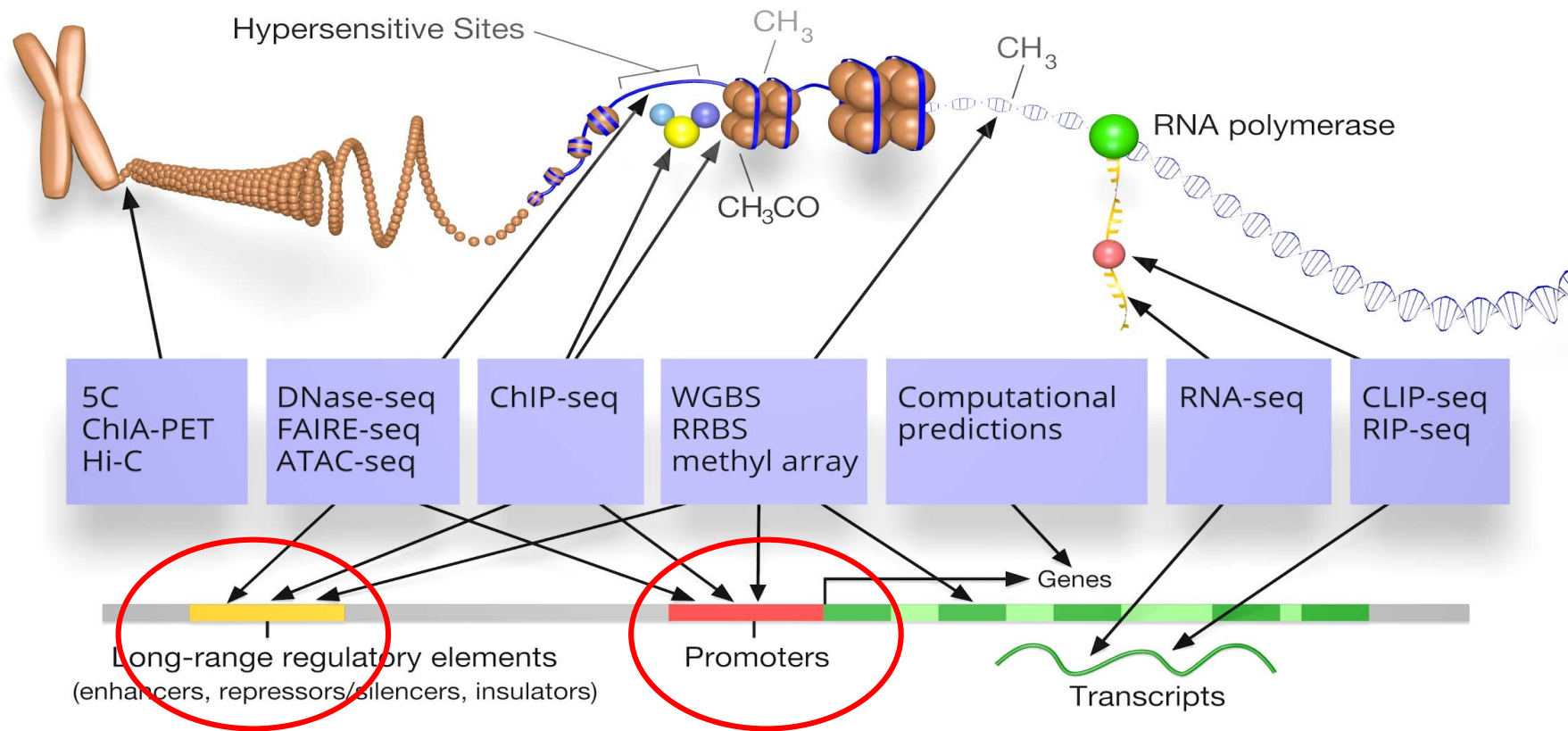


# Analisi di dati di ENCODE genomici ed epigenomici

Linguaggi di Programmazione:

- SQL
- R
- Python



Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

# Tipi di Dato

| Tables in snpImpactResource

+-----+

| CellLines

| FunctionalElement

| GeneSymbols

| GeneSymbols\_Symbols

| INDELS

| INDELS\_FunctionalElement

| INDELS\_PFM

| MotifDatabases

| PFM

| PFM\_GeneSymbols

| PFM\_MotifDatabases

| SNPs

| SNPs\_FunctionalElement

| SNPs\_PFM

| Symbols

+-----+

Tutte le tabelle del DataBase

In rosso sono le tabelle utilizzate

Server: MariaDB

Container: Singularity

# Dati quantitativi iniziali

- SNPs = 14810175
- Elementi Funzionali = 759
- CellLines = 238
- PFM = 5424 (SNPs\_PFM = 5352)

# 3 esempi di records

## SNPs

SNPID	rsid	chrom	pos	ref	alt	maf1000genomes	mafTOPMed
1	rs575272151	1	11008	C	G	NULL	NULL

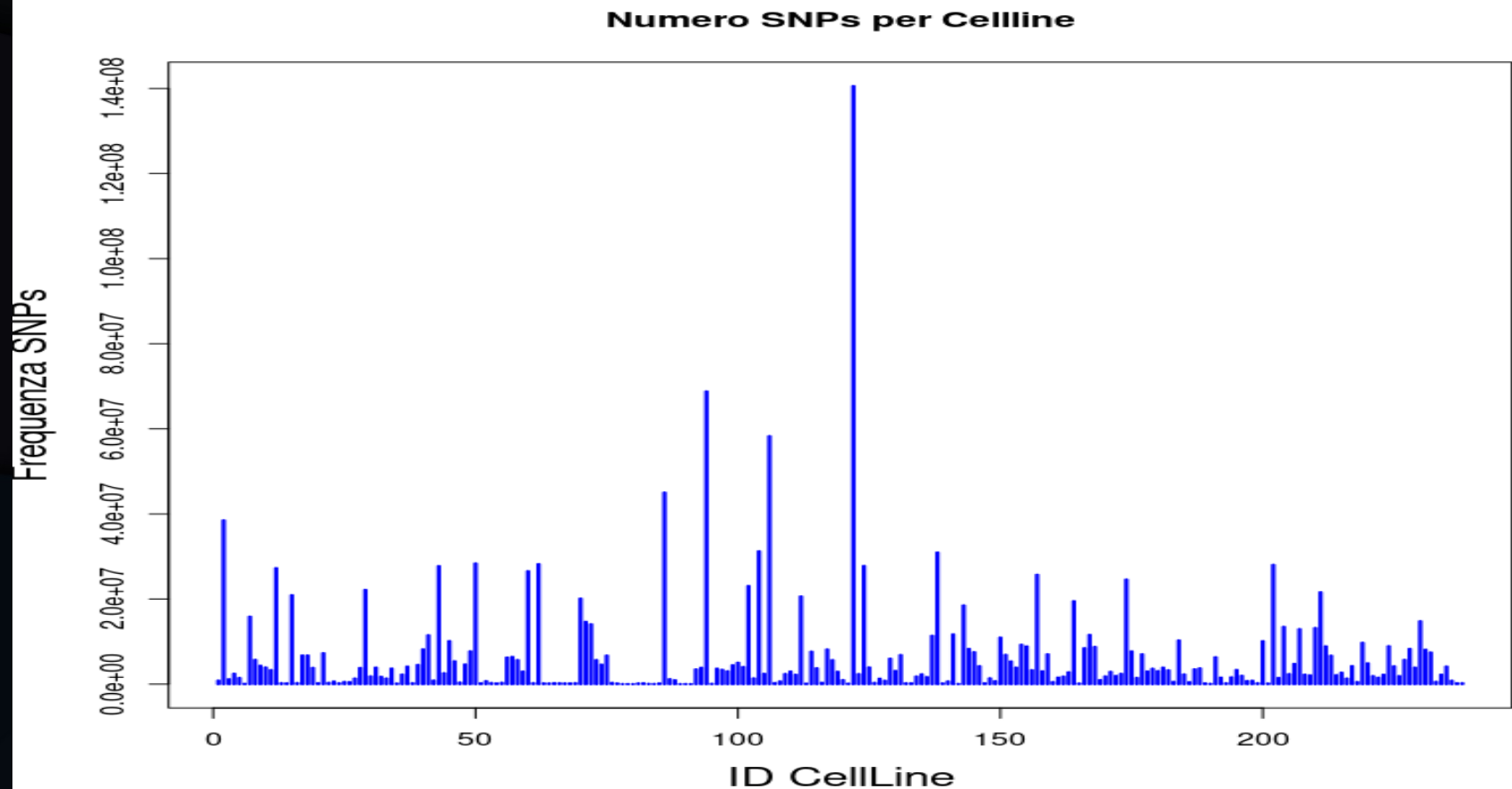
## SNPs\_FunctionalElement

SNPID	ElementID	CELLLINEID	countExperiments	fileType
1	33	122	1	broad

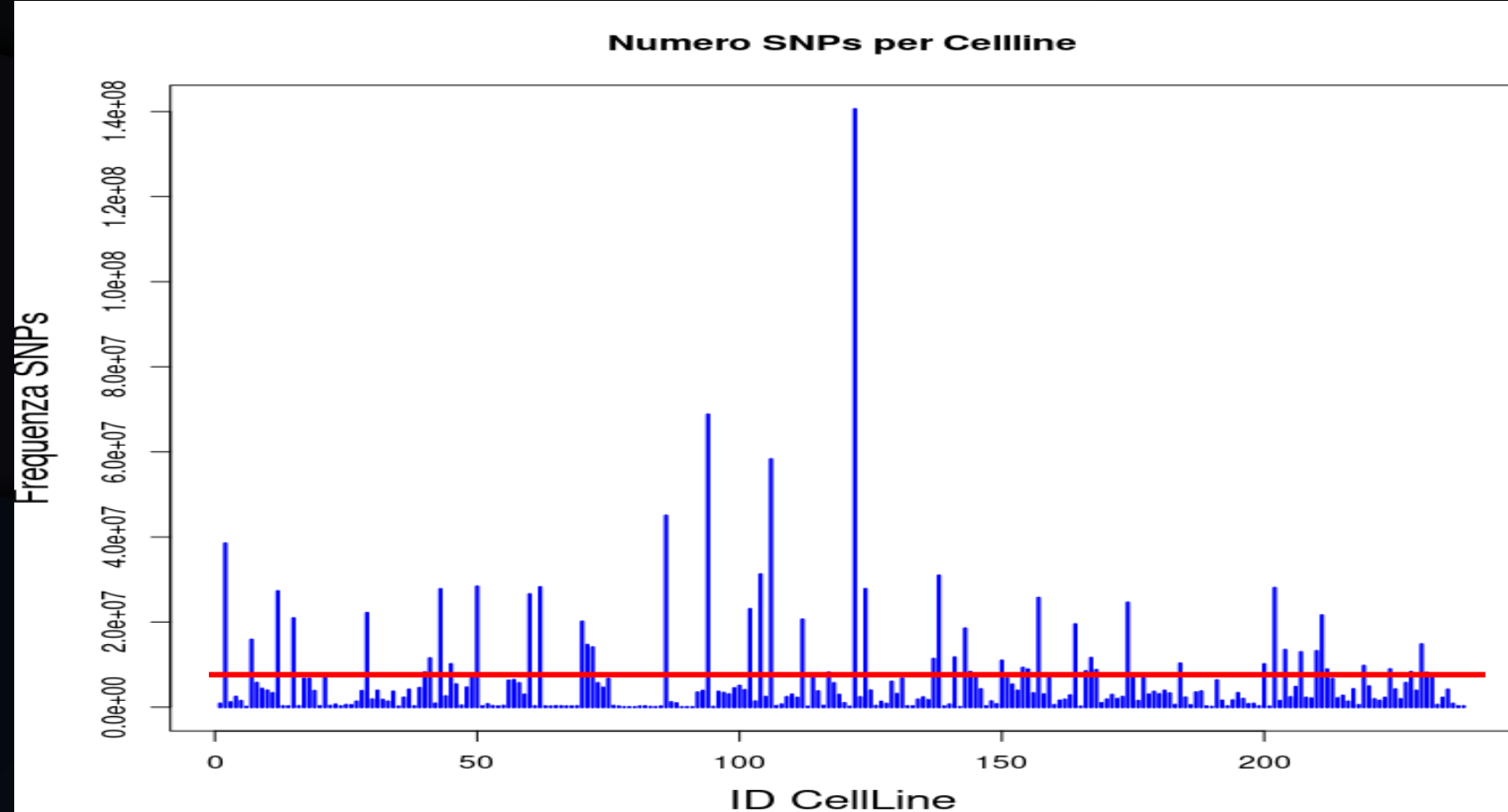
## PFM

SNPID	PFMID	start	strand	type	scoreRef	scoreALT
1	34384	11007		match	6.22488	6.80984

# Prima Query fatta al DataBase

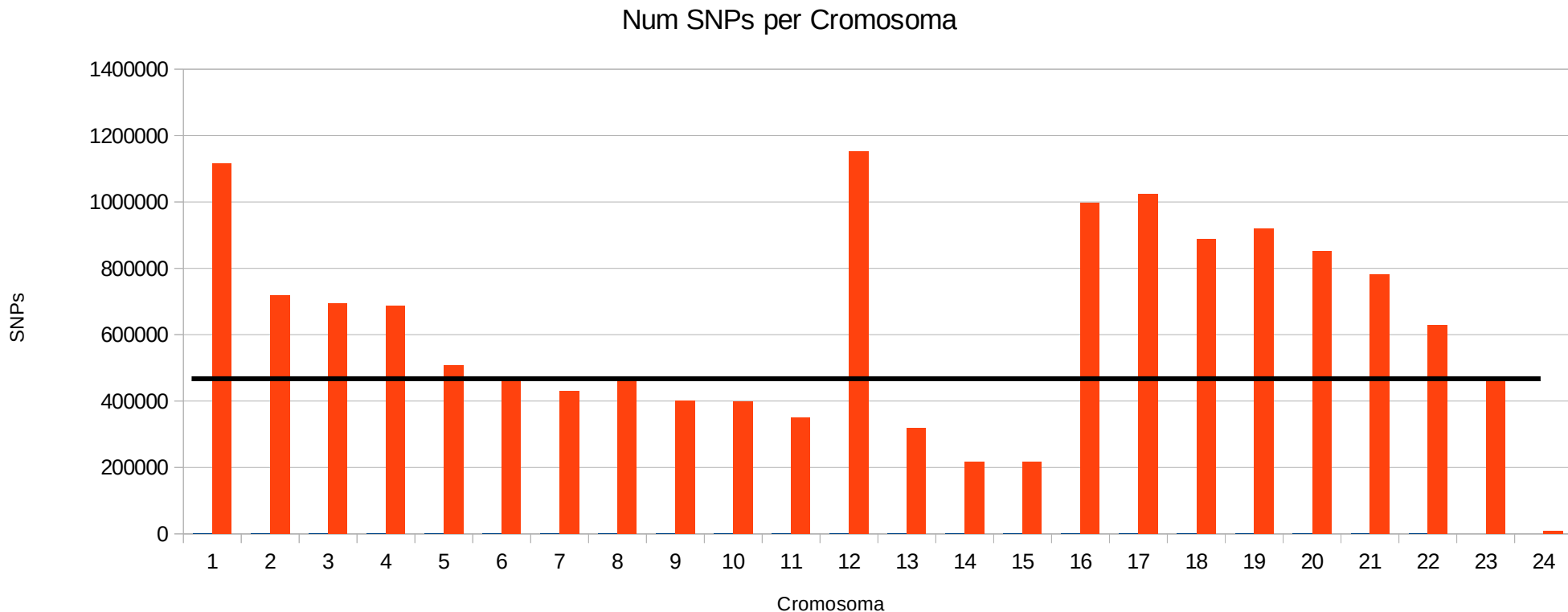


# Creare threshold di riferimento per i ricercatori



# SNPs per Chr

In questo caso non discrimino fra SNPs diversi, ne' per CellLine diverse

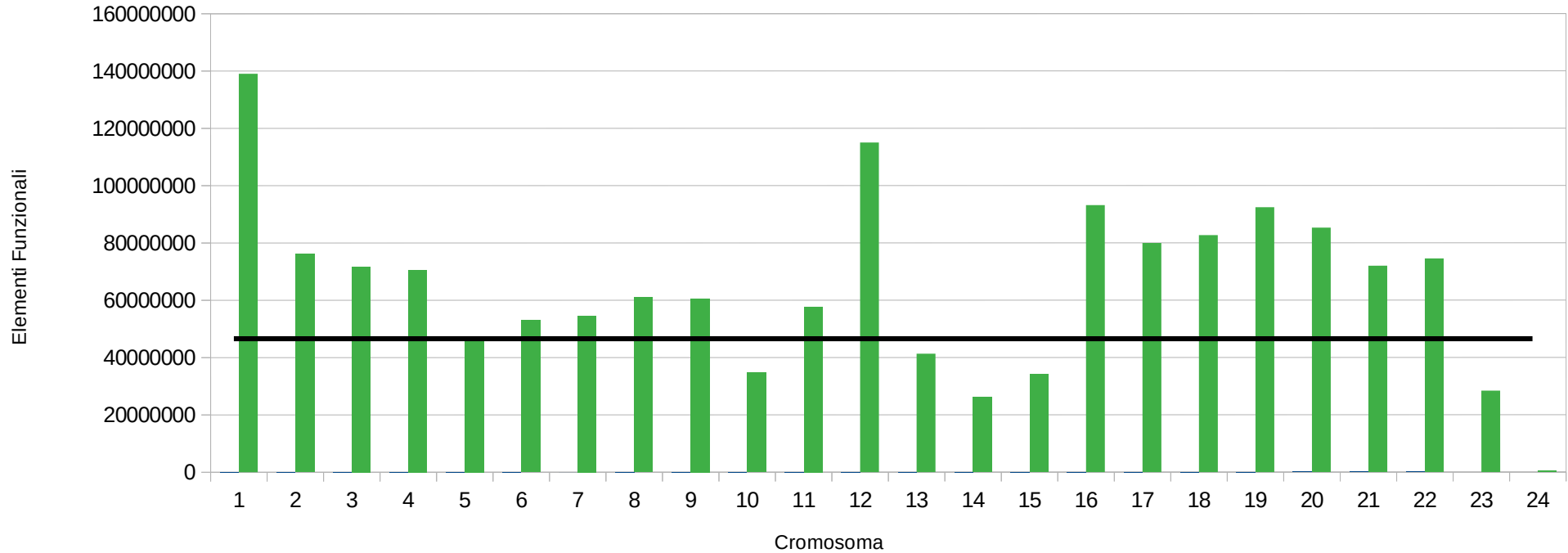


\* Chr 23=X, 24=Y

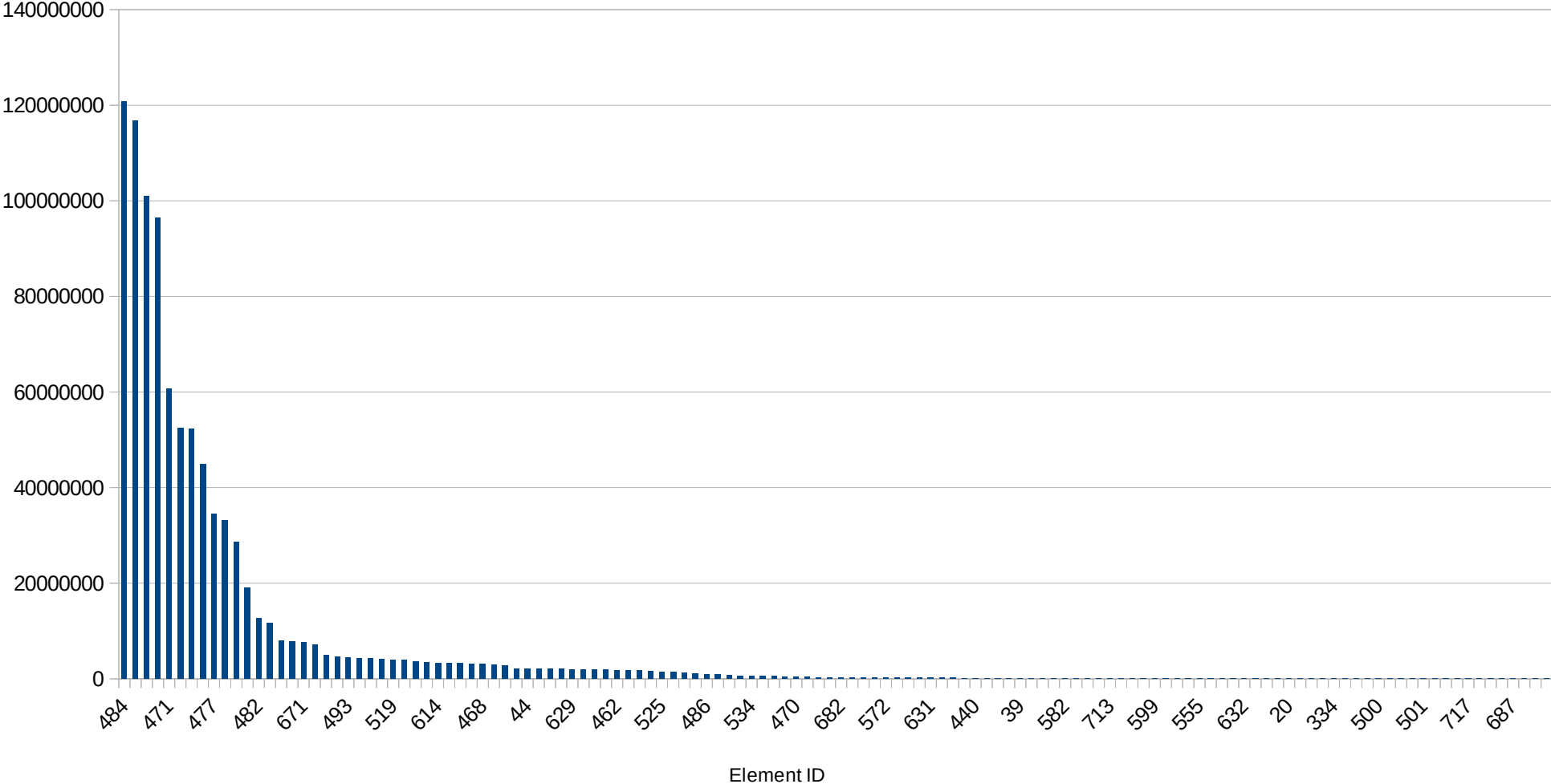


# Elementi Funzionali per Chr

Numero di Elementi Funzionali per Cromosoma



Numero di SNPs per Elemento Funzionale



# Esempio dei filtri per i dati

- File Type: broad / narrow or match/change/change
- Contare distinti elementi oppure contare la somma delle frequenze dei singoli elementi
- Non mettere i valori NULL
- Discriminare per la lunghezza delle PFM Probability Frequency Matrix

# Possibile filtro dati per i PFM

| type | file | scoreRef | scoreALT | name Element | length | maxScore | SNPs ID | chrom | pos

Uso di 3 tabelle: SNPs\_PFM, PFM, SNPs

Per ogni PFM e SNP ho uno score di riferimento, uno score massimo teorico ed uno score Alternativo nella cellula mutata.

Voglio filtrare gli score piu'interessanti con:

$| ( (\text{ScoreALT} - \text{ScoreRef}) / \text{maxScore}) * 100 | > 10 \%$

Filtro gli score che differiscono del + o - 10% rispetto a quello di riferimento.

SNPID	scoreRef	scoreALT	maxScore	PFMID
1	6.22488	6.80984	9.62551	34384

## Possibili informazioni utili ricavabili DataBase

Marcare le regioni dei cromosomi ad alta frequenza di mutazioni di SNPs.

Capire se un particolare SNPs (ex C → G) incide di più su alcuni geni piuttosto che su altri

Marcare le posizioni dei cromosomi a cui “si attaccano” i differenti Elementi Funzionali