac_broad_peaks.bed	<input type="checkbox"/>	perfect	0		H3K122ac_broad_peak...
<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> H3K27ac_broad_peaks.bed	<input checked="" type="checkbox"/>	negative	0		
<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> TSS (IGB)	<input type="checkbox"/>	negative	20000		TSS (IGB)

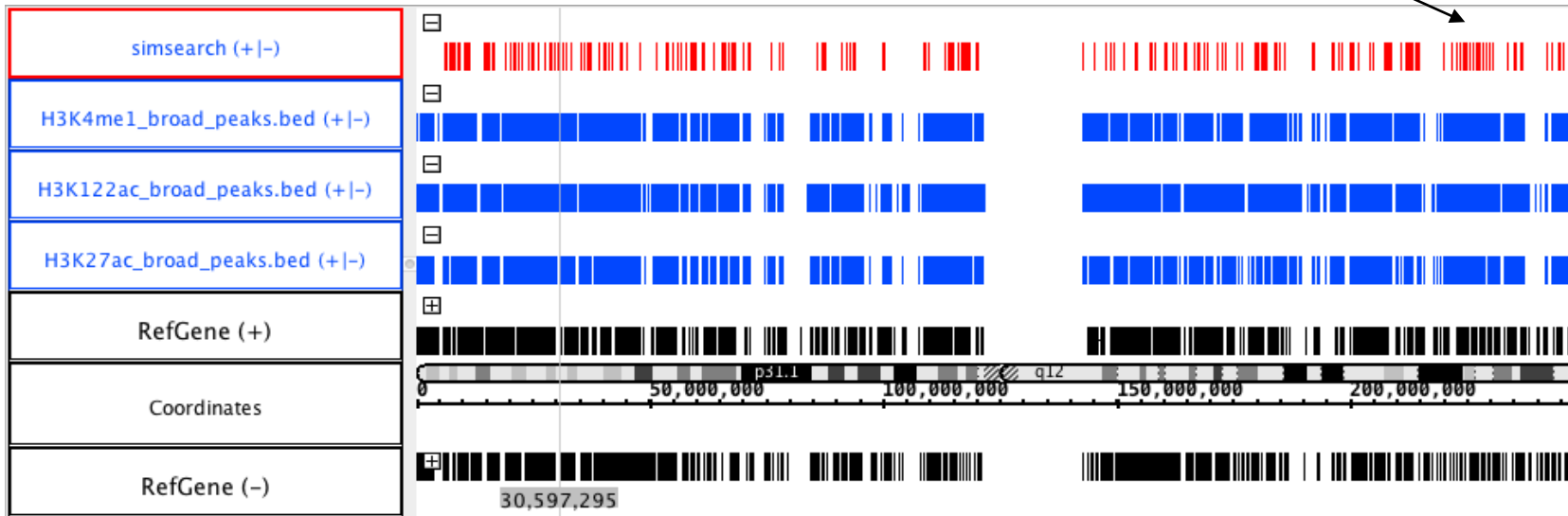


Discovering similar patterns in genomics tracks

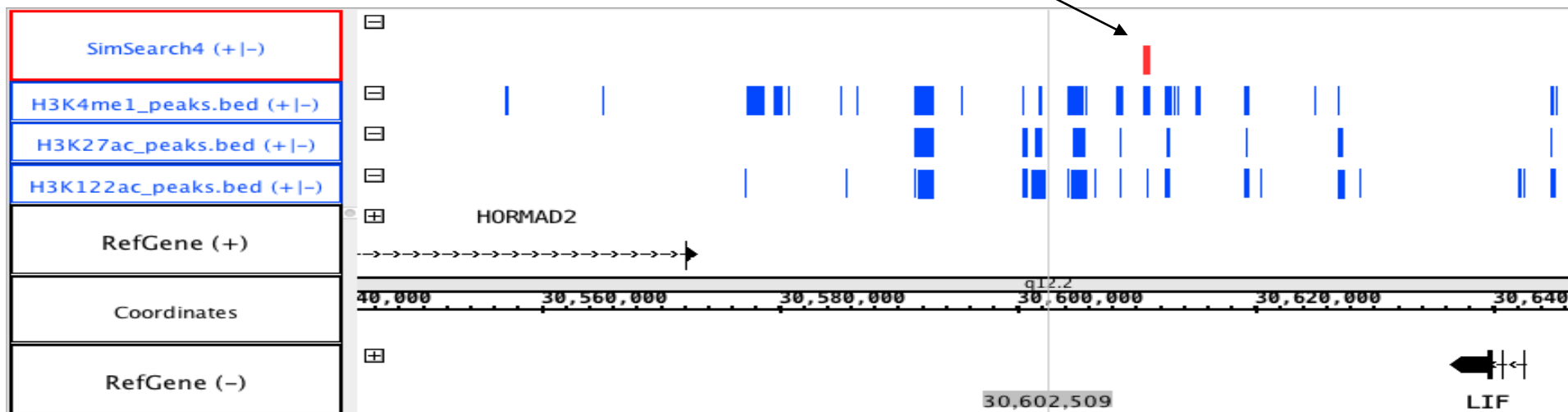


SimSearch: *editing the pattern*

In a **new track (red)**, see all active regions (matches) found



Similar match found **in human cells**





- Definition of a **method** for **finding patterns in genomic sequences** and of a **dynamic programming** algorithm for **efficient solution**
 - Experimental evaluation found it **accurate** and **efficient**
- Algorithm implemented as a **plugin** for **Integrated Genome Browser**
- Applicable to **any** (epi)genomic / transcriptomic **region data**, e.g.:
 - Histone Modifications (HM), Transcription Factor Binding Sites (TFBS), Single Nucleotide Polymorphisms (SNP), Differentially Expressed Genes (DEG), DNase I Hypersensitive Sites (DHS), Transcription Start Sites (TSS), or any other annotations
- Support for **understanding biomolecular mechanisms**
 - Response to treatments, pathology onset/development, ...
- Thanks to the **GenData 2020** PRIN project and **GeCo** ERC project!



Thank you for your attention!

Any question?

