# Supplementary notes

## Influence of the background model on the motif analysis

**Summary:**
We show here that the under-representation of the AT-rich motifs in the peak centers is not due to an inappropriate background model. We base this conclusion on several additional experiments we performed. First, using the same motif-scanning approach as in the main study, we tested several background models of various Markov orders, and local background models in sliding windows, to avoid over-fitting conditions. Furthermore, we show that the background model used recapitulates the specific composition of the peak dataset, and that the depletion is still detectable with these various background models.

As a complement, we used motif-discovery approaches: one that specifically looks for positional biases without explicit background model, and another approach to identify globally under-represented motifs. Here again, we still detect the local under-representation of AT-rich sequences.

We also performed an analyses on DNaseI hypersensitive sites derived from ENCODE data, which show the same local depletion of AT-rich sequences arguing that AT-depletion is a general feature of genomic regions involved in gene regulation.

## A – Datasets

The dataset used for this test is the ChIP-seq peaks in U2OS cells (ArrayExpress accession:E-MTAB-2731) (Fig. A1 & A2).
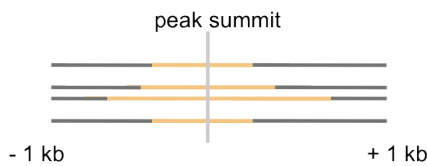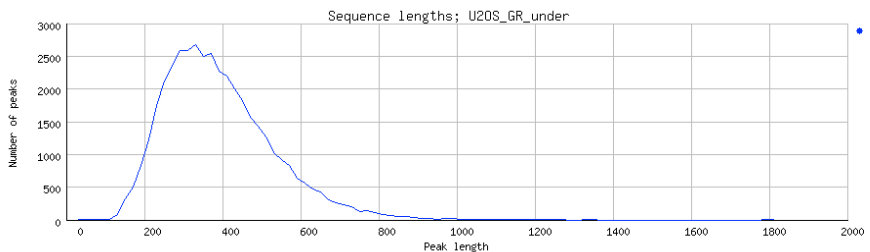


**Figure A1**: the "peak" dataset corresponds to the yellow fragments. The "2kb" dataset corresponds to the regions +/- 1kb centered on the peak summit.



Nb of peaks: 41402
Total seq. size: 16224 kb
Min length: 20 bp
Mean length: 391.864 bp
Max length: 1810 bp

**Figure A2**: Lengths of sequences in the "peak" dataset, and dataset statistics.

The composition profile of the 41402 peaks reveals a higher number of A and T nucleotides (Fig. A3). The dinucleotide profile highlights very low CpG occurences, typical of mammalian sequences [1]. For AA and TT (blue), which show the highest occurrences, we observe a local decrease at the center of the peaks. This local decrease is also detectable for TA (black) and AT (grey) even if these dinucleotides have lower occurrences overall. Note that the global nucleotide composition is roughly the same for each nucleotide (Fig. A3c).
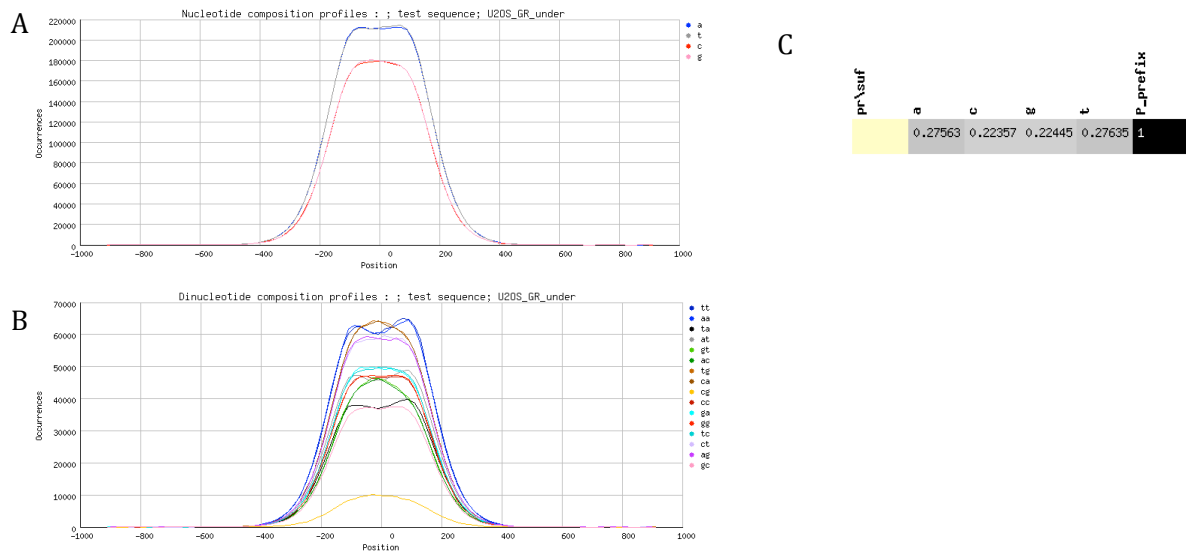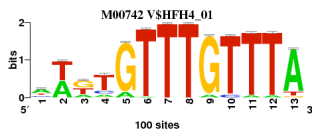
**Figure A3: A**: nucleotide composition profile. **B**: dinucleotide composition profile. **C**: Global composition of the peaks

## B – Motif scanning

The analyses of this section focus on motif scanning, with a PSSM (here Transfac M00742 – HFH4) that is under-represented in our study (Figure 1 of the main text) and is AT-rich (%GC: 0.22). We test the effect of the background model used for scanning, which most generally is a markov 0 (Bernoulli, mononucleotide) as most motif-scanning programs only supports Bernoulli models (e.g. FIMO from MEME suite). We used the RSAT *matrix-scan* program that allows higher-order markov models, and tested various order and various background datasets. Note that it is not possible to train a high-order background on a small dataset, to prevent over-fitting [2]. Increasing the order drastically increases the computing time.



### 1) Peak dataset (not repeat-masked)

Focusing on the ChIP-seq peaks with a background model calculated on these same peaks, we observe a local depletion of the HFH4 motif in the center of the peaks, even with a high-order background model, taking into account up to 4-mers frequencies (Fig. B1A,C,E). The background model correctly captures the particularities of the dataset (low CpG, high frequency of AA and TT (Fig B1B,D).
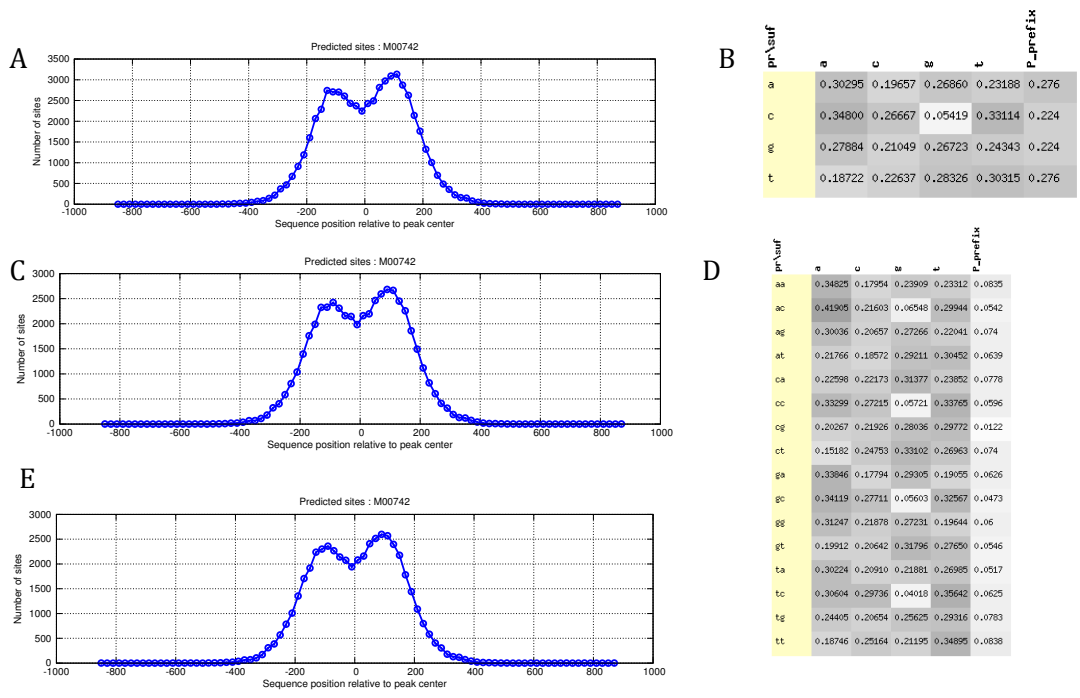
**Figure B1**: Motif profile on the peaks for matrix M00742 with background models calculated on the peaks. **A**: markov1 (dinucleotides). **B**: background model markov 1. **C**: markov2 (trinucleotides). **D**: background model markov 2. **E**. markov3 (4-mers)

### 2) 2kb dataset
#### a. peak background model
This is the model used in the main study, which shows an under-representation of several AT-rich (and few non-AT rich) motifs (Figure 1).

#### b. local window background model
To account for the local variations in nucleotide composition, we tested local background models with RSAT *matrix-scan*. A window is defined around the PSSM, in which the background model is calculated (excluding the sequence at the center, aligned with the PSSM). The window slides by 1bp, over each sequence. A very local background model can only be calculated for low order Markov models (to prevent over-fitting). We tested Markov 0 over in 75bp window (Fig.B2A) and Markov 1 in 200bp window (Fig.B2B). As these calculations are time-consuming (the background model is recalculated at each position of each sequences of the dataset, and the theoretical distribution of scores must be recalculated as well to compute the p-values), using a cluster was necessary even for just one matrix. Even with a local background, which takes the dinucleotide composition into account, we still observe the local depletion at the center. The number of hits is nevertheless lower (compare the Y axis between panels A an B with figure 1).
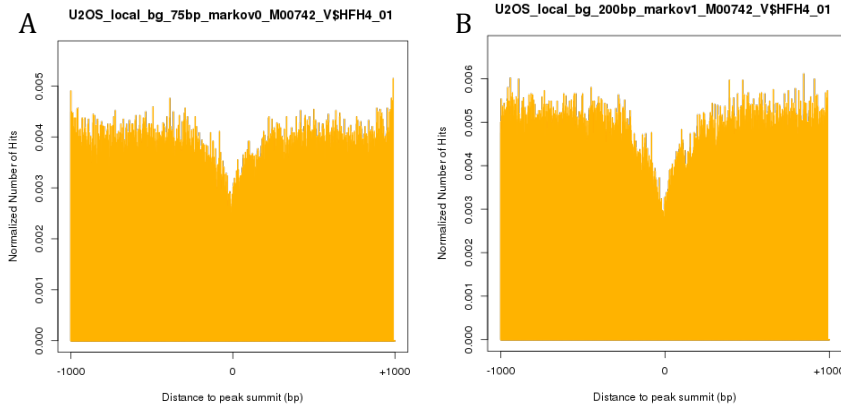
**Figure B2**: Motif profile on the 2kb dataset for matrix M00742 with local background models in sliding windows. **A**: 75bp window, Markov 0 (mononucleotides) **B**: 200bp window, Markov 1

## C – Motif discovery

The analyses of this section focus on motif discovery, without *a priori* knowledge of the searched motifs (contrary to motif scanning, there are no PSSMs).

### 1) Positional bias (no explicit background model)

The program RSAT *position-analysis* does not use an explicit background model, but finds exceptional k-mers with a bias in their position [1]. Applying this approach on the 2kb dataset reveals a local depletion of k-mers at the center of the dataset (within the peaks) (Fig.C1a). As repeats are usually excluded from the peaks due to the ChIP-seq processing (the read mapping step is usually done excluding the multi-mapping reads, hence excluding repeat regions), we also performed the same analysis on repeat-masked sequences (Fig.C1b). The results are similar, with a slight decrease in significance values. Although the most significantly under-represented 6-mers are AT-rich sequences without obvious specificity, some under-represented 6-mers are not AT-rich (Fig. C2).
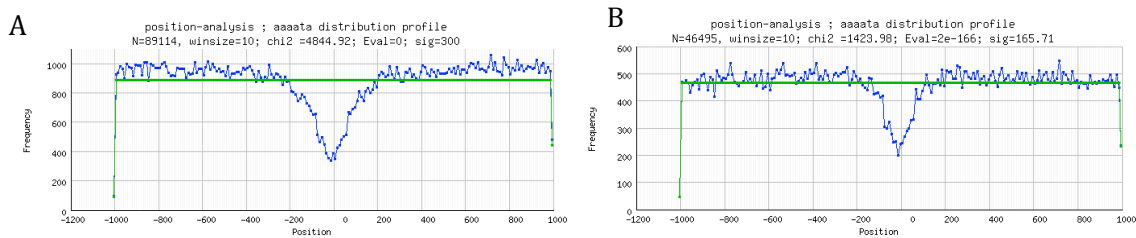


**Figure C1**: **A**: Local under-representation of the motif AAAATA in the 2kb dataset. **B**: same with repeat-masked sequences.

```
Sequence chi_value
AAAAAT  6546.3  AAAAAT|ATTTTT
AAAAAA  5296.7  AAAAAA|TTTTTT
AAAATA  4844.9  AAAATA|TATTTT
AAAATT  4504.9  AAAATT|AATTTT
TAAAAA  3880.9  TAAAAA|TTTTTA
AAATTA  3831    AAATTA|TAATTT
CAAAAA  3195.1  CAAAAA|TTTTTG
TTAAAA  3117.6  TTAAAA|TTTTAA
AATATA  2784    AATATA|TATATT
GCCTCC  2741.4  GCCTCC|GGAGGC
GGATTA  2671.8  GGATTA|TAATCC
ATATAT  2611.5  ATATAT|ATATAT
AGGCTG  2564.8  AGGCTG|CAGCCT
AAATAA  2532.9  AAATAA|TTATTT
ATTTTA  2496.8  ATTTTA|TAAAAT
AAATAT  2440    AAATAT|ATATTT
GATTAC  2398.9  GATTAC|GTAATC
TATATA  2280.1  TATATA|TATATA
```

**Figure C2**: k-mers with the highest significance value (300). The non AT-rich k-mers are indicated in bold.

## 2) Global under-represented motifs
### a. Peak-motifs on peak dataset

We added the possibility in the RSAT program *peak-motifs* to search for overall under-represented k-mers (only with the algorithm *oligo-analysis*). For 6-mers, we find various under represented motifs (Fig.C3A), some of which showing low complexity. For 7-mers, we find almost exclusively AT-rich motifs (Fig. C3B).
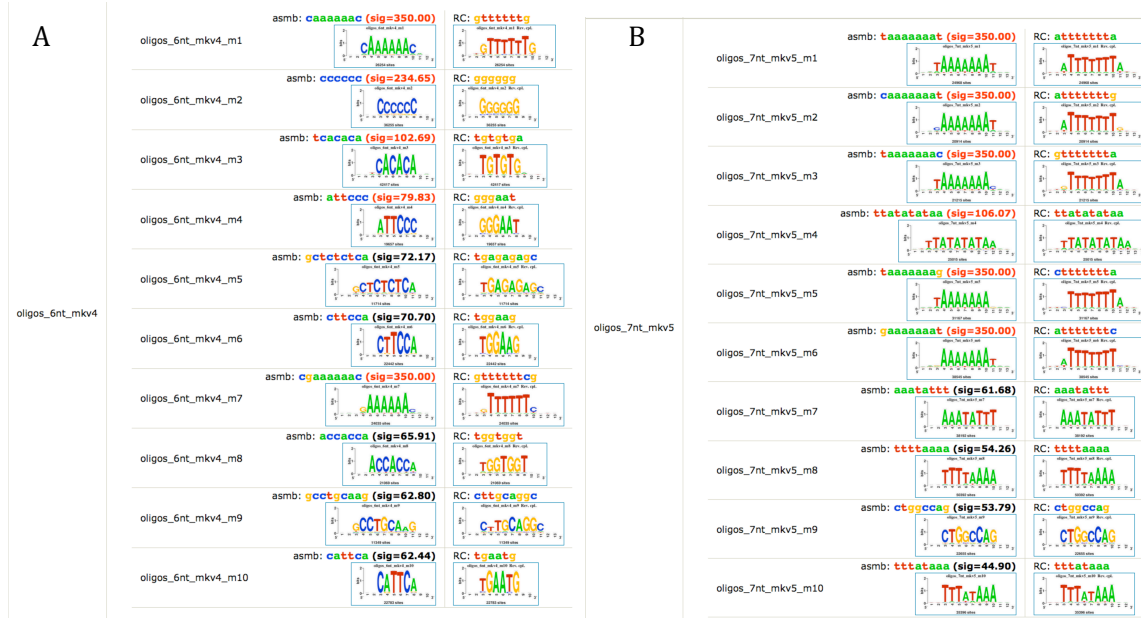


**Figure C3**: peak-motifs analysis for globally under-represented motifs **A**: 6-mers (automatic background model markov4). **B**: 7-mers (automatic background model markov5). The background models are theoretical markov models calculated on the peak dataset.

### b. Peak-motifs with two datasets: peaks and not-peaks

As control, we prepared a dataset with the same length/number of sequences as the peak dataset, but taken from genomic regions that are not located under a GR ChIP-seq peak in any of the cell lines examined (U2OS, A549, IMR90, K562, Nalm6), nor in the ENCODE blacklist. The dinucleotide profile does not feature the local decrease of AT-rich sequences (Fig.C4A). Consequently, the main difference found are these AT-rich regions in the center of the peaks (Fig.C4B).
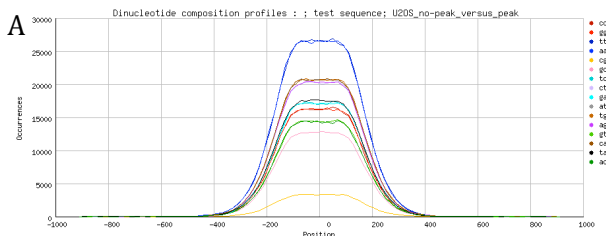


**Figure C4**: peak-motifs analysis for under-represented motifs between the peak dataset and the control dataset. The background model is an empirical background model calculated on the peak dataset, to search for under-represented motifs. **A**: dinucleotide composition profile of the control dataset.

### c. Peak-motifs on DNaseI dataset

We randomly extracted the same number of sequences from the DNaseI ENCODE dataset (wgEncodeRegDnaseClusteredV3) and performed a *peak-motifs* analysis searching for under-representation. We first observe that the dinucleotide profile showed a similar localized depletion of AA and TA as the one we saw for the GR-bound regions. The most significant under-represented motifs are AT-rich sequences, again resembling what we found for GR-bound peaks.
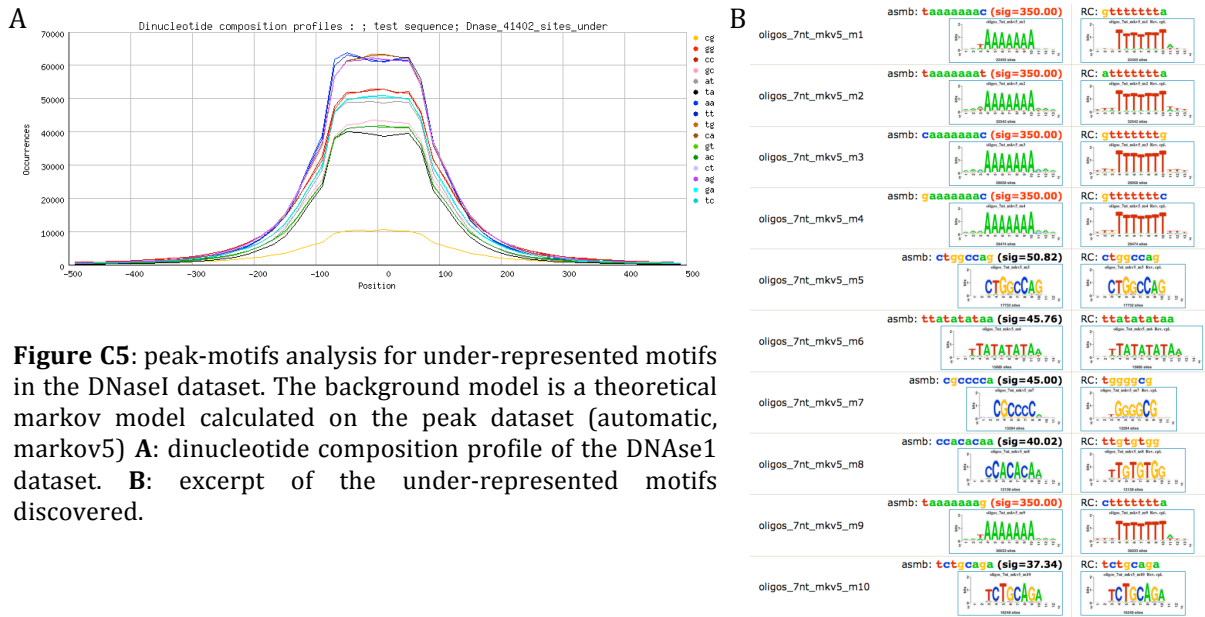
**Figure C5**: peak-motifs analysis for under-represented motifs in the DNaseI dataset. The background model is a theoretical markov model calculated on the peak dataset (automatic, markov5) **A**: dinucleotide composition profile of the DNAse1 dataset. **B**: excerpt of the under-represented motifs discovered.

### References

[1] Thomas-Chollier M, Darbo E, Herrmann C, Defrance M, Thieffry D, van Helden J (2012) A complete workflow for the analysis of full-size ChIP-seq (and similar) data sets using peak-motifs. *Nature Protocols* **7**, 1551–1568.
[2] Turatsinze J-V, Thomas-Chollier M, Defrance M, van Helden J (2008) Using RSAT to scan genome sequences for transcription factor binding sites and cis-regulatory modules. *Nature Protocols* **3**, 1578–1588.