**Supplementary Note 1: Derivation motivating the use of the outer sum as the background matrix.**

Here we describe the rationale for the strategy we used to combine enhancer-promoter 1D average signal into a background matrix.

We have noticed that different conditions in Micro-C libraries have very different densities of Micro-C paired-end tags mapping to each anchor, which may be caused by the effect of a perturbation on nucleosome occupancy and/ or other aspects of 1D chromatin structure. Our primary goal in the background-corrected APA is to normalize for these large changes in the distribution of 1D Micro-C signal between different conditions.

Our formulation assumes that we compute an aggregate background matrix to represent the background signal, which is an essential step when dealing with the sparsity of Micro-C data at higher resolutions. The background matrix helps to account for variations in signal density between different conditions in Micro-C libraries, as mentioned above. By incorporating the background matrix in our analysis, we can normalize the data and better understand the relationship between enhancer-promoter interactions and the underlying genomic architecture.

We assume that sequencing depth normalized signal reflects the probability of observing a read at a fixed position of the genome. We note that Micro-C signal is not a robust estimate of the probability of observing contacts, but we think this is a useful and simple way to think about the problem that illustrates the approach that we have proposed.

We can compute the probability of observing a 1D signal (i.e., the total Micro-C signal mapping to each position of the genome, regardless of where the mate pair maps) at any of the $n_p$ promoters as the sum of the probabilities over the set of promoters ($p1 .. p_{np}$):

$$p1 + p2 + ... p_{np}$$

And likewise the sum of observing the probabilities of any of the $n_e$ enhancers as:

$$e1 + e2 + ... e_{ne}$$

To achieve our goal (i.e., normalizing away changes in 1D Micro-C signal), we compute the probability of observing one end of a pair at any of the enhancers or at any of the promoters. We consider the signal between an enhancer or promoter and any other position in the genome, whether that read falls into the region defined by the APA for the partner promoter or enhancer or not.

In this case, for each cell *i, j* in the background matrix, the probability of observing a read at position *i, j* is the sum of the probabilities of observing a read at that position (*i*) relative to the enhancer or the probability of observing a read at position (*j*) relative to the promoter. To model this, we sum the probabilities:

$$p1\{j\} + p2\{j\} + ... p_{np}\{j\} + e1\{i\} + e2\{i\} + ... e_{ne}\{i\}$$

Or:

$$\sum_{y=1}^{n_p}\left(p_j^y\right) + \sum_{x=1}^{n_e}(e_i^x)$$

This reduces to the outer sum, which we have used in our manuscript.

A commonly used alternative in some fields is to use the outer product. We therefore considered the alternative assumptions that we would have to make to justify the use of the outer product in the same basic setup. Rather than assuming that a read can fall in position *i* or position *j* of the enhancer or promoter, we calculate the probability that a read falls into position *i* of the enhancer *and* falls into position *j* of the promoter based on the 1D distribution of reads. In this case, the probability that a read falls into position *i* (at the enhancer) and position *j* (promoter) is the product of the probabilities across every pair of enhancer-promoter candidates. The computation would be:

p1{i} * e1{j} + p1{i} * e2 {j} ... p1{i} * e$_{ne}$ {j} + p2 {i} * e1{j} + …
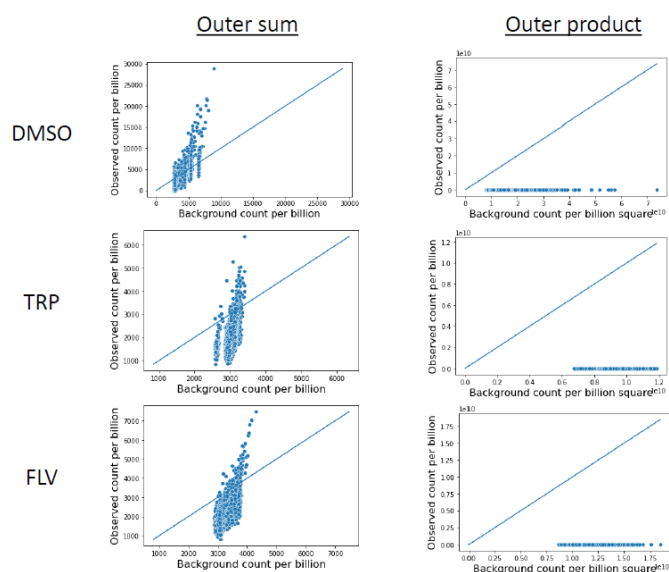
Or:

$$\sum_{x=1}^{n_e}\sum_{y=1}^{n_p}\left(p_j^y\right) * (e_i^x)$$
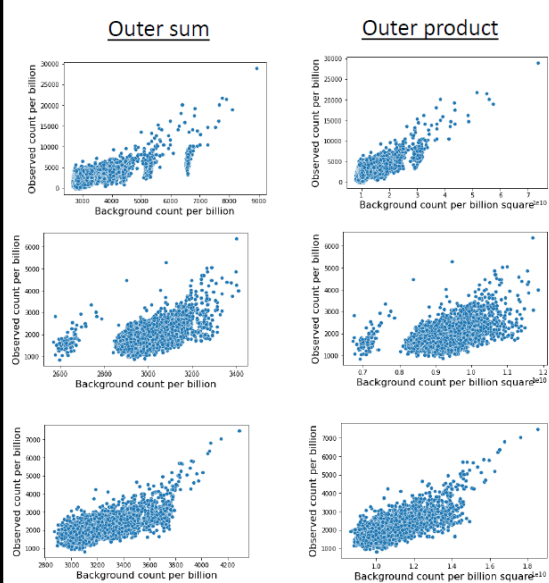
This is the outer product.

Using the outer product in this context has several limitations. One major concern is that this null model assumes that any of the enhancers can interact with any promoter with equal probability, as it involves summing over the products of all pairs. However, this assumption is not consistent with Hi-C/Micro-C data, as the genome's topology restricts or significantly reduces interactions between enhancer-promoter (E-P) pairs on different chromosomes or those far apart on the same chromosome. Real interactions are also influenced by known biological structures, such as topological domains, 1D distance decay, and other constraints. Moreover, technical factors, including the need to exclude enhancer-promoter pairs that would overlap at the resolution of the APA, prevented us from including some candidate pairs in the observed APA. These biological and technical factors would lead the background model to correct for potential interactions that are excluded from our observed matrix for technical or biological reasons if we used the outer product. Based on these considerations, normalizing the data for interactions between candidate enhancer and promoter pairs that could never interact in practice seems inappropriate.

Consistent with the intuition noted above, we observed that the outer product formulation was on a very different scale as the observed data. Conversely, the scale of the outer sum was much more similar to the observed data. Both of these observations are shown in the scatterplots below (see left):

## Compared to x=y

## General trend

### Outer sum    ### Outer product

### Outer sum    ### Outer product

DMSO

TRP

FLV

Taken together, these considerations provide some empirical motivation in support of the outer sum formulation.