



GMQL in action to solve biological problems

Pietro Pinoli

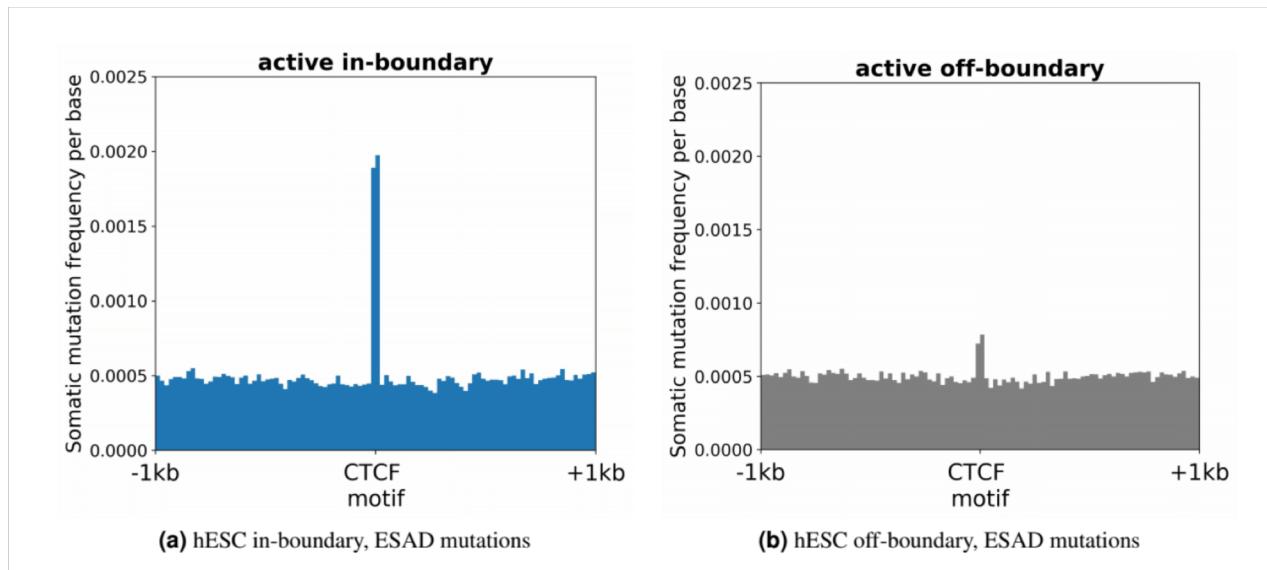
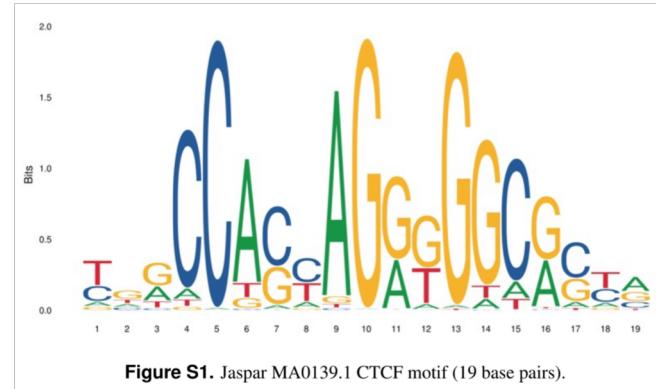
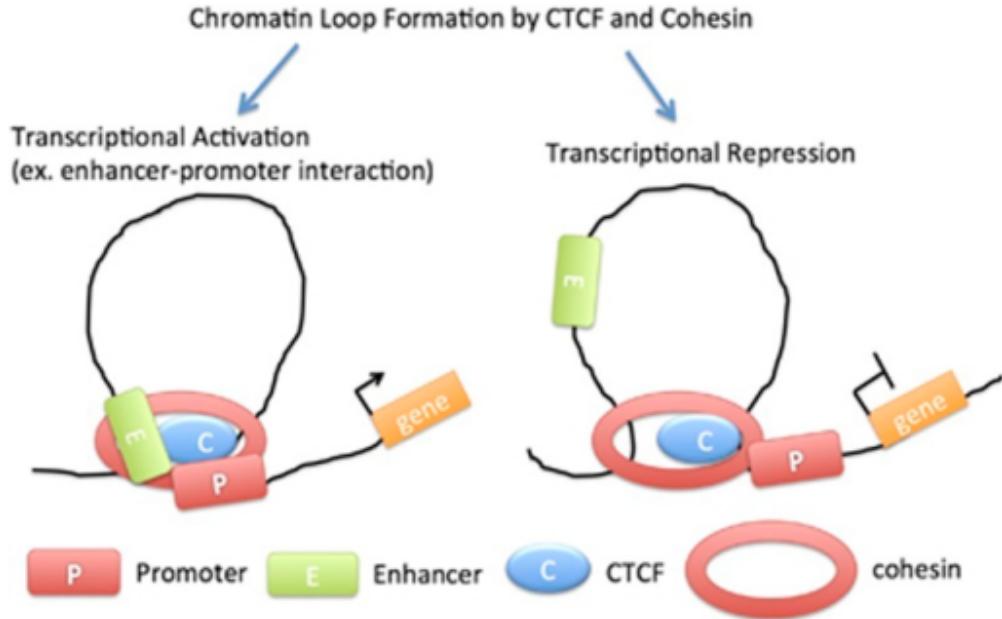
Stefano Perna

Luca Nanni

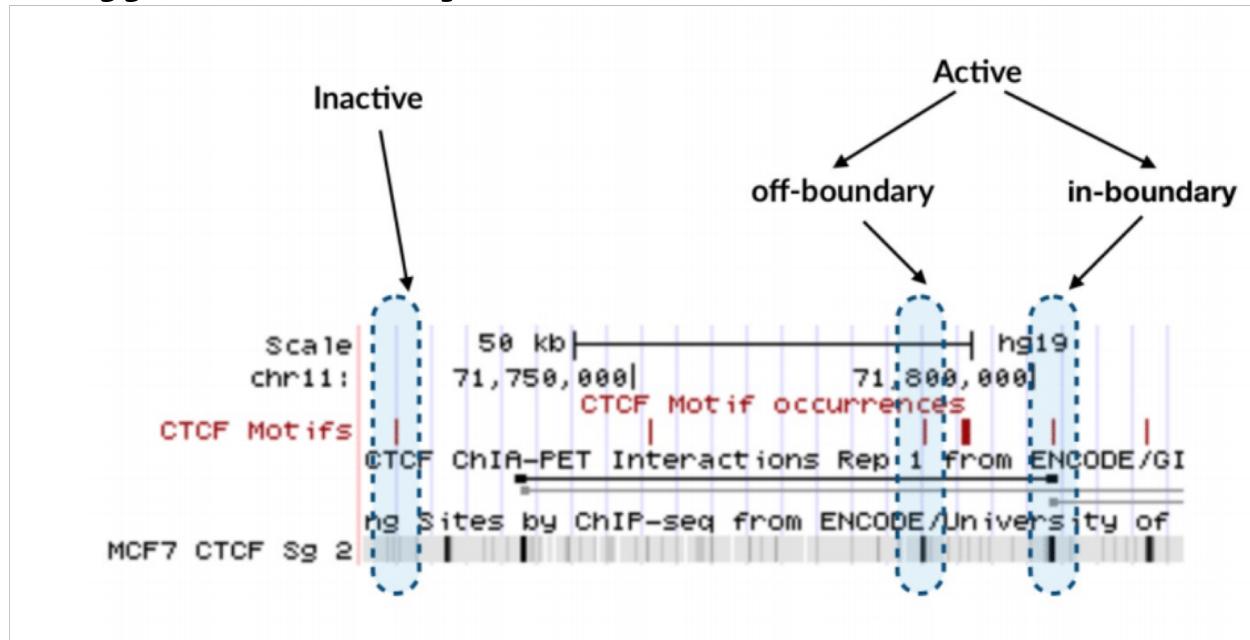
Pietro Pinoli

Pan-cancer analysis of somatic mutations and epigenetic alterations in insulated neighbourhood boundaries.

Problem Statement

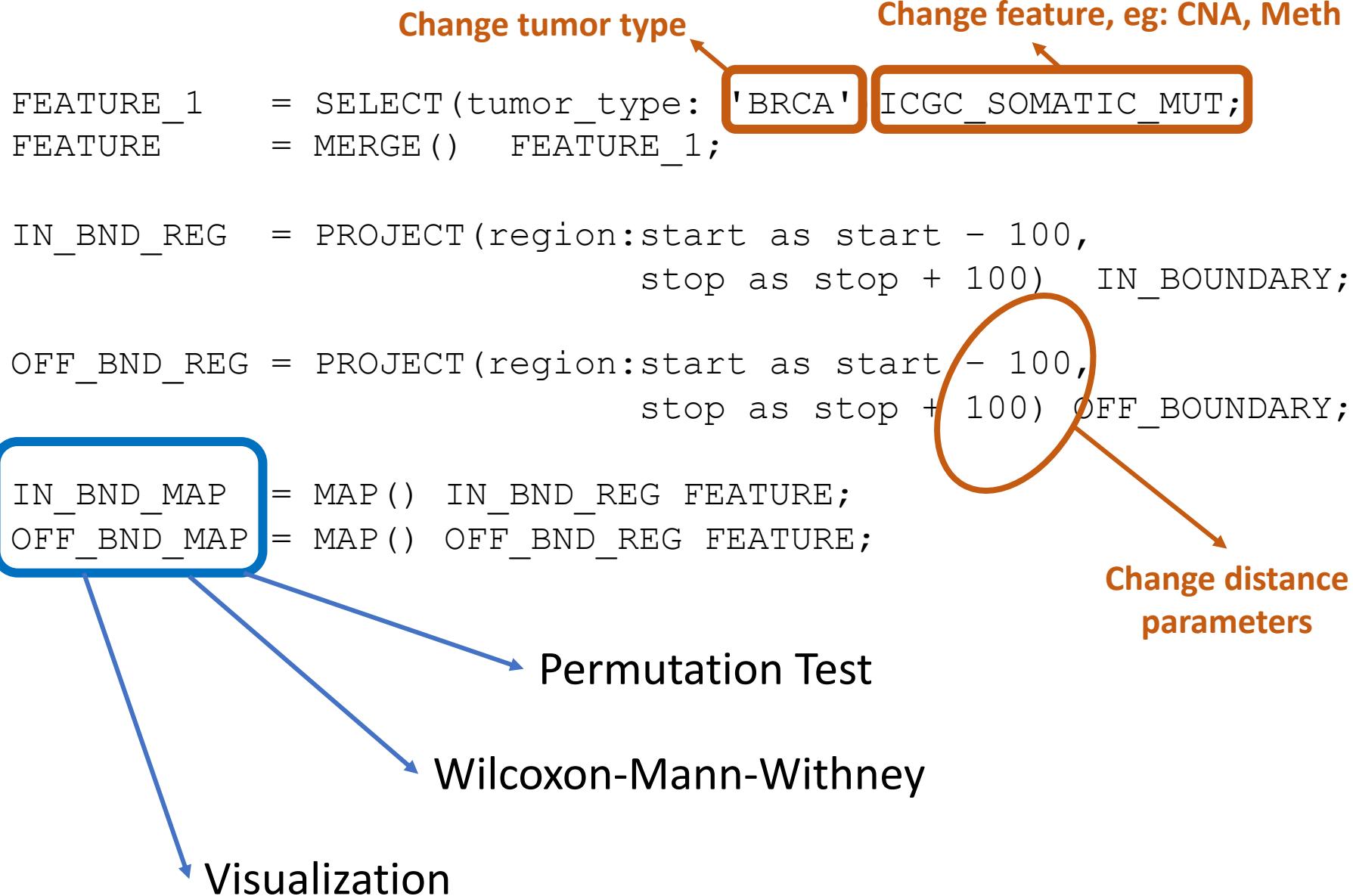


Split CTCF motifs in *active in-boundary* and *active off-boundary*



```
MOTIFS      = SELECT() MOTIFS_BioStrings;
CHIP        = SELECT(cell_line == 'MCF7') ENCODE_NARROW;
ACTIVE       = JOIN(distance < 0; output: left) MOTIFS CHIP;
BOUNDARIES   = SELECT(cell_line == 'MCF7') CHIA_PET;
OVERLAP      = MAP() ACTIVE BOUNDARIES;
IN_BOUNDARY  = SELECT(region: count > 0 ) OVERLAP;
OFF_BOUNDARY = SELECT(region: count == 0) OVERLAP;
```

Differential analysis *in-boundary* vs. *off-boundary*



Results

With ~ 20 lines of GMQL query,
Sistematic analysis on:

- 3 sets of neighbourhoods
- 26 tumor types
- 3 genomic features
(mutations, DNA Methylation,
Copy Number Variation)

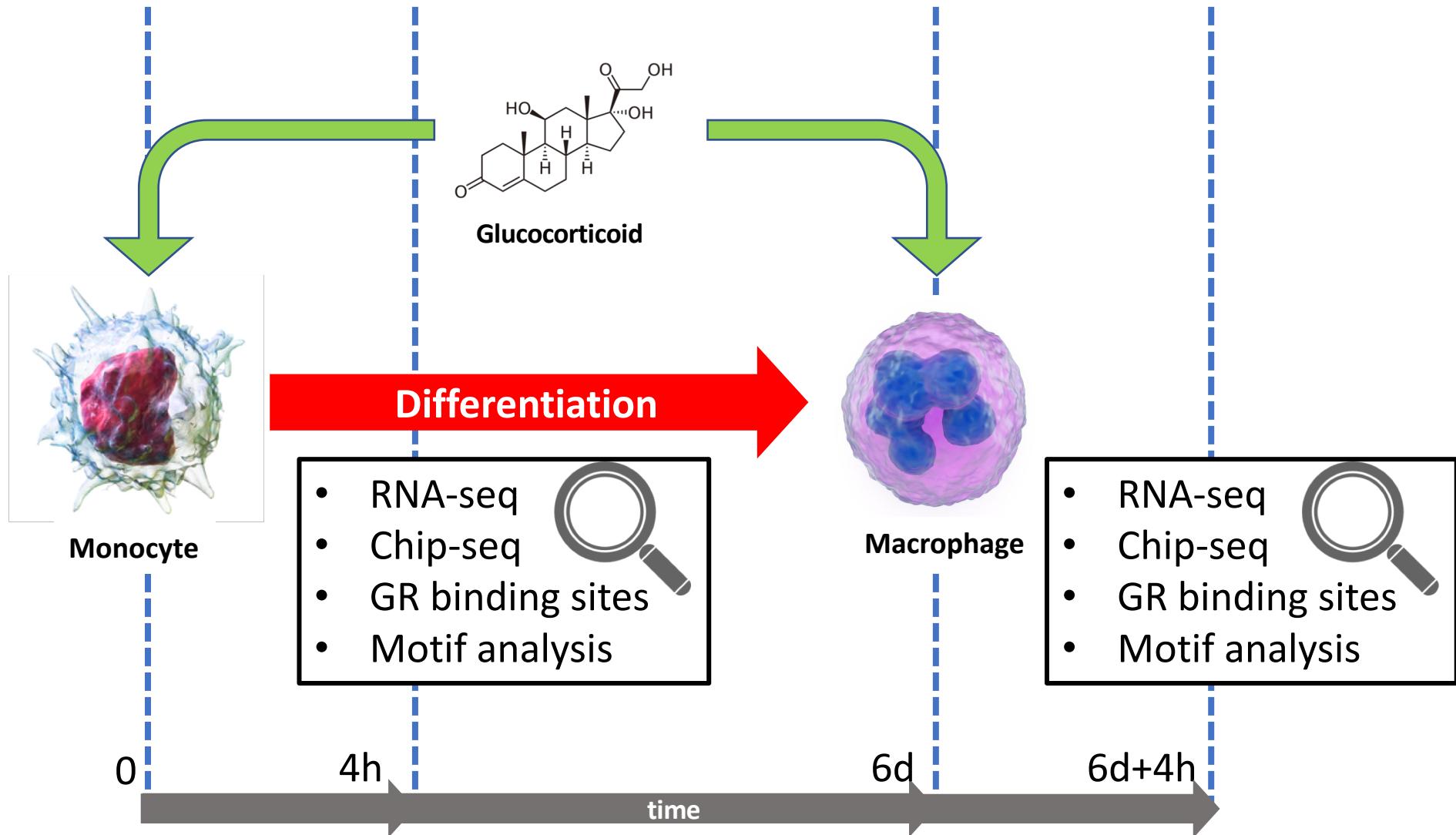
	Somatic mut.			DNA meth.			CNA		
Tumour	h	M	H	h	M	H	h	M	H
BLCA	-	-	-	N	N	N	N	N	N
BOCA	N	N	Y	-	-	-	-	-	-
BRCA	Y	Y	Y	Y	Y	Y	Y	Y	Y
BTCA	N	Y	N	-	-	-	-	-	-
COCA	N	N	N	-	-	-	-	-	-
EOPC	N	N	Y	-	-	-	-	-	-
ESAD	Y	Y	Y	-	-	-	-	-	-
GACA	Y	Y	Y	-	-	-	-	-	-
GBM	-	-	-	-	-	-	N	Y	N
HNSC	-	-	-	Y	Y	Y	N	N	N
KIRC	-	-	-	N	N	N	-	-	-
KIRP	-	-	-	-	-	-	N	N	N
LIHC	-	-	-	N	N	N	Y	Y	Y
LIRI	Y	Y	Y	-	-	-	-	-	-
LUAD	-	-	-	N	N	N	Y	Y	Y
LUSC	-	-	-	-	-	-	Y	Y	Y
MALY	N	N	N	-	-	-	-	-	-
MELA	Y	Y	Y	-	-	-	-	-	-
OV	N	Y	Y	-	-	-	Y	Y	Y
PACA	N	Y	Y	-	-	-	-	-	-
PRAD	-	-	-	Y	N	N	N	N	N
RECA	Y	Y	Y	-	-	-	-	-	-
SKCA	N	N	Y	-	-	-	-	-	-
SKCM	-	-	-	Y	Y	Y	-	-	-
THCA	-	-	-	N	N	N	N	N	N
UCEC	-	-	-	Y	Y	Y	N	N	N

Luca Nanni

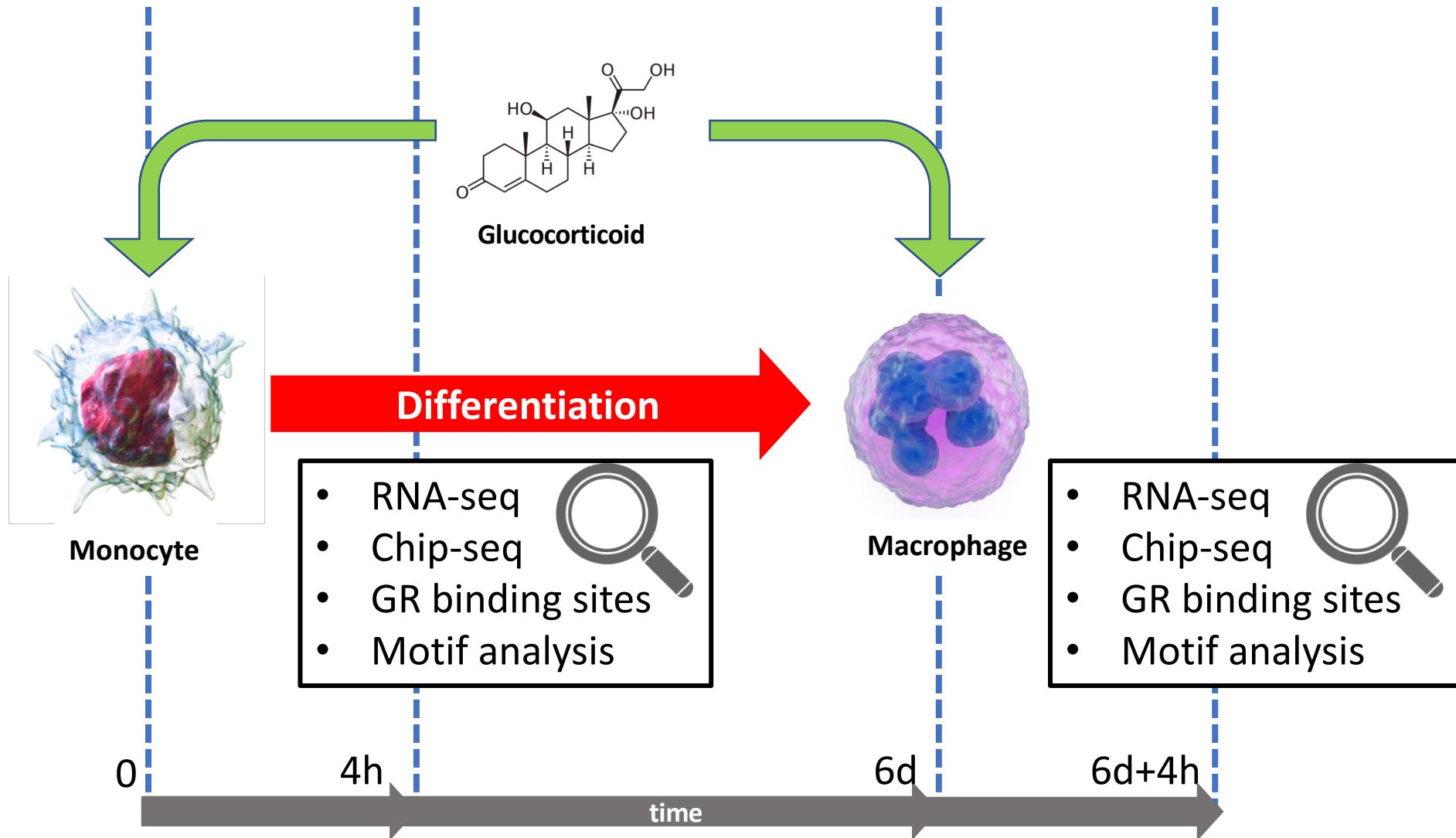
*Extensive epigenomic integration of the
glucocorticoid response in primary human
monocytes and in vitro derived macrophages*

Cheng Wang, Luca Nanni, Boris Novakovic, Wout Megchelenbrink, Tatyana Kuznetsova, Hendrik G. Stunnenberg, Stefano Ceri & Colin Logie

Experimental setup

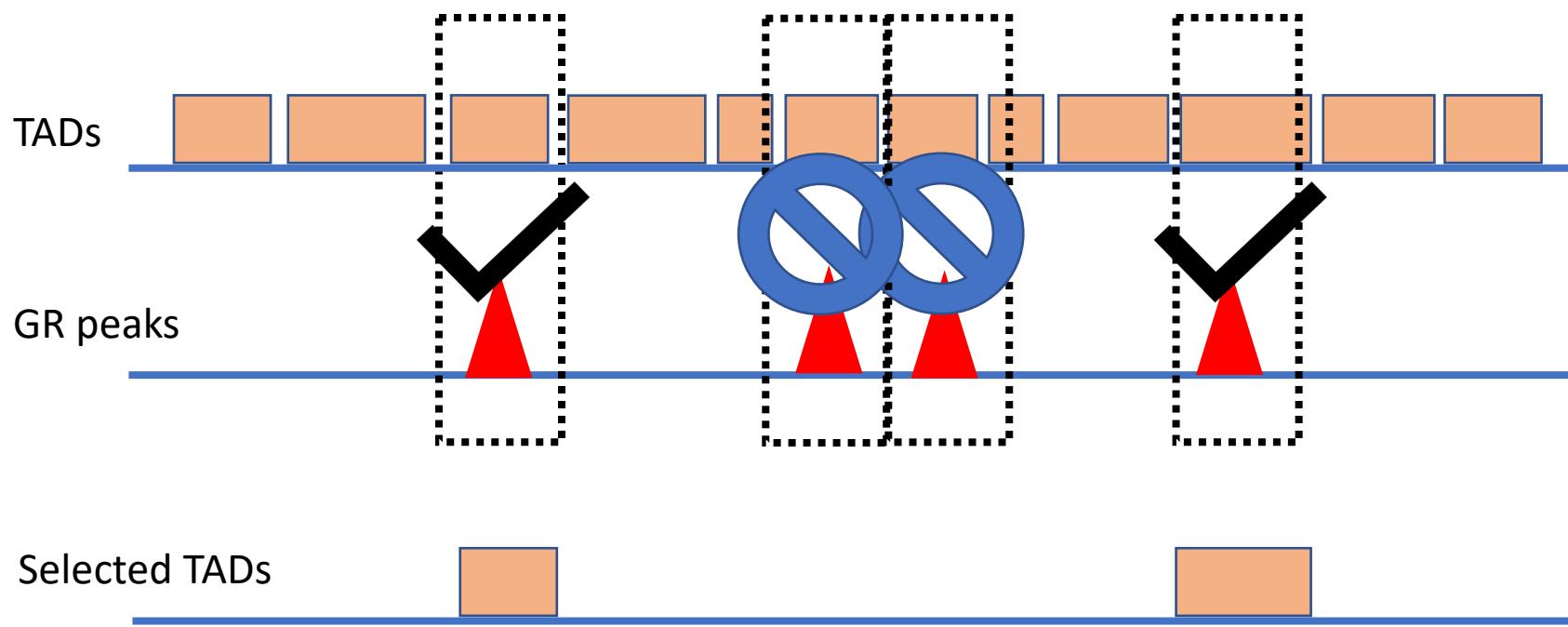


Experimental setup



Adjacency analysis

Question: Extract the TADs bearing at least one GR peak and whose two upstream and two downstream TAD neighbors have no GR peak.



Adjacency analysis

Question: Extract the TADs bearing at least one GR peak and whose two upstream and two downstream TAD neighbors have no GR peak.

```
tads = gmql.load_from_path("./tads/")

tad_tad_UP = tads.join(tads,
                       [gl.UP(), gl.DLE(5e6), gl.DG(0), gl.MD(2)],
                       output="LEFT", refName="LEFT",
                       expName="RIGHT")

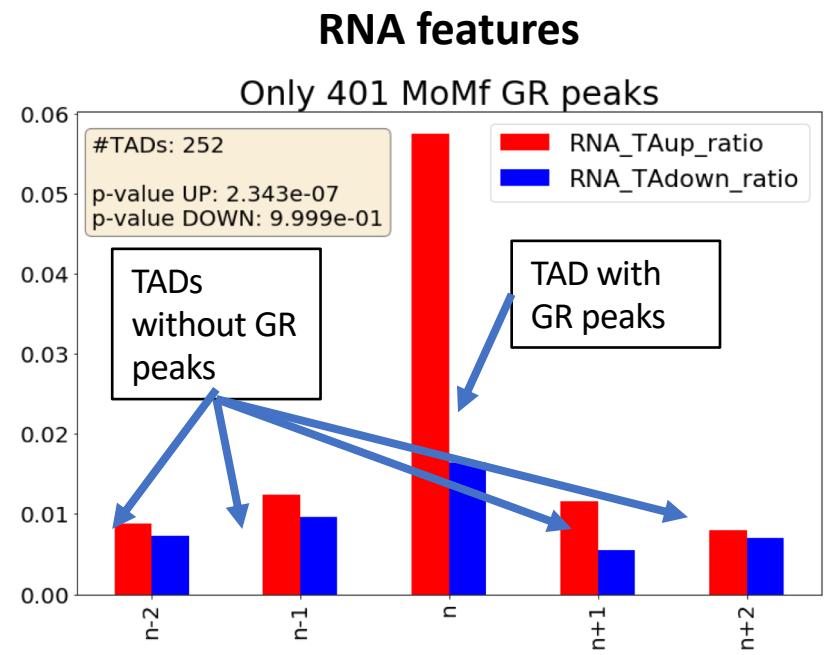
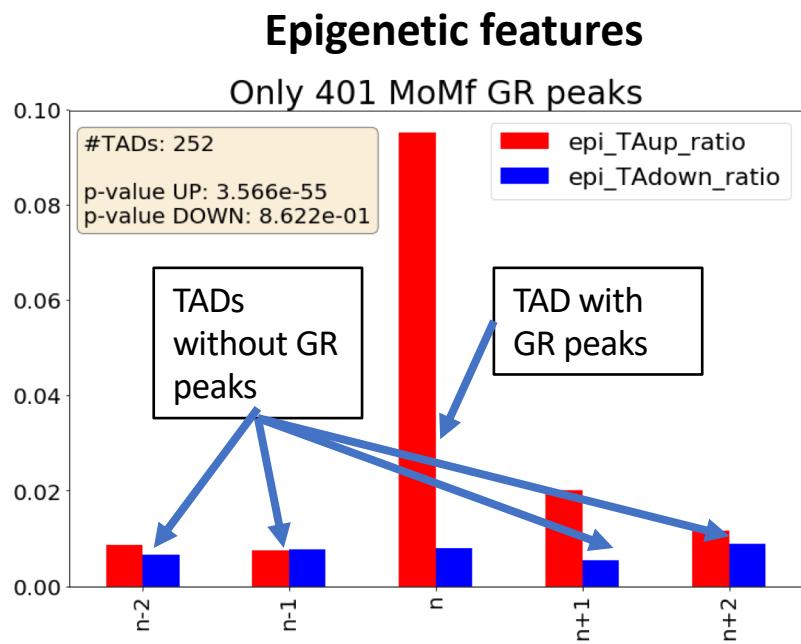
tad_tad_DOWN = tads.join(tads,
                          [gl.DOWN(), gl.DLE(5e6), gl.DG(0), gl.MD(2)],
                          output="LEFT", refName="LEFT",
                          expName="RIGHT")
```

Adjacency analysis

Question: Extract the TADs bearing at least one GR peak and whose two upstream and two downstream TAD neighbors have no GR peak.

```
tad_groups = tad_tad.group(  
    by=['LEFT.TAD_index'],  
    agg={'near_GRs': gl.SUM("RIGHT.n_GRs") })  
  
sel_tads = tad_groups.reg_select(  
    (tad_groups.LEFT_n_GRs > 0) &  
    (tad_groups.near_GRs == 0))  
  
result = sel_tads.materialize()
```

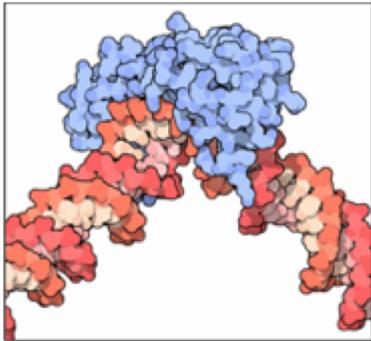
Adjacency analysis



Stefano Perna

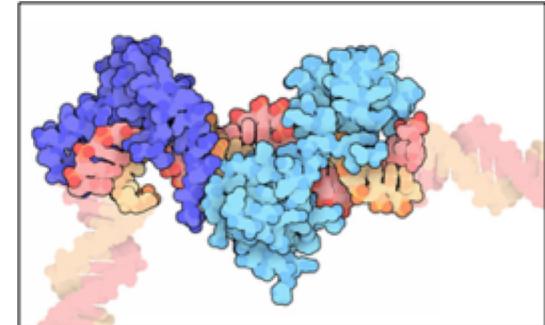
*Implementing a transcription factor
interaction prediction system using the
Genometric Query Language.*

Background



- Transcription Factors (TFs) bind and modify DNA accessibility

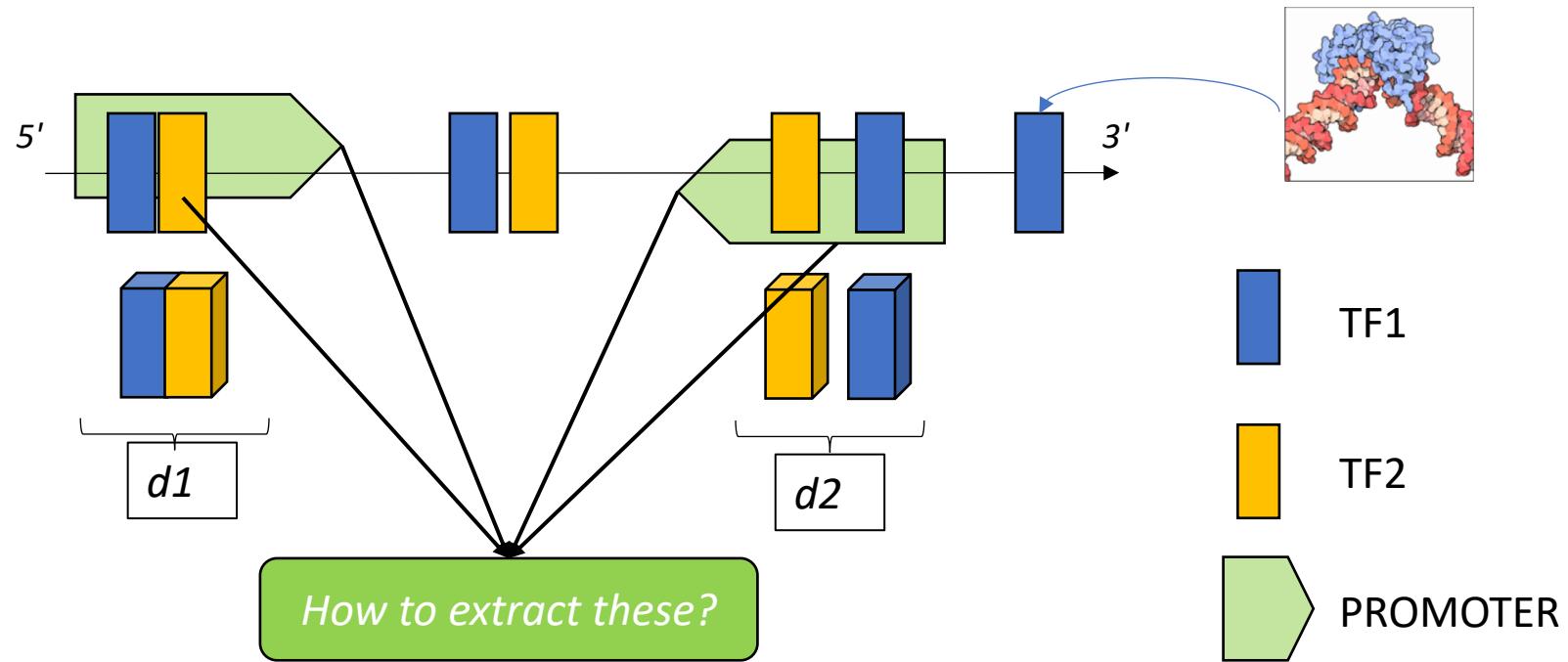
- TFs act alone or in groups
- Wet-lab investigation is combinatorial



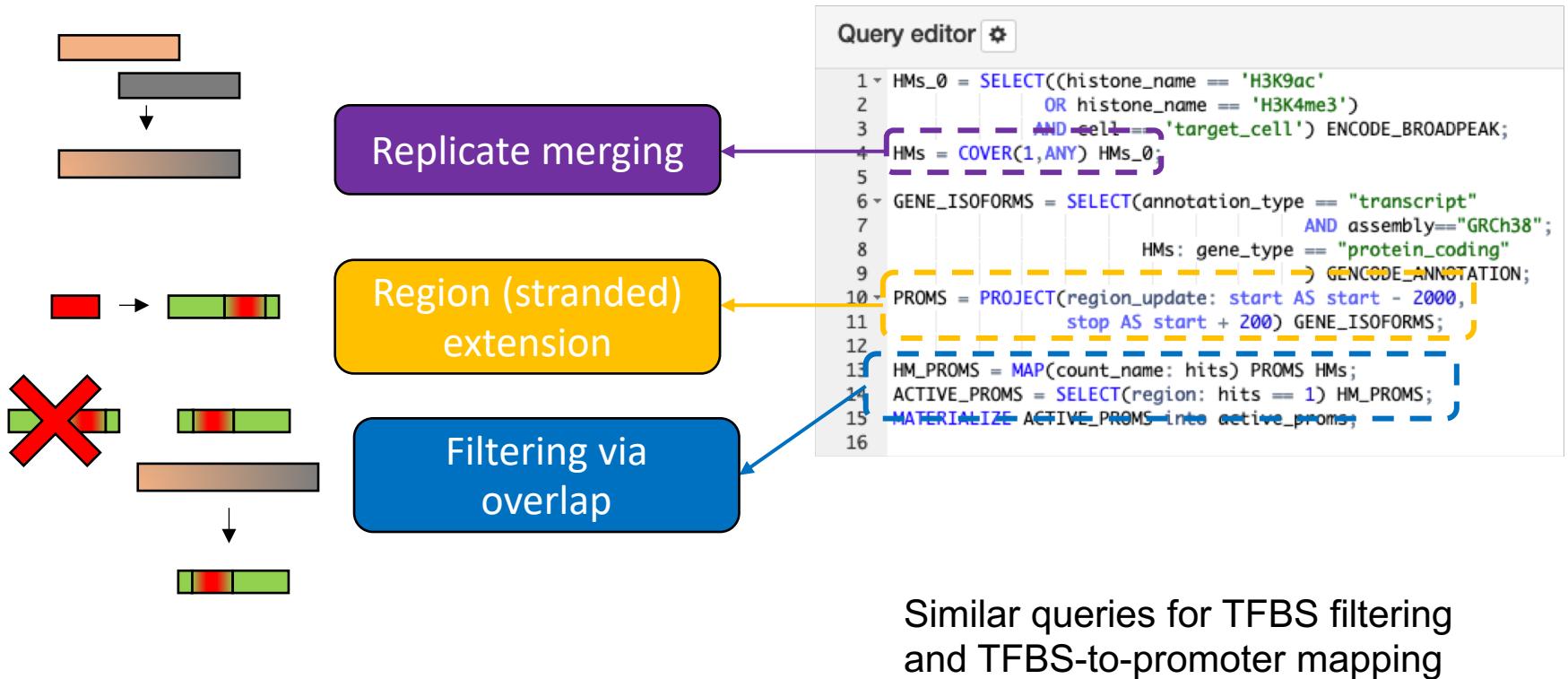
Perna S, Canakoglu A, Pinoli P, Ceri S, Wong LS. Implementing a transcription factor interaction prediction system using the Genometric Query Language. In Hiroshi Mamitsuka, editor, Data Mining for Systems Biology, chapter 6, pages 63-81, Springer Protocols, 2018.

Data & Modeling

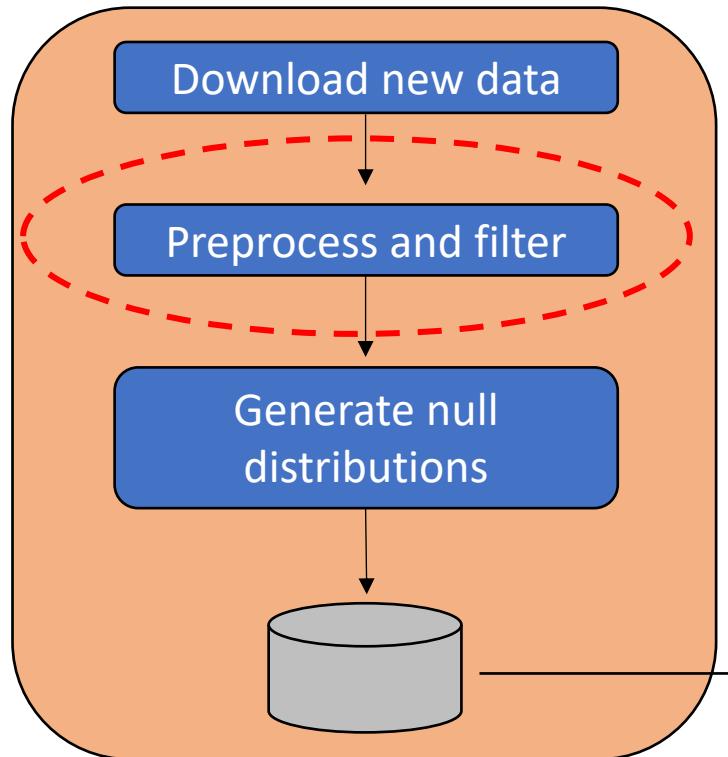
ChIP-seq experiment → Binding site information



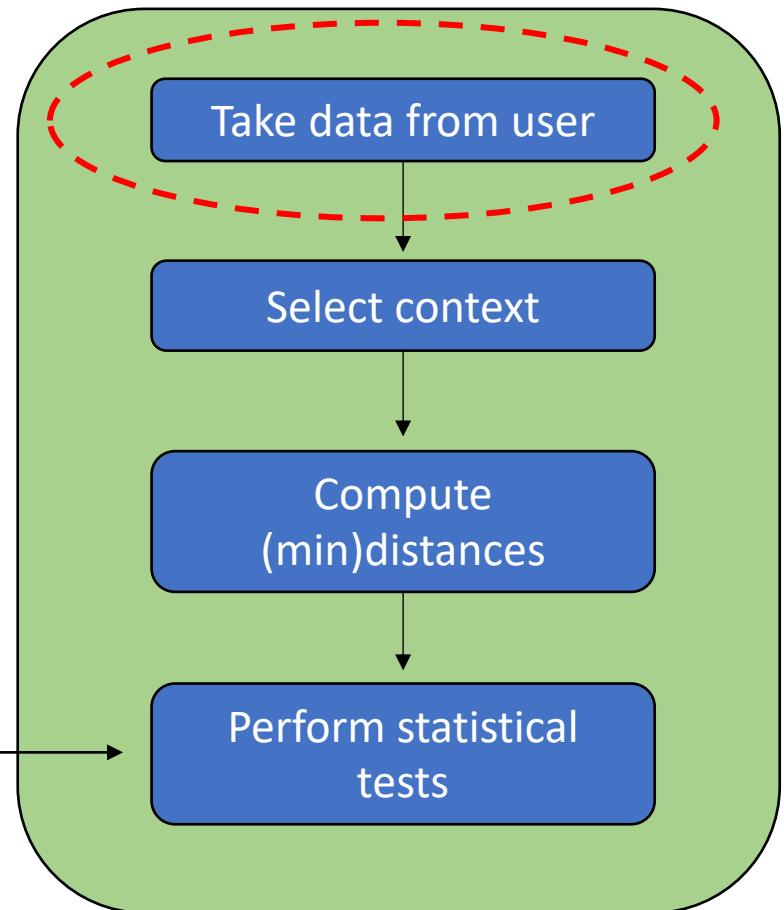
ETL and filtering with GMQL



Workflow

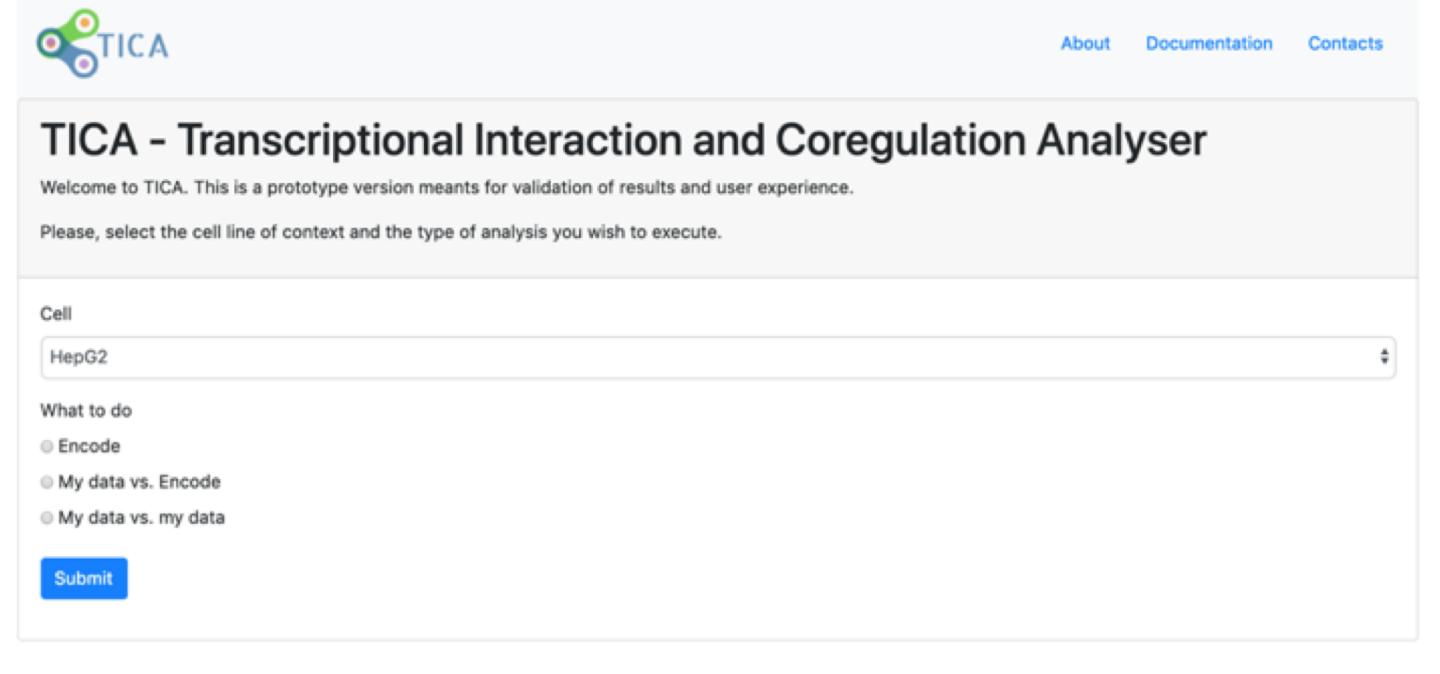


Costly, done few times



Quick, on demand

Web Service



The screenshot shows the TICA - Transcriptional Interaction and Coregulation Analyser web interface. At the top right are links for "About", "Documentation", and "Contacts". The main title "TICA - Transcriptional Interaction and Coregulation Analyser" is displayed prominently. Below it, a message says "Welcome to TICA. This is a prototype version means for validation of results and user experience." A instruction "Please, select the cell line of context and the type of analysis you wish to execute." is present. The "Cell" dropdown menu is set to "HepG2". Under "What to do", there are three radio button options: "Encode", "My data vs. Encode", and "My data vs. my data", with "Encode" being selected. A blue "Submit" button is located at the bottom left of the form area.

Web service on top of (Py)GMQL queries

Perna S, Canakoglu A, Pinoli P, Ceri S, Wong LS. Implementing a transcription factor interaction prediction system using the Genometric Query Language. In Hiroshi Mamitsuka, editor, Data Mining for Systems Biology, chapter 6, pages 63-81, Springer Protocols, 2018.

Web Service



[About](#) [Documentation](#) [Contacts](#)

TICA - Parameter Input

Please, input the parameters you wish to use in your analysis.

Select one or more TFs

- ARID3A
- ARID4B
- ARNT
- ATF1

You can select multiple TFs at a time by clicking and dragging on the list.

Select one or more TFs

- ARID3A
- ARID4B
- ARNT
- ATF1

You can select multiple TFs at a time by clicking and dragging on the list.

Maximum distance in couples [bp]

Maximum distance allowed for mindist couples, measured in bps.

How many mindistance couples are needed?

Minimum number of mindist couples required to accept a candidate.

Fraction of couples colocating in a promoter?

Minimum fraction of mindist couples which must collocate in a promoter.

Which tests do you want to use?

- Average
- Median Absolute Deviation
- Median
- Right tail size

Perna S, Canakoglu A, Pinoli P, Ceri S, Wong LS. Implementing a transcription factor interaction prediction system using the Genometric Query Language. In Hiroshi Mamitsuka, editor, Data Mining for Systems Biology, chapter 6, pages 63-81, Springer Protocols, 2018.

TICA - Results

Analysis complete. Please review results below.

Testing on *SIN3A TAF1 TBP* vs. *SIN3A TAF1 TBP* on cell line hepg2 returned the following results:

Export as CSV											
Name tf1	Name tf2	Couples	Couples Tss	Average	Average Passed	Median	Median Passed	Mad	Mad Passed	Tail 1000	Tail 1000 Passed
SIN3A	TAF1	104	0.654	111.644	Passed	44.0	Passed	35.0	Passed	0.019	Failed
SIN3A	TBP	878	0.61	105.196	Passed	47.0	Passed	39.0	Passed	0.003	Passed
TAF1	SIN3A	104	0.654	111.644	Passed	44.0	Passed	35.0	Passed	0.019	Failed
TAF1	TBP	1235	0.593	98.548	Passed	41.0	Passed	34.0	Passed	0.002	Passed
TBP	SIN3A	878	0.61	105.196	Passed	47.0	Passed	39.0	Passed	0.003	Passed
TBP	TAF1	1235	0.593	98.548	Passed	41.0	Passed	34.0	Passed	0.002	Passed

