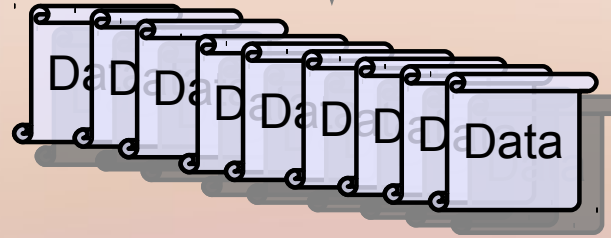


# Efficient comparison of sets of intervals with NC-lists

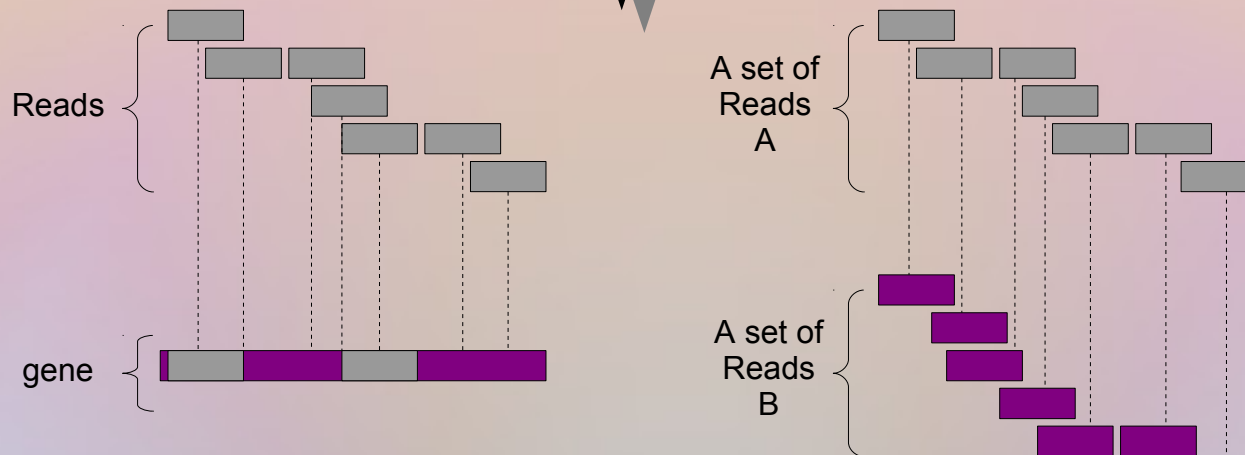
*Yufei Luo, Matthias Zytnicki, Hadi Quesneville*  
***URGI INRA Versailles***

- Introduction
- Methods
- Results

# High-throughput sequencing



Genomic coordinates

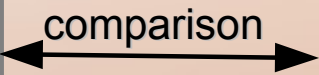


Are there **optimized algorithms** for the fast comparison of large amount of sequenced data?

RNA-seq data



reads



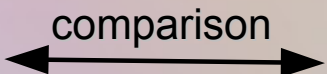
annotations



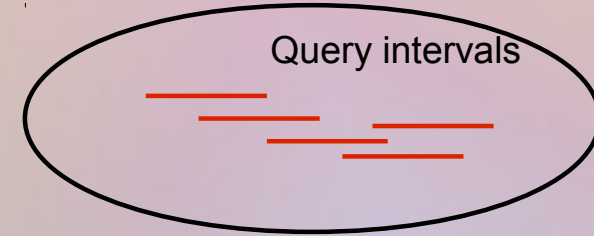
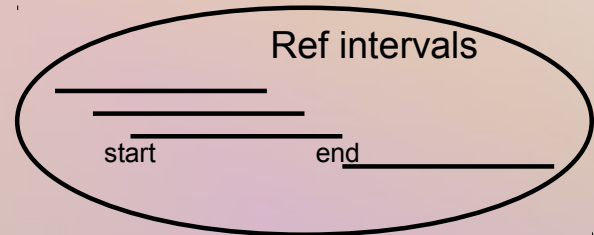
Gene expression



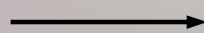
intervals



intervals



NC-List



$O(n+\log N)$ ,



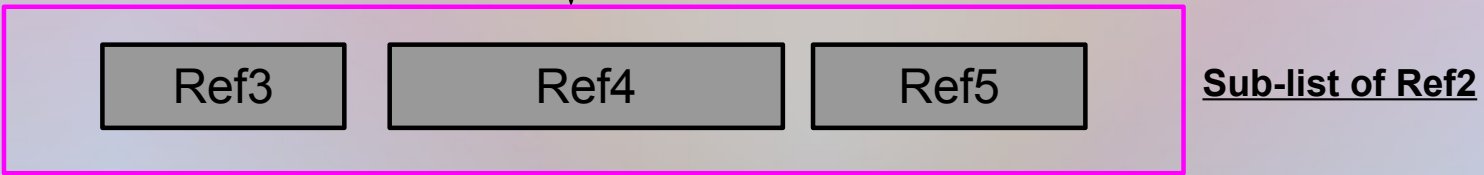
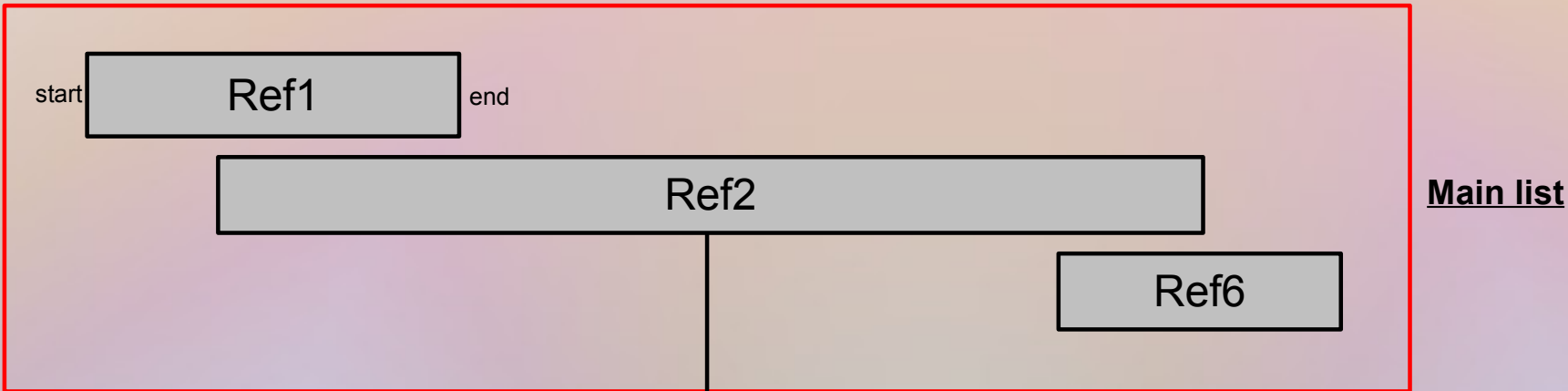
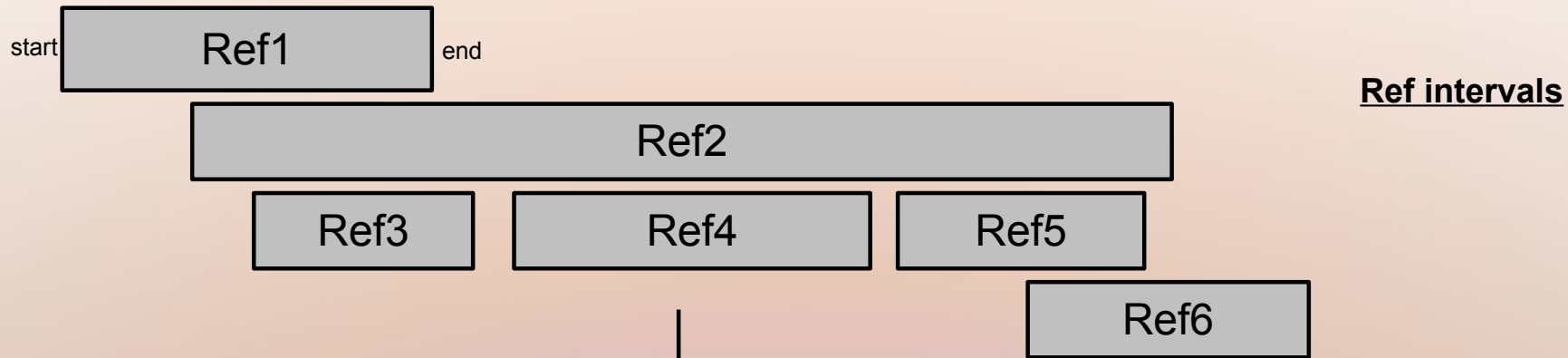
Is it true?

Alexander V. Alekseyenko and Christopher J. Lee. *Bioinformatics*, 2007

$N = \text{number of reference intervals ;}$   
 $n = \text{number of overlaps}$



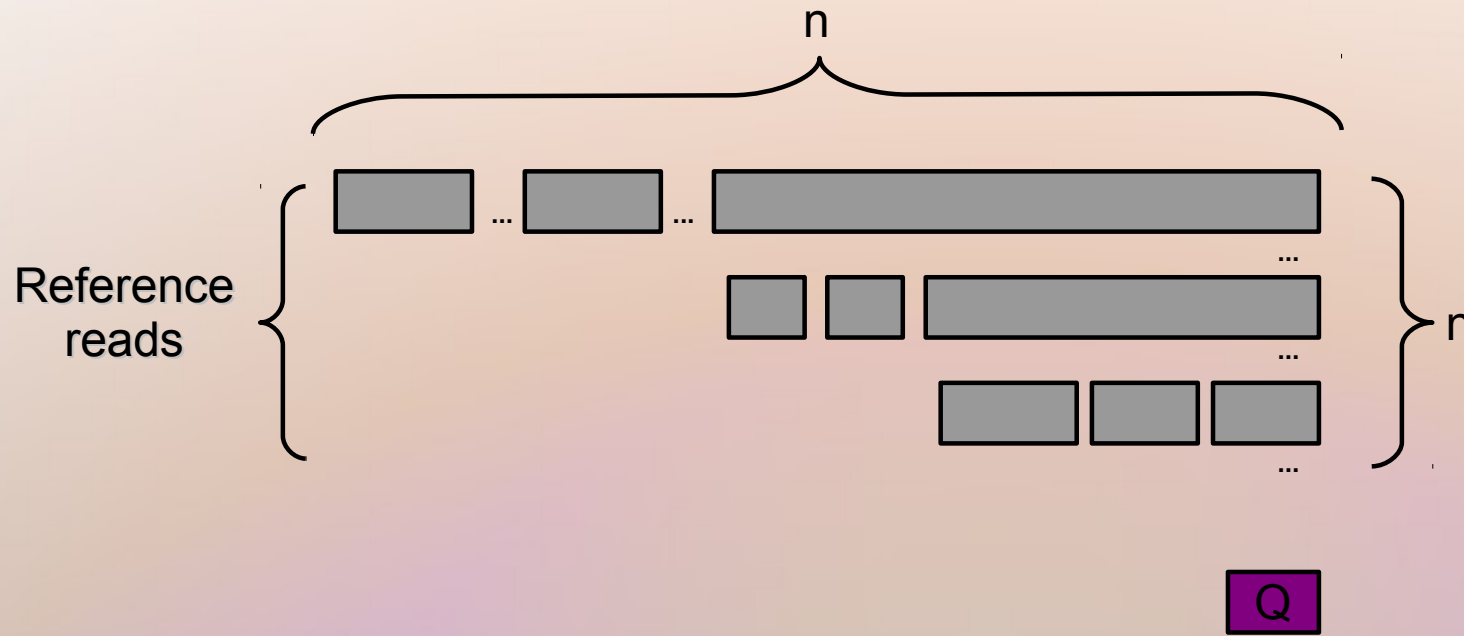
# Nested Containment List (NC-List)



$O(n+\log N)$ ,  
 $N$  = number of reference intervals ;  
 $n$  = number of overlaps



## Example:

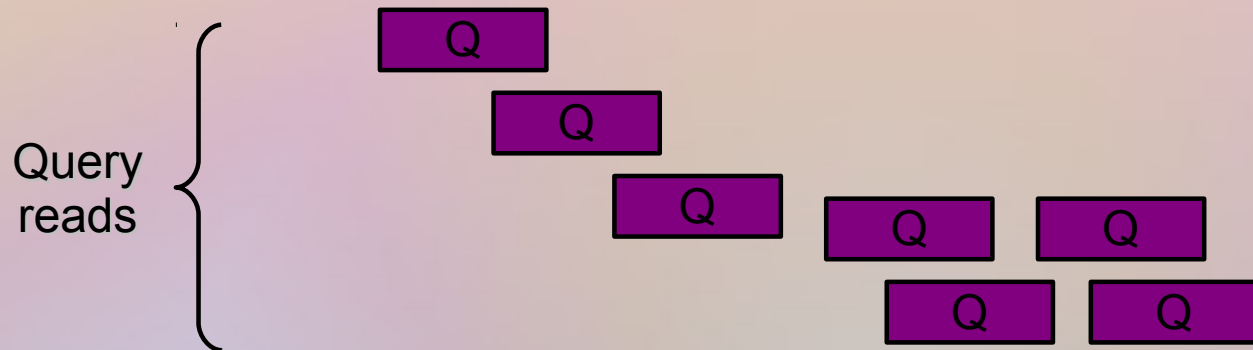
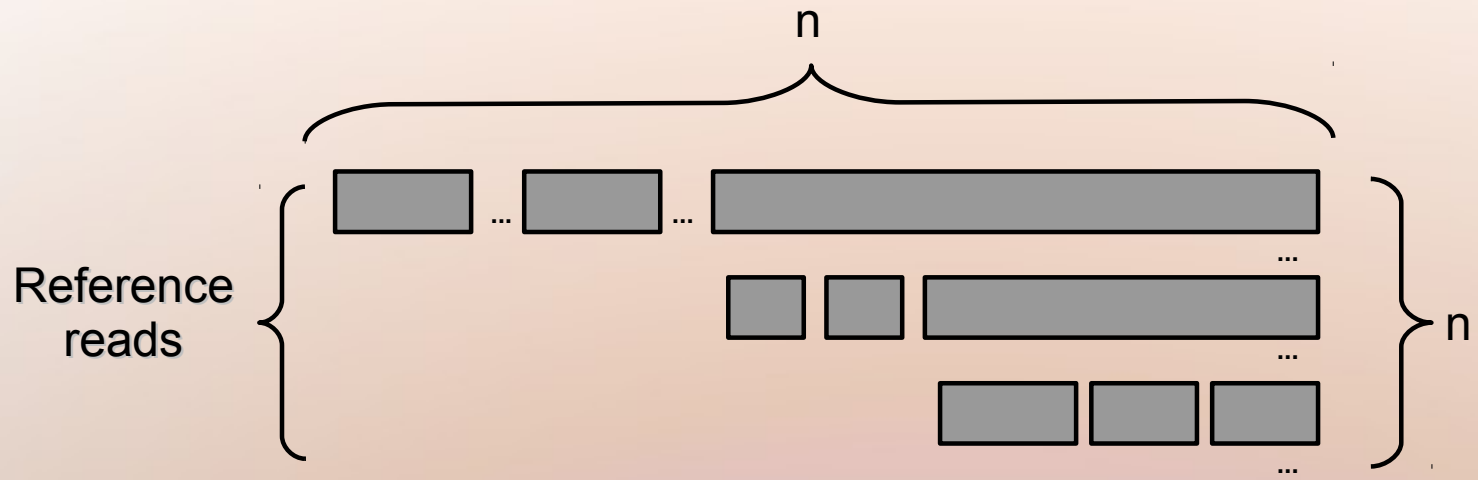


#Reference reads =  $n^2$ ; #Overlaps =  $n$

Complexity =  $O(\log \# \text{Reference reads} + \# \text{Overlaps})$

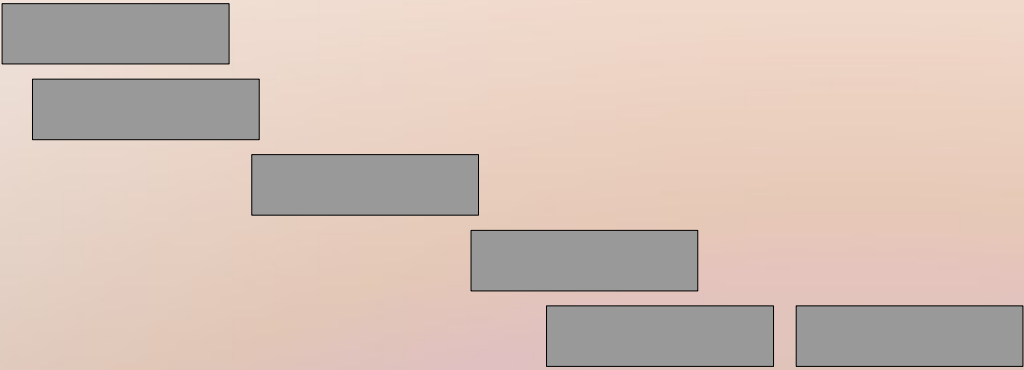
=  $O(\log(n^2) + n) = O(2\log n + n)$

$\sim O(n \log n) \gg O(\log n + n)$  ⚡

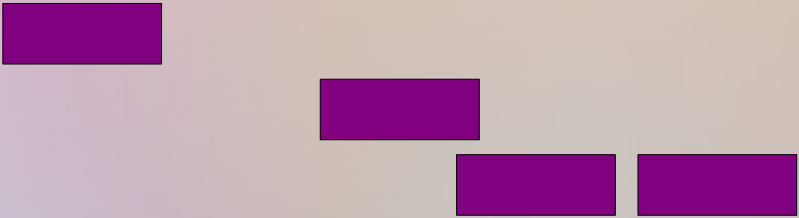


A new algorithm?

New idea



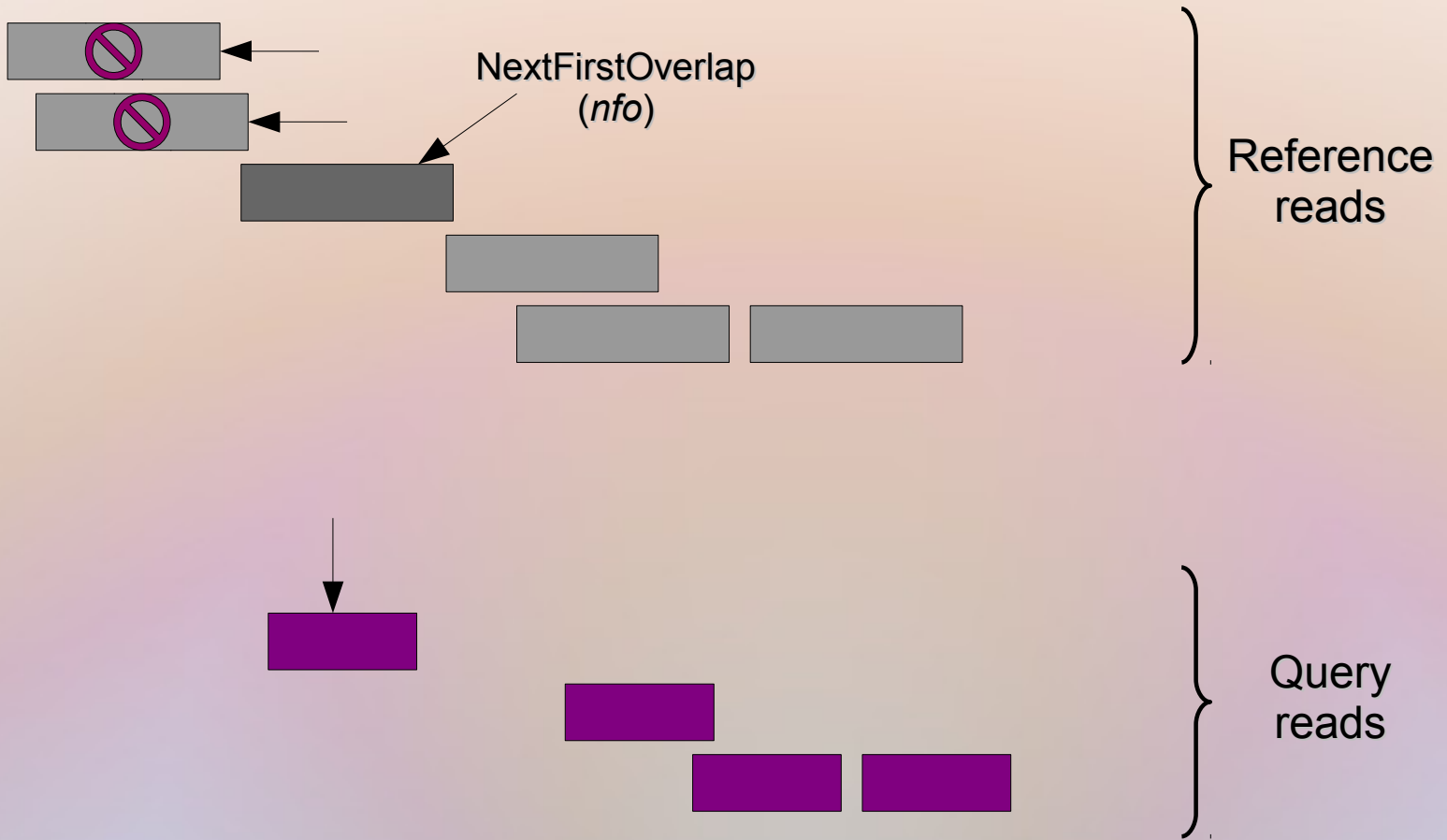
} Reference reads



} Query reads



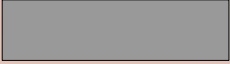
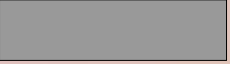
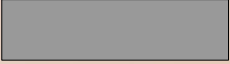
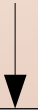
New idea



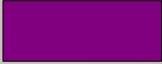
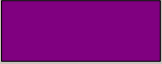
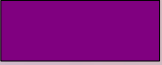
New idea



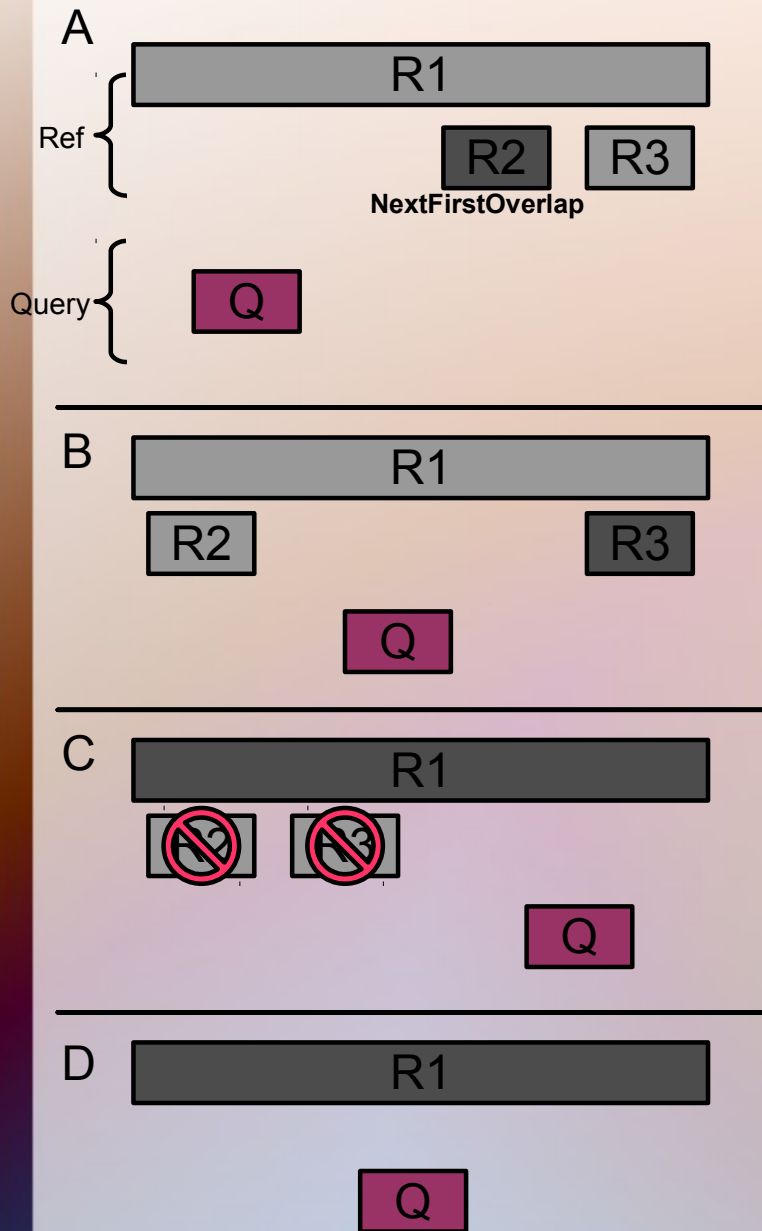
Compare from here !!!



} Reference reads



} Query reads



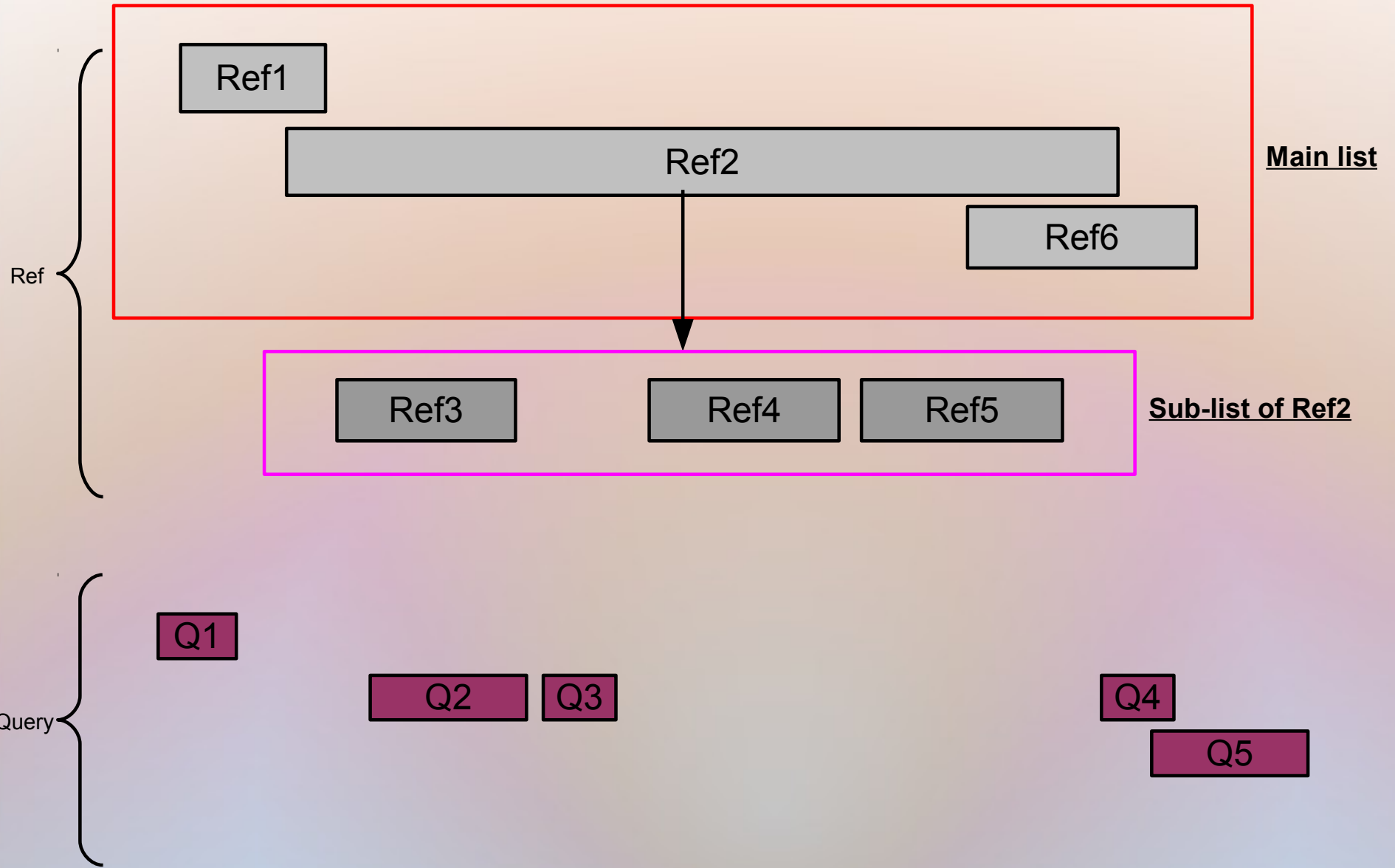
## Function simpleFindOverlap

```

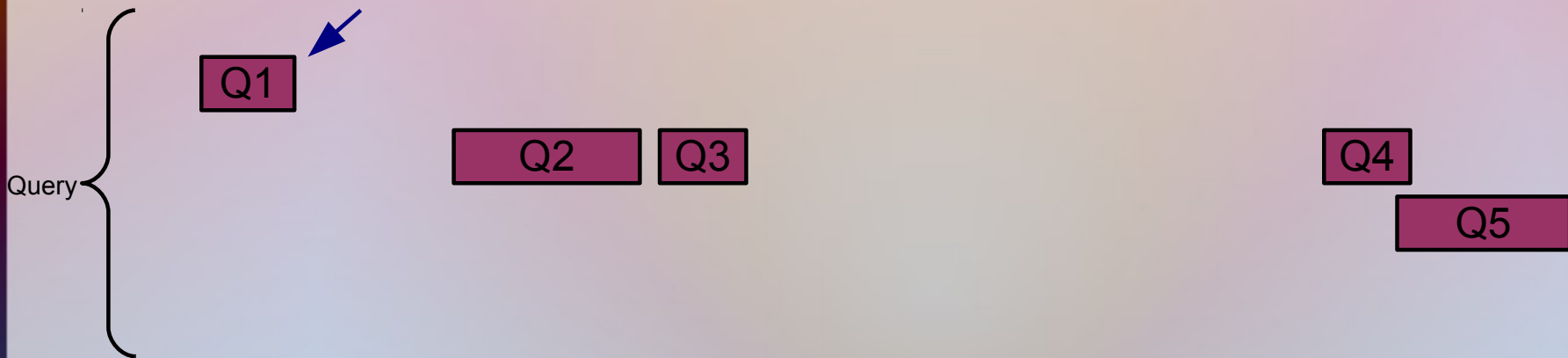
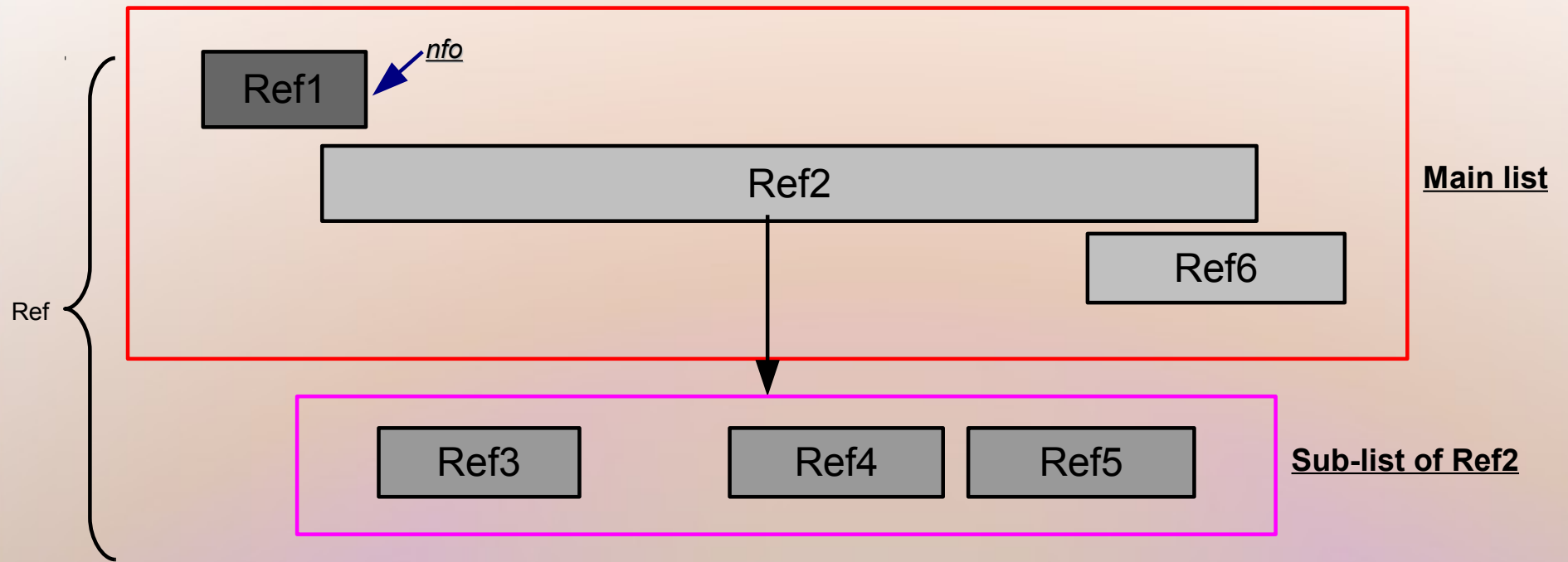
1  $nfo \leftarrow R[0]$ ;
2 foreach  $q$  in  $Q.sorted()$  do
3    $r \leftarrow nfo$ ;
4    $nfo \leftarrow \text{None}$ ;
5   while true do
6     if  $r < q$  then  $r \leftarrow r.next$ ;
7     else if  $r <> q$  then
8       if  $nfo = \text{None}$  then  $nfo \leftarrow r$ ;
9        $M.add(r)$ ;
10       $r \leftarrow r.firstChild$  or  $r.next$ ;
11    else break;
```

!! Compare firstly with all parents of  $nfo$  for the next Query.

Example:

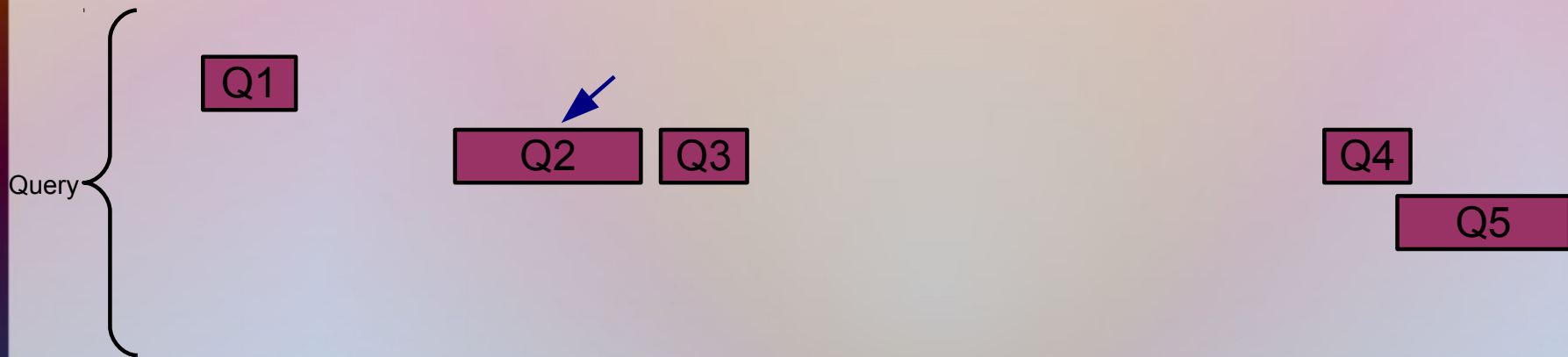
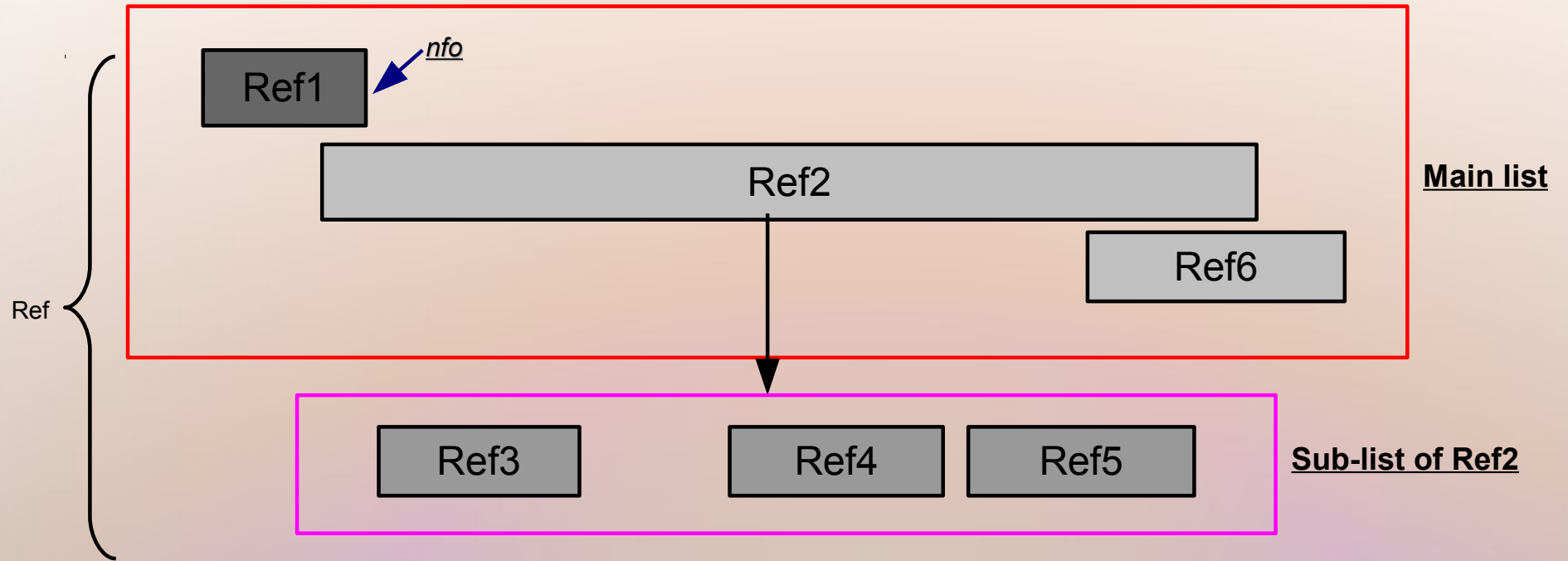


Example:



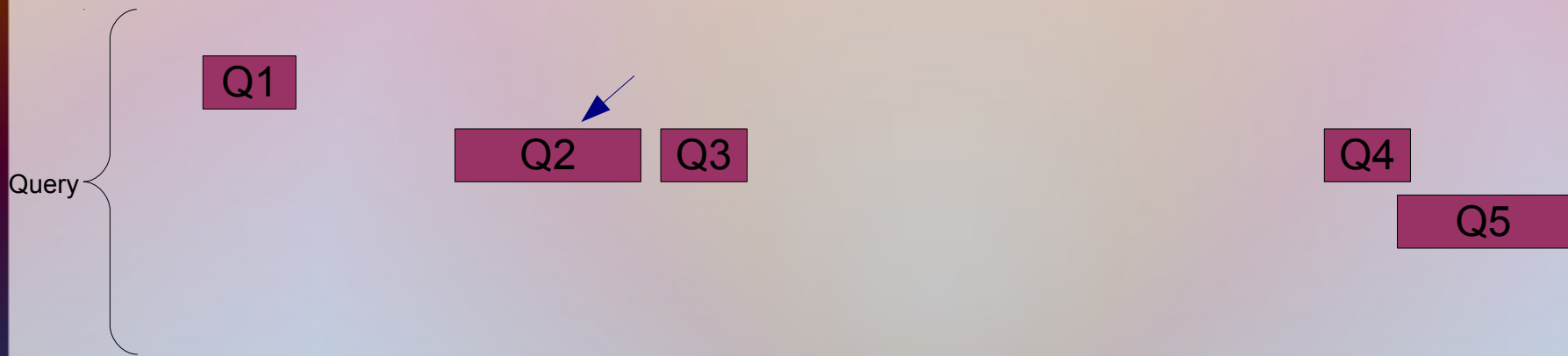
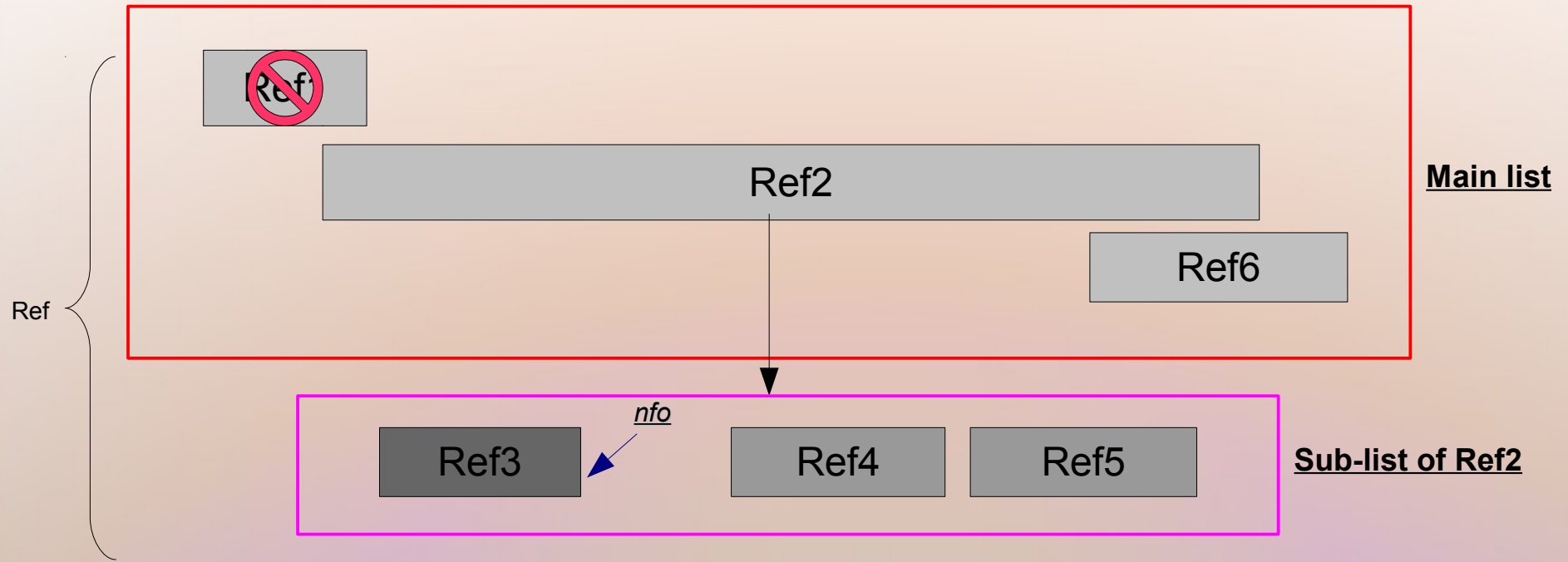
Q1 is overlap with: Ref1

Example:



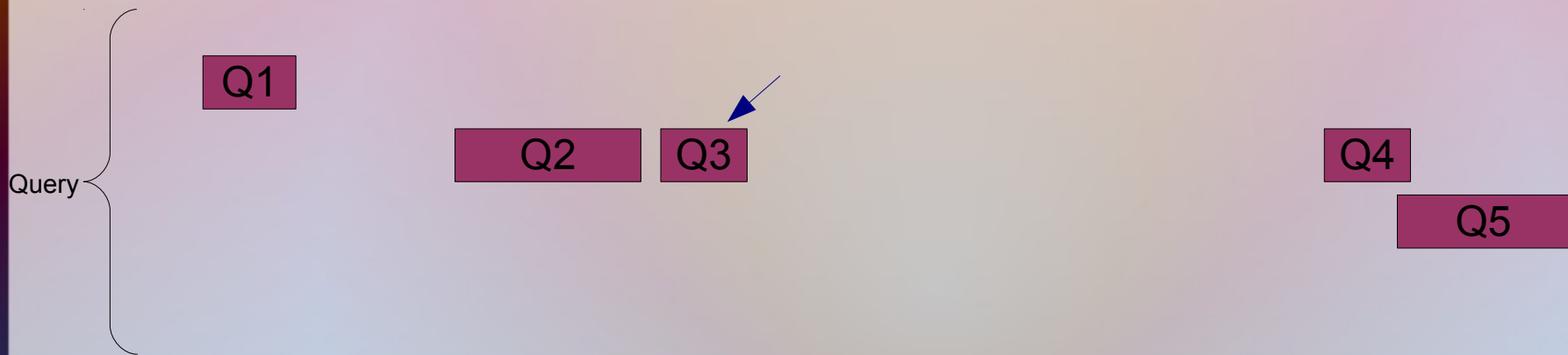
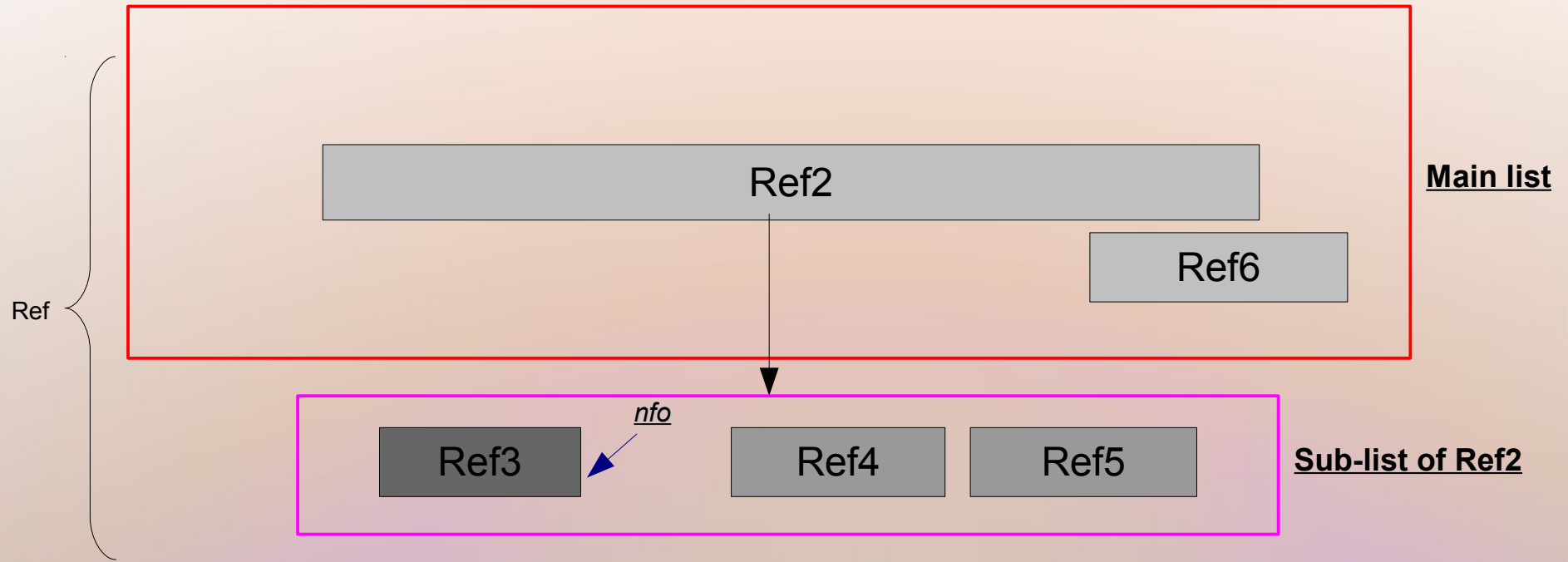
Q1 is overlap with: Ref1

Example:



Q2 is overlap with: Ref2, Ref3

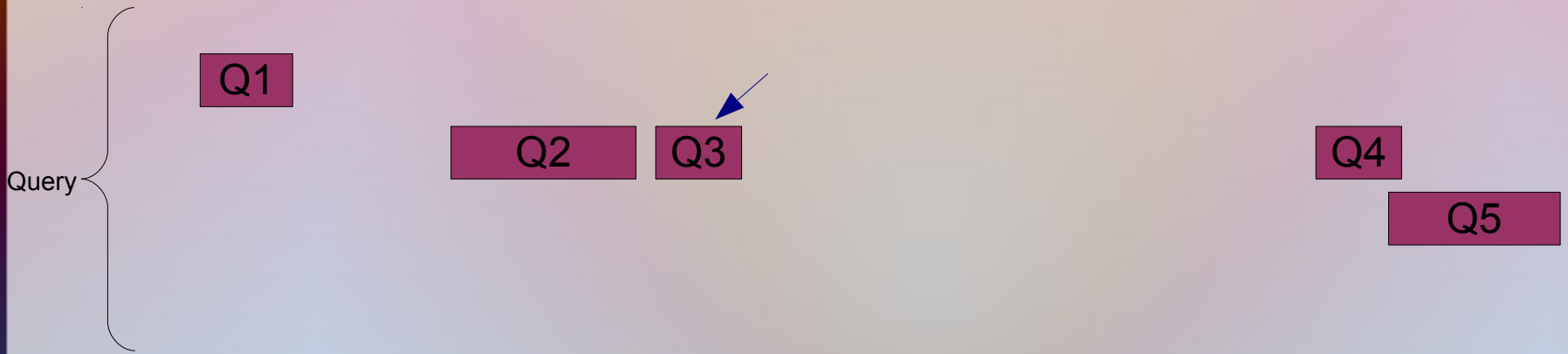
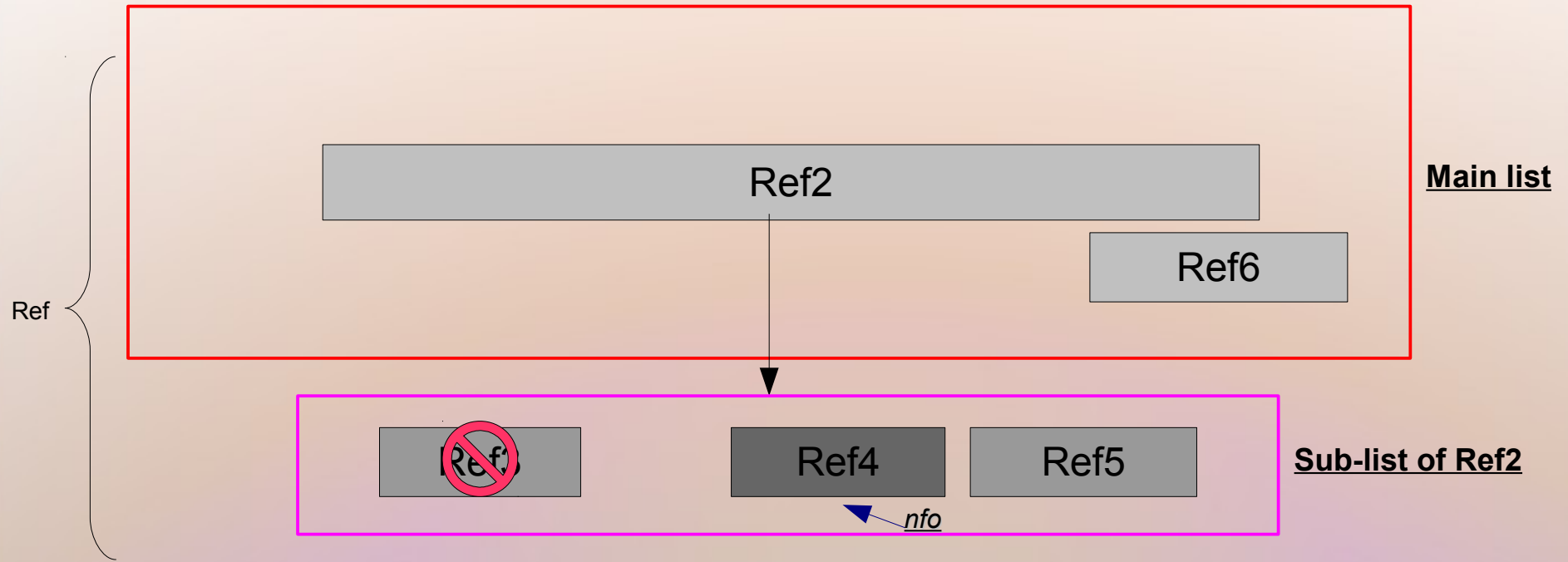
Example:



Q2 is overlap with: Ref2, Ref3

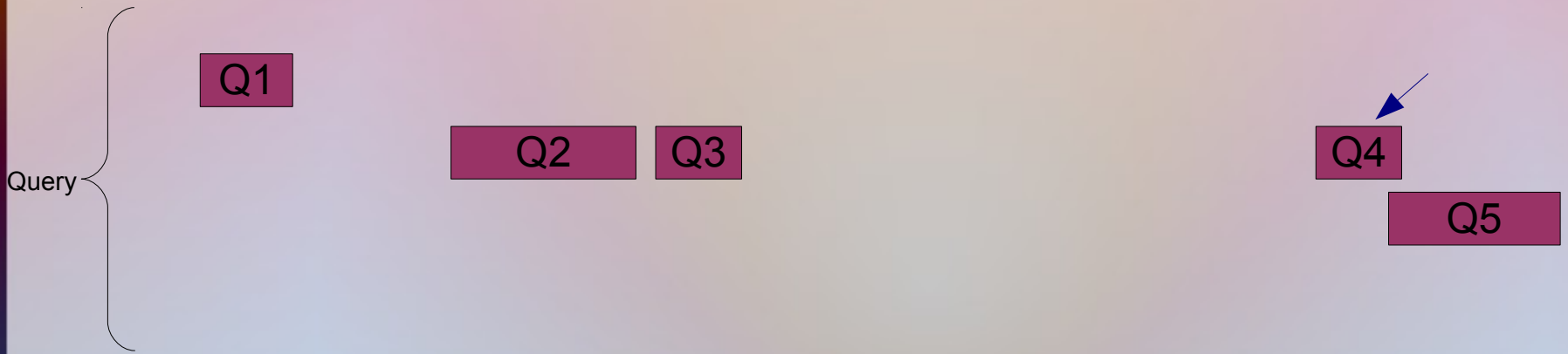
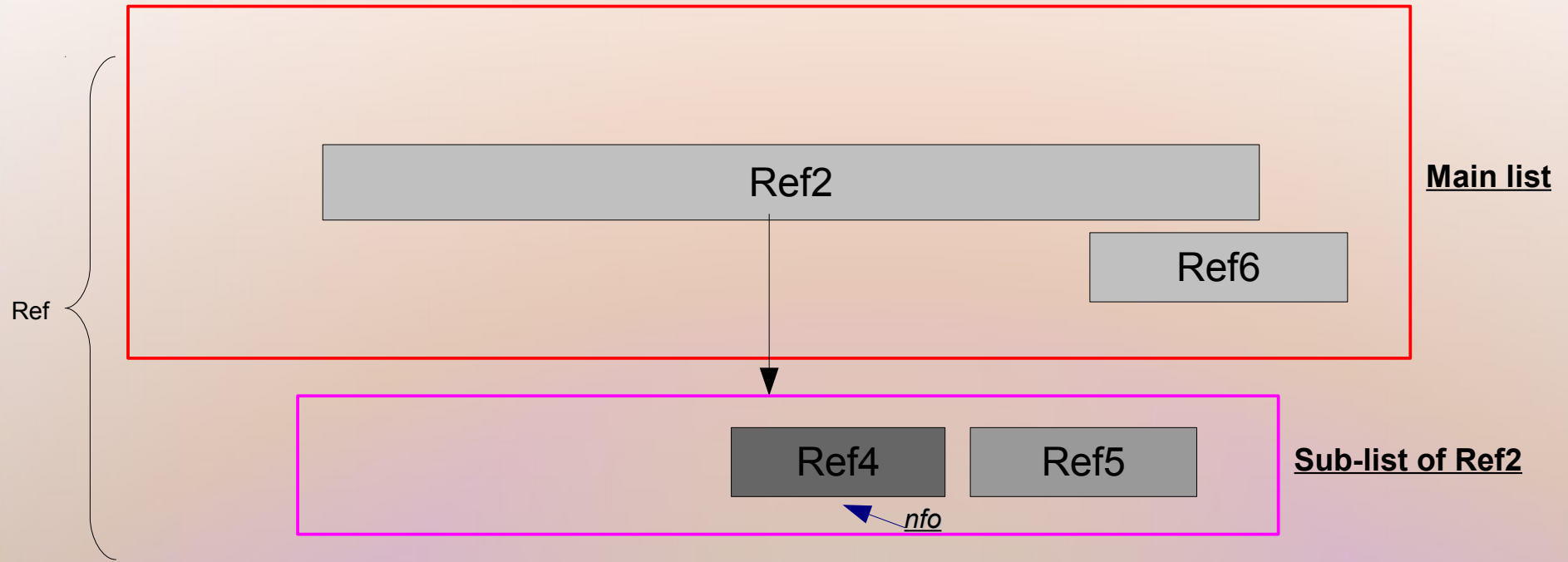


Example:



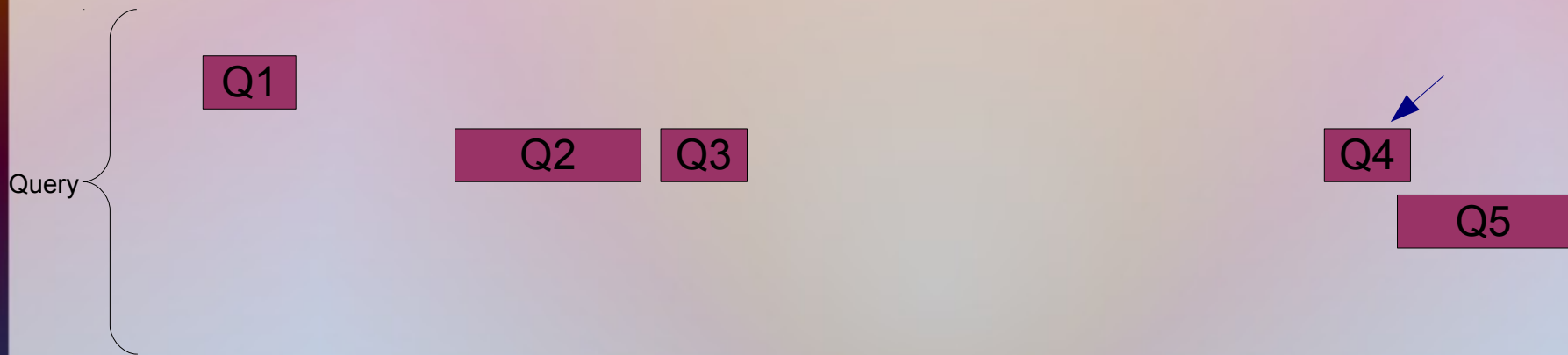
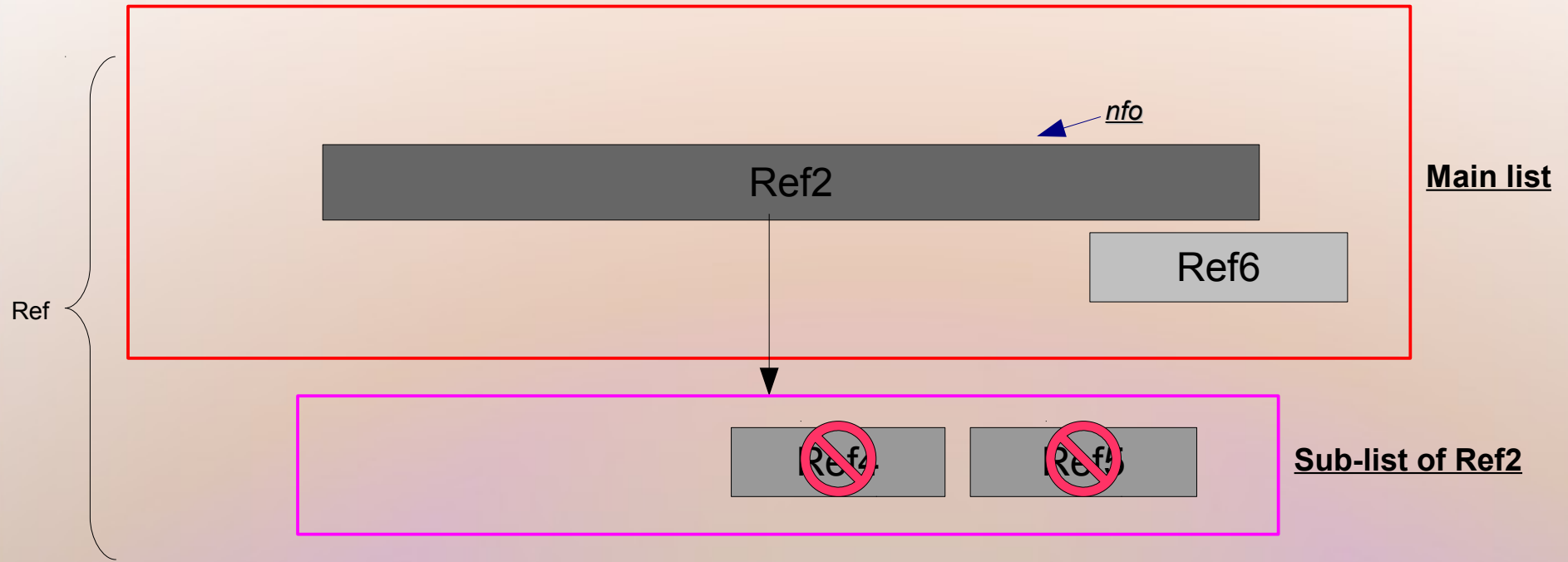
Q3 is overlap with: Ref2

Example:



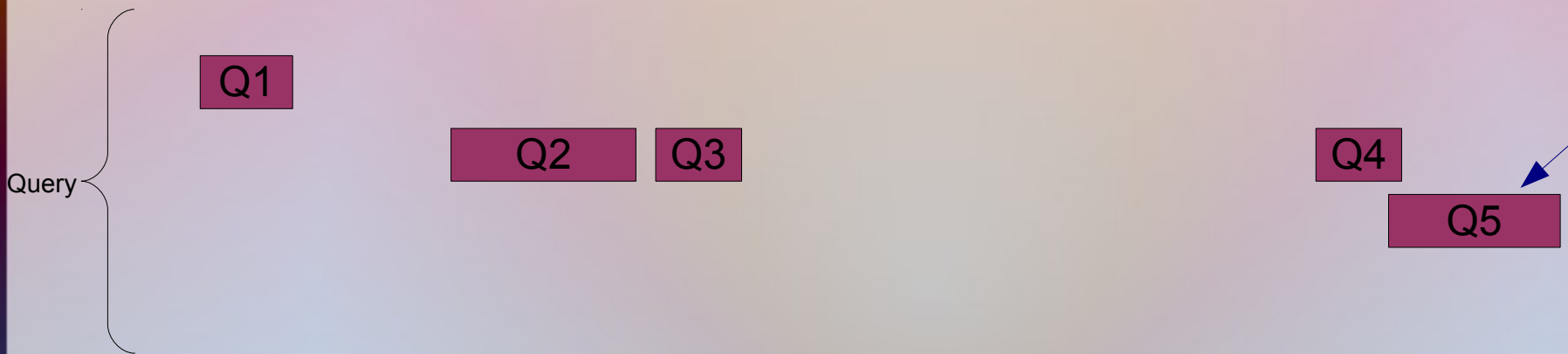
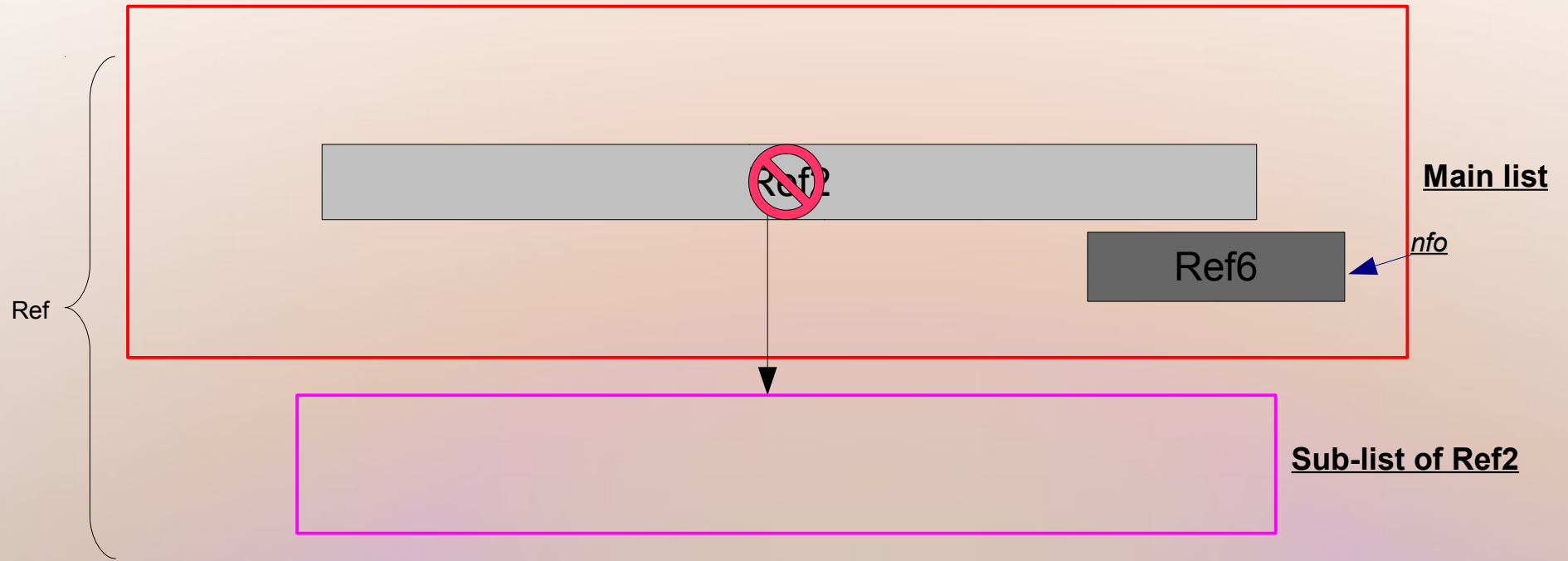
Q3 is overlap with: Ref2

Example:



Q4 is overlap with: Ref2, Ref6

Example:



Q5 is overlap with: Ref6

## The complexity:

$$O(\#R + \#Q + \#M)$$

*Where R is the set of reference intervals;*

*Q is the set of query intervals;*

*M is the set of reference intervals overlapping with  
the query intervals.*

# Results

## *Compare with binning algorithm*

Organism	Reference: RefSeq genes	Query: reads	Number of overlaps	Results
Drosophila melanogaster	18453	24274 454 technology, unique mapping	37761	Our algorithm: 10s
		292085 454 technology, Non-unique mapping	467815	Binning: 29s
		3013398 Illumina Unique mapping	3716626	Our algorithm: 1m41s
				Binning: 4m27s
				Our algorithm: 15m18s
				Binning: 46m45s

# perspectives

Comparaison with other tools  
Implementation in S-MART<sup>1</sup> and Galaxy

<sup>1</sup> Matthias Zytnecki and Hadi Quesneville, *PloS ONE*, 2011

# acknowledgment

**APLIBIO**



**eBIO**



**INRA**

