# nature

## Accelerated Article Preview

# Cell type directed design of synthetic enhancers

Ibrahim I. Taskiran, Katina I. Spanier, Hannah Dickmänken, Niklas Kempynck, Alexandra Pančíková, Eren Can Ekşi, Gert Hulselmans, Joy N. Ismail, Koen Theunis, Roel Vandepoel, Valerie Christiaens, David Mauduit & Stein Aerts

# Cell type directed design of synthetic enhancers

Ibrahim I. Taskiran[1,2,3], Katina I. Spanier[1,2,3], Hannah Dickmänken[1,2,3], Niklas Kempynck[1,2,3], Alexandra Pančíková[1,2,3,4], Eren Can Ekşi[1,2,3], Gert Hulselmans[1,2,3], Joy N. Ismail[1,3,#], Koen Theunis[1,2,3], Roel Vandepoel[1,2,3], Valerie Christiaens[1,2,3], David Mauduit[1,2,3], and Stein Aerts[1,2,3,*]

1 Laboratory of Computational Biology, VIB Center for AI & Computational Biology (VIB.AI), Leuven, Belgium.
2 VIB-KULeuven Center for Brain & Disease Research, Leuven, Belgium.
3 Department of Human Genetics, KU Leuven, Leuven, Belgium.
4 VIB-KULeuven Center for Cancer Biology, Leuven, Belgium.
# Current address: UK Dementia Research Institute at Imperial College London, London, UK

* Correspondence to stein.aerts@kuleuven.be

## Summary

Transcriptional enhancers act as docking stations for combinations of transcription factors (TFs) and thereby regulate spatiotemporal activation of their target genes. It has been a long-standing goal in the field to decode the regulatory logic of an enhancer and to understand the details of how spatiotemporal gene expression is encoded in an enhancer sequence. Here, we show that deep learning models can be used to efficiently design synthetic, cell type specific enhancers, starting from random sequences, and that this optimization process allows for a detailed tracing of enhancer features at single-nucleotide resolution. We evaluate the function of fully synthetic enhancers to specifically target Kenyon cells or glial cells in the fruit fly brain using transgenic animals. We further exploit enhancer design to create "dual-code" enhancers that target two cell types, and minimal enhancers smaller than 50 base pairs that are fully functional. By examining the state space searches towards local optima, we characterise enhancer codes through the strength, combination, and arrangement of TF activator and TF repressor motifs. Finally, we apply the same strategies to successfully design human enhancers, which adhere to similar enhancer rules as *Drosophila* enhancers. Enhancer design guided by deep learning leads to better understanding of how enhancers work and shows that their code can be exploited to manipulate cell states.

## Main

Cell type specific expression of a target gene is achieved when a unique combination of TFs activates a specific enhancer; while this enhancer remains either passively ("default-off"[1,2]) or actively repressed in other cell types (e.g., via repressor binding[3] or co-repressor/polycomb recruitment). Typically, when an enhancer is translocated to another chromosome or to an episomal plasmid, it maintains cell type specific control of its nearby reporter gene[4,5]. Therefore, its regulatory capacity is contained within the enhancer DNA sequence and has co-evolved to respond uniquely to a specific trans-environment in a cell type. A thorough understanding of how enhancer activation is encoded in its DNA sequence is important, as it is a key component for the modelling and prediction of gene expression[6,7]; for the interpretation of non-coding genome variation[8,9]; for the improvement of gene therapy; and for the reconstruction and manipulation of dynamic gene regulatory networks underlying developmental, homeostatic, and disease-related cell states.

44   Many complementary approaches and techniques have been used to decode enhancer logic[4]. These
45   include studies of individual enhancers by mutational analysis[10–12], in vitro TF binding (e.g.,
46   electrophoresis mobility shift assay), cross-species conservation[13], and reporter assays. The upscaling
47   of such studies led to the identification of common features of co-regulated enhancers [14–16]. These
48   experimental findings also triggered the improvement of computational methods for the prediction
49   of *cis*-regulatory modules, whereby feature selection and parameter optimization led to new insights
50   into how binding sites cluster and how their strength (or binding energy) impacts enhancer
51   function[11,12,17–20]. Wider adoption of genome-wide profiling of chromatin accessibility[21], single-cell
52   chromatin accessibility[22–24], histone modifications[25,26], TF binding[27], and enhancer activity[15,28] led to
53   significantly larger training sets of co-regulated enhancers that could then be used for *a posteriori*
54   discoveries of TF motifs and enhancer rules, aided by the growing resources of high-quality TF
55   motifs[29,30]. Additional mechanistic insight has been provided by thermodynamic modelling of
56   enhancers[31,32], in vivo imaging of enhancer activity[33], the analysis of genetic variation through eQTL
57   and caQTL analysis[2,34], and high-throughput in vitro binding assays[35,36]. Recently, the enhancer biology
58   field embraced the use of convolutional neural networks (CNN) and network-explainability techniques
59   that again provided a significant leap forward in terms of prediction accuracy and syntax
60   formulation[6,37–44].
61   An orthogonal strategy to decode enhancer logic is to engineer synthetic enhancers from scratch. This
62   approach has the advantage that the designer knows exactly which features are implanted, so that
63   the minimal requirements for enhancer function can be revealed. Recent work showed the promise
64   of CNN-driven enhancer design by successfully designing yeast promoters[45], and by using a CNN to
65   select high-scoring enhancers for S2 cells, from a large pool of random sequences[38]. Here we tackle
66   the next challenge in enhancer design, namely to design enhancers that are cell type specific. To this
67   end, we used previously trained deep learning models for which we have already validated the
68   accuracy of nucleotide-level interpretation and motif-level predictions[8,39] (Supplementary Note 1).
69   Using these enhancer models as a guide (or 'oracle'), we tested three different sequence design
70   approaches[46,47] (Fig. 1).
71
72   **In silico evolution**
73   As a first strategy for enhancer design, we created synthetic enhancers to specifically target Kenyon
74   cells (KC) in the mushroom body of the fruit fly brain, using a nucleotide-by-nucleotide sequence
75   evolution approach[45] (Methods). This approach starts from a 500 bp random sequence that is evolved
76   from scratch (EFS) in silico towards a chosen cell type through multiple iterations. Prediction scores
77   are calculated using DeepFlyBrain[39], a deep learning model trained on differentially accessible regions
78   across multiple cell-types of the *Drosophila* brain and that can recognize motif-level nucleotide
79   arrangements for many cell-types (Supplementary Note 1). At each iteration we performed saturation
80   mutagenesis[9,44,48] whereby all nucleotides were mutated one by one, and each sequence variation was
81   scored by DeepFlyBrain to select the mutation with the greatest positive delta score for the KC class
82   (among 81 classes representing different cell types that the model learned to predict). We performed
83   this procedure starting from 6,000 GC-adjusted random sequences and observed that after 15
84   iterations, DeepFlyBrain KC prediction scores increased from around the minimal score (0) to nearly
85   the maximum score (1), while remaining low for other cell types (Fig. 2a, Extended Data Fig. 1a,b). We
86   found this greedy search to provide a good balance between computational cost and ability to
87   efficiently yield high-scoring sequences, compared to alternative state space searches (Extended Data
88   Fig. 2a-d, Methods).

89   Next, we investigated the initial (random) sequence and the specific paths that are followed through
90   the search space towards local optima. For only a small fraction (3%) of random sequences the
91   prediction score remained below 0.5 even after 15 mutations (Extended Data Fig. 1c). These
92   sequences were mostly characterized by more instances of repressor binding sites together with an
93   increased number of mutations required to generate sufficient activator binding sites. A second
94   observation is that even though 500 bp space is given to the model, the selected mutations
95   accumulated in about 200 bp space, preferentially at the center of the random sequence (Extended
96   Data Fig. 1d,e).
97   We investigated the consequences of each mutation on shaping the enhancer code using
98   DeepExplainer-based contribution scores (Fig. 2b, Methods). This revealed that initial random
99   sequences harbor several short repressor binding sites by chance and these are preferentially
100  destroyed during the first iterations (Extended Data Fig. 1f,g). These repressor sites contribute
101  negatively to the KC class prediction and represent candidate binding sites for KC specific repressor
102  TFs such as Mamo and CAATTA[39]. The nucleotides with the highest impact represent mutations that
103  destroy a repressor binding site and simultaneously generate a binding site the key activators
104  Eyeless (Ey), Mef2 or Onecut. Eventually, DeepExplainer highlighted multiple candidate activator
105  binding sites, whereby Ey, Mef2, and Onecut sites dominate (Fig. 2b and Extended Data Fig. 1f,g).
106  To test whether the in silico evolved enhancers can drive reporter gene expression in vivo, we
107  randomly selected 13 sequences after 10 or 15 iterations (Fig. 2c and Supplementary Fig. 1, 2) and
108  integrated them into the fly genome with a minimal promoter and a GFP reporter gene (Methods).
109  Investigating the GFP expression pattern by confocal imaging showed that 10 out of these 13 tested
110  synthetic enhancers were active specifically in the targeted cell-type, the Kenyon cells (Fig. 2d and
111  Extended Data Fig. 1h). Some enhancers did not show activity after 10 mutations but became active
112  after an additional five mutations (Fig. 2d, Extended Data Fig. 1i,j and Supplementary Fig. 3). The three
113  enhancers without GFP signal in KC were found to also be Dachshund negative, indicating the potential
114  loss of KC (Extended Data Fig. 1k). Using assay for transposase accessible chromatin by sequencing
115  (ATAC-seq) on the brains of the transgenic lines, we verified that the synthetic enhancers become
116  accessible when integrated into the genome (Extended Data Fig. 1l), as predicted by the model.
117  We also generated transgenic lines to test enhancers at different steps during the evolutionary design
118  process (Supplementary Fig. 4, 5). We found that random sequences, or sequences with only few
119  mutations remain inactive, while enhancer activity is initiated when repressor sites are removed and
120  Ey and Mef2 sites are generated; and activity further increases with more and stronger instances of
121  activator motifs (Extended Data Fig. 1m,n).
122  To demonstrate that enhancers can be generated for other cell types, we started from the same
123  random sequences as above and evolved them into perineurial glia (PNG) enhancers (Extended Data
124  Fig. 2e). After 15 mutations, putative PNG repressor sites have been destroyed and activator sites have
125  been generated (Fig. 2e and Supplementary Fig. 6). We validated six designed sequences by creating
126  transgenic GFP reporter flies, and confirmed that four were positive, as they drive GFP specifically in
127  perineurial glial cells (Fig. 2f and Extended Data Fig. 2f). Because the same random sequence was
128  evolved into either KC or PNG enhancers, this experiment underscores that the chosen mutations, and
129  the candidate binding sites they destroy or generate, causally underlie the activity of these synthetic
130  enhancers.
131  Given that KC enhancers can arise from random sequences after 10 or 15 mutations, we hypothesized
132  that certain genomic regions may require even fewer mutations to acquire KC enhancer activity. We
133  scanned the entire fly genome and identified regions with high prediction scores but without

134  chromatin accessibility in KC (Extended Data Fig. 2g,h, Methods). By applying sequence evolution to
135  these sequences, three out of four sequences became positive KC enhancers with only six mutations
136  (Fig. 2g,h, Extended Data Fig. 2i,j and Supplementary Fig. 7). When the negative enhancer was further
137  evolved, with an additional five mutations, it also became positive (Fig. 2g and Extended Data Fig. 2i,j).
138  This suggests that KC enhancers, and likely other cell type enhancers as well, can arise de novo in the
139  genome with few mutations.
140  To summarize the changes that happened during the design process, we performed motif discovery
141  across all 6,000 sequences, at each step of the optimization path (Extended Data Fig. 1f,g). This
142  confirmed that repressor sites are often present in random sequences and that they are preferentially
143  destroyed during the first steps of the search algorithm. To experimentally test that these short
144  repressor sites functionally cause repression, we selected three positive synthetic enhancers and
145  three of the near-enhancers rescued from the genome and evolved these to become non-functional
146  by manually choosing the mutations that decrease the prediction score by creating repressor binding
147  sites (Extended Data Fig. 2i and Supplementary Fig. 8, 9). We avoided mutating any of the predicted
148  activator sites (Fig. 3a); thus, placed repressor motifs in between activator sites. New transgenic lines
149  with these sequences integrated into the genome confirm that all tested enhancers have entirely lost
150  their activity (Fig. 3b). This shows that a sufficient number of repressor sites can dominate over a
151  functional combination of activator sites.
152  The sequence evolution strategy thus represents an intuitive and efficient approach to generate cell
153  type specific enhancers and to characterize their functional constituents.
154

155  **Multiple cell type codes**
156  A single enhancer can be active in multiple, different cell types[49], and our earlier work suggested that
157  this can be achieved by enhancers that contain multiple codes for different cell types, intertwined
158  within a single ~500 bp sequence[39]. Based on this finding, we wondered whether a genomic enhancer
159  that is active in a single cell type, could be synthetically augmented to become also active in a second
160  cell type. To test this, we started with two optic lobe enhancers (*amon* and *CG15117*) that are
161  accessible and active in T4/T5 and T1 neurons respectively[39] and whose activity per cell type is also
162  predicted correctly by DeepFlyBrain (Fig. 3c-e, Extended Data Fig. 3a-c). We then performed in silico
163  evolution on these enhancers towards KC, while simultaneously maintaining a high prediction score
164  for the original cell type. After 13 and 14 mutations, the enhancers were also predicted as KC
165  enhancers, but retained T4 and T1 binding sites. Testing the augmented sequences in vivo with a GFP
166  reporter confirmed the spatial expansion of the enhancer activity to KC (Fig. 3f-g, Extended Data Fig.
167  3c-f, Supplementary Fig. 10, Methods).
168  Reciprocally, enhancers active in multiple cell types may be pruned towards a single cell-type code.
169  We searched for genomic enhancers that score high for multiple cell types (Fig. 3h-l). We selected a
170  *Pkc53e* enhancer that is accessible and active in both optic lobe T neurons and KCs and predicted
171  correctly by the model. This time, we drove the in silico evolution to maintain the KC prediction score,
172  while decreasing the T neurons prediction score (Methods). After nine mutations, the sequence was
173  predicted to have only KC activity (Fig. 3m). Nucleotide contribution scores show that the most
174  important binding sites for KCs were unaffected after nine mutations while the activator binding sites
175  were destroyed and new repressor binding sites were created for T neurons (Extended Data Fig. 3g).
176  Testing the final sequence in vivo confirmed the spatial restriction of the enhancer activity (Fig. 3n).
177  Together, our results suggest that, guided by the DeepFlyBrain model, intertwined enhancer codes
178  can be independently dissected and altered.

179

**Motif implantation**

180 As a second strategy, we used a classical motif implantation approach to design KC enhancers. The

181 rationale behind this strategy is based on our results above: nucleotide-by-nucleotide sequence

182 evolution showed that all the selected mutations were associated with the creation or destruction of

183 a TF binding site, rather than affecting contextual sequence between motif instances (Fig. 2b,e,h,

184 Extended Data Fig. 3d,e,g). This suggested that a combination of appropriately positioned activator

185 motifs, without the presence of repressor motifs, would be sufficient to create a cell type-specific

186 enhancer. Furthermore, we reasoned that by applying this design strategy to thousands of random

187 sequences we could gain additional insight into the KC enhancer logic. To this end, we iteratively

188 implanted strong TF binding site instances in 2,000 random sequences, selecting locations with the

189 highest prediction score towards the KC class. We first implanted a single binding site for one of the

190 four key activators of KC enhancers, namely Ey, Mef2, Onecut, and Sr[39] and then specific combinations

191 of sites in a particular implantation order (Extended Data Fig. 4a, Methods). This revealed that Ey and

192 Mef2 had the strongest effect on the prediction score, while Onecut and Sr increased the prediction

193 score only marginally (Fig. 4a). Implanting Ey and Mef2 consecutively increased the score more than

194 the sum of their individual contribution and their implantation order did not affect the final score.

195 Adding Onecut and then Sr on top of Ey and Mef2 sites increased the scores even further until it

196 reached the level that we obtained above after 15 mutations through in silico sequence evolution (Fig.

197 4a). We could also observe some minor preferences in the motif flanking sequence (e.g. Mef2 is

198 flanked by T or G in 5' and A or C in 3'; Extended Data Fig. 4a)

199 We also found that high-scoring configurations consisted of activator sites that are positioned close

200 together within a distance usually smaller than 100 bp (Fig. 4b,c, Extended Data Fig. 4b). When the Ey

201 and Mef2 pair were implanted on the same strand, we observed strong preference for a 5 bp distance

202 (or 4 bp when implanted on opposite strands) between the two binding sites whereby Mef2 was

203 located upstream of Ey (Fig. 4b, Extended Data Fig. 4c). For the Ey and Onecut pair, there was a strong

204 preference for a 3 bp space and Onecut preferred the downstream side of Ey (Fig. 4c, Extended Data

205 Fig. 4d).

206 We investigated the nucleotide contribution scores before and after motif implantations for an

207 example sequence with high prediction score where motifs were inserted close together (Fig. 4d,e,

208 Supplementary Fig. 11). The initial random sequence contained multiple repressor binding sites and

209 the Ey binding site implantation destroyed the strongest repressor binding site. Mef2 and Onecut

210 implantations followed the predicted spacing relative to Ey, with a distance of 5 bp and 3 bp,

211 respectively. This can explain why implantation of motifs at random locations yields lower scoring

212 sequences (Fig. 4a). Even though some repressor binding sites were still present at further distances,

213 their relative negative contribution was decreased after the activator binding site implantations (Fig.

214 4e). Testing this designed 500 bp sequence in vivo confirmed specific activity in KC (Fig. 4f).

215 Introduction of mutations to generate repressor sites close to the implanted motifs (none of the

216 activator sites was modified) resulted in complete loss of enhancer activity in vivo, suggesting

217 dominance of repressor motifs (Fig. 4d,e,g). Furthermore, a 49 bp subsequence, containing just the

218 three binding sites, resulted in the same activity and specificity in vivo (Fig. 4h,i, Supplementary Fig.

219 12). We further confirmed the robustness of the motif implanting design by validating in vivo a second

220 500 bp sequence displaying increased spacing between motifs (Extended Data Fig. 4e,f,g). This result

221 suggests that a functional KC enhancer can be created via motif-by-motif implantation with just these

223 three binding sites and its size can be decreased to the minimal length required to contain these
224 binding sites.

225 As a third strategy for enhancer design, we used Generative Adversarial Networks (GAN) that have
226 been shown to be powerful generators in different fields[43,48], including the generation of functional
227 genomic sequences[46]. This method was less interpretable than in silico evolution or motif implanting
228 but still allowed for the generation of functional and specific enhancers (Supplementary Note 2).

229

## Human enhancer design

231 We used our previously trained and validated melanoma deep learning model, DeepMEL2[8]
232 (Supplementary Note 1) with the same three strategies as before, to design human melanocyte, or
233 melanocyte-like melanoma (MEL) enhancers. Like the *Drosophila* experiments, we started from GC-
234 adjusted random sequences (Extended Data Fig. 5a) and, by following the nucleotide-by-nucleotide
235 sequence evolution approach, we evolved them into sequences with high prediction scores for the
236 MEL class. This process drove the generation of activator binding sites (SOX10, MITF, TFAP2) and the
237 destruction of ZEB motifs to resemble MEL genomic enhancers; the prediction scores started to
238 plateau after 15 mutations (Fig. 5a, Extended Data Fig. 5b,c). We randomly selected 10 regions that
239 were evolved from scratch (EFS-1-10) with 15 mutations and tested their activity with a luciferase
240 assay in vitro, in a MEL cell line (MM001) (Fig. 5b,c and Methods). Seven out of 10 tested enhancers
241 showed activity in the range of previously characterized positive control (native) enhancers and none
242 of them showed activity in a cell line that represents another melanoma cell state (mesenchymal-like,
243 MM047) where the MEL-specific TFs (SOX10, MITF, and TFAP2) are not expressed (Fig. 5d, Extended
244 Data Fig. 5d). When we integrated these synthetic enhancers into the genome of the MM001 cell line
245 using lentiviral vectors (Methods), they generated an ATAC-seq peak, while neither the random
246 sequences nor the evolved sequence when integrated in a non-MEL cell line are accessible (Fig. 5e,
247 Extended Data Fig. 5e,f).

248 Next, we tested the activity of a series of synthetic sequences, along the design path, from a random
249 sequence to an active enhancer (Extended Data Fig. 6, Supplementary Fig. 13, 14). This shows that the
250 predicted activity by DeepMEL2 correlates with the luciferase reporter activity in vitro (Fig. 5f,
251 Extended Data Fig. 5g), suggesting that the steps of increased activity are not biased to our DeepMEL2
252 model, but reflect biological activity. Functional in silico evolved enhancers lost their activity, and
253 accessibility, when ZEB sites were generated in proximity of activator sites (Fig. 5e,f, Extended Data
254 Fig. 5g, 8), and this repressive mechanism depended on the number and the strength of repressor
255 sites (Extended Data Fig. 8a,b-e, Supplementary Fig. 15). We confirmed that the same principles of
256 repression apply to genomic enhancers, using the MEL enhancer in an *IRF4* intron as example, and
257 through ChIP-seq we identified ZEB2 as the actual repressor TF (Fig. 5g,h, Supplementary Note 3).
258 Mutating the endogenous ZEB2 site in the *IRF4* enhancer causes a significant increase in activity, while
259 mutations that generate additional ZEB2 sites (without touching activator sites) decrease its activity
260 (Fig. 5i., Supplementary Note 3).

261 These findings could be further corroborated by scoring all sequences during the optimization process
262 with two other deep learning models, namely a newly trained ChromBPNet model[50] on bulk MM001
263 ATAC-seq data (Methods) and the previously published Enformer model, for which the SK-MEL-5
264 ATAC-seq class represents the MEL state[6]. The Enformer model has a receptive field of 200 kb and can
265 be used to predict both enhancer activity and target gene expression in the context of an entire gene
266 locus. To simulate whether our synthetic enhancers do function like genomic enhancers in a complex
267 locus, we replaced the *IRF4* enhancer studied above with synthetic enhancers, thus performing an in

268    silico CRISPR experiment. Replacement of the *IRF4* enhancer by a random sequence results in no
269    predicted accessibility, while replacement by different synthetic enhancers along their design path
270    gradually obtains increased prediction scores for accessibility, H3K27Ac signal, and CAGE gene
271    expression (Fig. 5j,k, Extended Data Fig. 7b). Since Enformer contains more than 600 chromatin
272    accessibility (DNase Hypersensitivity) output classes, across a wide variety of cell types, we used it to
273    assess the specificity of our designed enhancers, and found high prediction scores for only four classes,
274    each representing either melanocytes or melanocyte-like melanoma cell states (Fig. 5l, Extended Data
275    Fig. 7a). The ChromBPNet model shows continuous increases of predicted enhancer activity along the
276    optimization path (Fig. 5m). Again, all three models correctly predict that synthetic enhancers, after
277    they reach their highest activity level, can be switched off entirely by introducing point mutations that
278    generate ZEB binding sites (Fig. 5j,k,m, Extended Data Fig. 7a,b). Furthermore, changing the location
279    of the enhancer relative to the TSS did not alter its functionality, suggesting that the enhancers are
280    not dependent on the local sequence context around the *IRF4* enhancer location to be functional
281    (Extended Data Fig. 7c). As a final example of in silico evolution, we identified a human 'near-enhancer'
282    and rescued its activity with only 4 mutations (Extended Data Fig. 9a-d).
283    We also applied the motif implantation strategy to design human enhancers. We implanted SOX10,
284    MITF, and TFAP2 binding sites to 2,000 random sequences of 500 bp. While implanting only MITF or
285    TFAP2 resulted in a small increase in the prediction score, implanting SOX10 alone had the strongest
286    effect (Fig. 5n). Adding MITF and then TFAP2 on top of SOX10 sites increased the prediction scores to
287    0.6 on average. The prediction scores continued increasing even further after adding another set of
288    SOX10, MITF, and TFAP2 binding sites (Fig. 5n). We did not observe a preferential location for the
289    implantation of MITF or TFAP2 relative to SOX10, however both binding sites were located within 100
290    bp of SOX10 (Fig. 5o). The second SOX10 binding site was placed further away at a 200-250 bp distance
291    relative to the first SOX10 (Fig. 5o). We selected four sequences with either single or double SOX10,
292    MITF, and TFAP2 implanted sites and tested their activity with luciferase assays. All enhancers showed
293    activity in the range of native enhancers and adding the binding sites twice consistently increased the
294    activity of the enhancers (Fig. 5p, Extended Data Fig. 10a,b,c). Replacing the implanted binding sites
295    with their weaker versions taken from a native enhancer (IRF4) decreased the activity of the enhancers
296    dramatically (Extended Data Fig. 10a,b,c). To confirm that the activity of the enhancers was driven by
297    the implanted binding sites, we cut the sequences from the most upstream binding site to the most
298    downstream binding site. These subsequences (116-164 bp) were also active with a slight change in
299    their activity levels (Extended Data Fig. 10a,b,c). Finally, instead of choosing the best location for MITF
300    and TFAP2 implantation, we implanted them at the closest location to the SOX10 binding site that
301    would result in a positive change in the prediction score. These minimal enhancers (51-64 bp) were as
302    active as their longer (500 bp) version (Extended Data Fig. 10a,b,c).
303    Finally, we applied the GAN-based sequence generation approach to the generation of human
304    enhancers and obtained similar performances as with the *Drosophila* GAN-generated enhancers
305    (Supplementary Note 2).
306    In conclusion, these results show that enhancer design strategies are adaptable to different biological
307    systems and even other species including human.
308

309    **Discussion**
310    Understanding the code of transcriptional regulation and utilising this knowledge to design synthetic
311    enhancers has been a persistent challenge. We successfully designed synthetic enhancer sequences
312    in human and fly guided by deep learning models. By combining a stepwise enhancer design approach

alongside model interpretation techniques, we followed the trajectories of in silico enhancer emergence in *Drosophila* and human, towards local optima. Nucleotide-by-nucleotide evolution revealed that the selected mutations predominantly destroy candidate repressor TF binding sites and create candidate activator sites. Mostly, ten iterative mutations were sufficient to convert a random sequence into a cell type-specific functional enhancer. Similarly, for native yeast promoter sequences, it was recently shown that only four mutations could dramatically increase or decrease their activities[45]. This evolutionary design process may represent an optimized version of natural evolution of genomic enhancers. We found that the fly and human genomes contain "near-enhancers" that require few mutations to become functional.

The location, orientation, strength, and number of TF motifs within a single enhancer, and their distance to other motifs are important features determining an enhancer code that is unique to each cell type. This array of well-arranged TF binding sites constitutes a docking platform for a specific combination of TFs. Their cooperative binding makes the enhancer accessible/active at different levels and in different cell types. We found certain enhancers to be active in multiple cell types. Besides the trivial possibility whereby two cell types share a common set of TFs that bind to a common set of sites (e.g., different KC subtypes), we showed that some enhancers have evolved multiple intertwined codes (e.g., KCs and T neurons). We could prove this by either removing a code from a native dual-code enhancer or adding a second code to a native single-code enhancer.

The consequence of this motif-driven enhancer model is that it allows for enhancer design by motif implantation. Several studies have used motif implantation in an attempt to reconstitute enhancer activity, but successes of accurate in vivo activity have been limited[51,52]. More recently, motif embedding has also been used in combination with deep learning models[38,42,53] with the advantage that many different motif implanting scenarios can be tested in silico, before performing experimental validation[38,42,43,53], as compared to high-throughput testing of random implantations[28,54,55]. By exploiting motif implantation further, particularly by scoring each possible implant position, as well as combinations of motifs, we could reveal motif synergies (e.g., Ey + Mef2; or SOX10 + MITF), as well as preferred orientations and distances between motifs, motif strengths, and motif copy number. A minimal fly brain enhancer designed with three abutting motif instances illustrates that functional enhancers can be created without further sequence context. Compared to random insertions of motif instances[52,56], deep learning guided implantation has the capacity to take the entire enhancer sequence into account. Consequently, what makes an enhancer is not only the optimal combination of motifs used (including each motif's strength and copy number), but also the optimal balance between repressor and activator motifs, and the optimal motif arrangement.

Two out of 13 Kenyon cell enhancers remain negative while one is inconclusive. Nevertheless, this leads to a conservative success rate >75%. We also envision several routes for further improvement in enhancer design. Firstly, whereas our examples focused on adult cell types, we did not consider temporal changes. It thus remains to be investigated whether developmental enhancers with highly dynamic and complex output functions can be decoded and designed along the same principles. Studies of the *shavenbaby* enhancer in *Drosophila* showed that its output is affected by mutations in most of its nucleotides[57]. This may be due to a densely packed motif content, like our minimal enhancer, or to yet unknown sequence features. It may be interesting to investigate such developmental enhancers with deep learning models [INSERT CITATION TO FURLONG&STARK BACK-TO-BACK]. Additionally, we observed slight variations in the GFP output pattern of (genomic and synthetic) enhancers. Incorporating such high-resolution variations in the training data may yield models with improved spatial and quantitative resolution. Lastly, the repressor motifs identified by

358  our models recruit TFs that cause a decrease in chromatin accessibility. However, this is likely not true
359  for all transcriptional repressors (e.g., binding sites of the REST repressor overlap with accessible
360  chromatin[58]). A future challenge will be to take repressor motifs into account that do not decrease
361  chromatin accessibility. To train such models, additional enhancer activity data or gene expression
362  data will be needed.
363  The successful application of enhancer design on both fly brain and human cancer cells has shown that
364  simple, yet powerful strategies guided by deep learning models are adaptable to different organisms
365  or systems. Our proof-of-concept study is an encouraging step forward towards the development of
366  organism-wide deep learning models. Such models will facilitate the generation of synthetic
367  enhancers during development, disease, and homeostasis; and will further improve our understanding
368  and control of the genomic cis-regulatory code.
369
370

**References**

371  **References**

372  1. Zaret, K. S. & Carroll, J. S. Pioneer transcription factors: establishing competence for

373  gene expression. *Genes Dev* **25**, 2227–2241 (2011).

374  2. Jacobs, J. *et al.* The transcription factor Grainy head primes epithelial enhancers for

375  spatiotemporal activation by displacing nucleosomes. *Nat Genet* **50**, 1011–1020 (2018).

376  3. Payankaulam, S., Li, L. M. & Arnosti, D. N. Transcriptional repression: conserved and

377  evolved features. *Curr Biol* **20**, R764-771 (2010).

378  4. Davidson, E. H. *Genomic regulatory systems: development and evolution.* (Academic

379  Press, 2001).

380  5. Pennacchio, L. A. *et al.* In vivo enhancer analysis of human conserved non-coding

381  sequences. *Nature* **444**, 499–502 (2006).

382  6. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-

383  range interactions. *Nat Methods* **18**, 1196–1203 (2021).

384  7. Linder, J., Srivastava, D., Yuan, H., Agarwal, V. & Kelley, D. R. Predicting RNA-seq

385  coverage from DNA sequence as a unifying model of gene regulation. 2023.08.30.555582

386  Preprint at https://doi.org/10.1101/2023.08.30.555582 (2023).

387  8. Atak, Z. K. *et al.* Interpretation of allele-specific chromatin accessibility using cell state-

388  aware deep learning. *Genome Res.* gr.260851.120 (2021) doi:10.1101/gr.260851.120.

389  9. Zhou, J. & Troyanskaya, O. G. Predicting effects of noncoding variants with deep

390  learning–based sequence model. *Nature Methods* **12**, 931–934 (2015).

391  10.  Yuh, C. H., Bolouri, H. & Davidson, E. H. Genomic cis-regulatory logic: experimental

392  and computational analysis of a sea urchin gene. *Science* **279**, 1896–1902 (1998).

393  11.  Patwardhan, R. P. *et al.* Massively parallel functional dissection of mammalian

394  enhancers in vivo. *Nature Biotechnology* **30**, 265–270 (2012).

395  12.  Kheradpour, P. *et al.* Systematic dissection of regulatory motifs in 2000 predicted

396  human enhancers using a massively parallel reporter assay. *Genome Res.* **23**, 800–811

397  (2013).

398  13.  Hare, E. E., Peterson, B. K., Iyer, V. N., Meier, R. & Eisen, M. B. Sepsid even-

399      skipped Enhancers Are Functionally Conserved in Drosophila Despite Lack of Sequence

400      Conservation. *PLOS Genetics* **4**, e1000106 (2008).

401   14.    Kvon, E. Z. *et al.* Genome-scale functional characterization of Drosophila

402      developmental enhancers in vivo. *Nature* **512**, 91–95 (2014).

403   15.    Arnold, C. D. *et al.* Genome-wide quantitative enhancer activity maps identified by

404      STARR-seq. *Science* **339**, 1074–1077 (2013).

405   16.    Zinzen, R. P., Girardot, C., Gagneur, J., Braun, M. & Furlong, E. E. M. Combinatorial

406      binding predicts spatio-temporal cis-regulatory activity. *Nature* **462**, 65–70 (2009).

407   17.    May, D. *et al.* Large-scale discovery of enhancers from human heart tissue. *Nat

408      Genet* **44**, 89–93 (2011).

409   18.    Narlikar, L. *et al.* Genome-wide discovery of human heart enhancers. *Genome Res

410      * **20**, 381–392 (2010).

411   19.    Ghandi, M., Lee, D., Mohammad-Noori, M. & Beer, M. A. Enhanced Regulatory

412      Sequence Prediction Using Gapped k-mer Features. *PLOS Computational Biology* **10**,

413      e1003711 (2014).

414   20.    Kantorovitz, M. R. *et al.* Motif-Blind, Genome-Wide Discovery of cis-Regulatory

415      Modules in Drosophila and Mouse. *Developmental Cell* **17**, 568–579 (2009).

416   21.    Meuleman, W. *et al.* Index and biological spectrum of human DNase I hypersensitive

417      sites. *Nature* **584**, 244–251 (2020).

418   22.    Lareau, C. A. *et al.* Droplet-based combinatorial indexing for massive-scale single-

419      cell chromatin accessibility. *Nat Biotechnol* **37**, 916–924 (2019).

420   23.    Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human

421      immune cell development and intratumoral T cell exhaustion. *Nat Biotechnol* **37**, 925–936

422      (2019).

423   24.    Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin

424      Accessibility. *Cell* **174**, 1309-1324.e18 (2018).

425   25.    Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome.

426      *Nature* **489**, 57–74 (2012).

427   26.   Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and

428         characterization. *Nat Methods* **9**, 215–216 (2012).

429   27.   Yan, J. *et al.* Transcription Factor Binding in Human Cells Occurs in Dense Clusters

430         Formed around Cohesin Anchor Sites. *Cell* **154**, 801–813 (2013).

431   28.   Smith, R. P. *et al.* Massively parallel decoding of mammalian regulatory sequences

432         supports a flexible organizational model. *Nat Genet* **45**, 1021–1028 (2013).

433   29.   Weirauch, M. T. *et al.* Determination and Inference of Eukaryotic Transcription Factor

434         Sequence Specificity. *Cell* **158**, 1431–1443 (2014).

435   30.   Rauluseviciute, I. *et al.* JASPAR 2024: 20th anniversary of the open-access

436         database of transcription factor binding profiles. *Nucleic Acids Research* gkad1059 (2023)

437         doi:10.1093/nar/gkad1059.

438   31.   He, X., Samee, M. A. H., Blatti, C. & Sinha, S. Thermodynamics-Based Models of

439         Transcriptional Regulation by Enhancers: The Roles of Synergistic Activation,

440         Cooperative Binding and Short-Range Repression. *PLOS Computational Biology* **6**,

441         e1000935 (2010).

442   32.   Parker David S., White Michael A., Ramos Andrea I., Cohen Barak A., & Barolo

443         Scott. The cis-Regulatory Logic of Hedgehog Gradient Responses: Key Roles for Gli

444         Binding Affinity, Competition, and Cooperativity. *Science Signaling* **4**, ra38–ra38 (2011).

445   33.   Fukaya, T., Lim, B. & Levine, M. Enhancer Control of Transcriptional Bursting. *Cell*

446         **166**, 358–368 (2016).

447   34.   Deplancke, B., Alpern, D. & Gardeux, V. The Genetics of Transcription Factor DNA

448         Binding Variation. *Cell* **166**, 538–554 (2016).

449   35.   Jolma, A. *et al.* DNA-Binding Specificities of Human Transcription Factors. *Cell* **152**,

450         327–339 (2013).

451   36.   Zhu, F. *et al.* The interaction landscape between transcription factors and the

452         nucleosome. *Nature* **562**, 76–81 (2018).

453   37.   Minnoye, L. *et al.* Cross-species analysis of enhancer logic using deep learning.

454         *Genome Res.* **30**, 1815–1834 (2020).

455 38. de Almeida, B. P., Reiter, F., Pagani, M. & Stark, A. DeepSTARR predicts enhancer

456 activity from DNA sequence and enables the de novo design of synthetic enhancers. *Nat*

457 *Genet* **54**, 613–624 (2022).

458 39. Janssens, J. *et al.* Decoding gene regulation in the fly brain. *Nature* 1–7 (2022)

459 doi:10.1038/s41586-021-04262-z.

460 40. Koo, P. K., Majdandzic, A., Ploenzke, M., Anand, P. & Paul, S. B. Global importance

461 analysis: An interpretability method to quantify importance of genomic features in deep

462 neural networks. *PLOS Computational Biology* **17**, e1008925 (2021).

463 41. Karollus, A., Mauermeier, T. & Gagneur, J. Current sequence-based models capture

464 gene expression determinants in promoters but mostly ignore distal enhancers. *Genome*

465 *Biology* **24**, 56 (2023).

466 42. Avsec, Ž. *et al.* Base-resolution models of transcription-factor binding reveal soft

467 motif syntax. *Nat Genet* **53**, 354–366 (2021).

468 43. Toneyan, S., Tang, Z. & Koo, P. K. Evaluating deep learning for predicting

469 epigenomic profiles. *Nat Mach Intell* **4**, 1088–1100 (2022).

470 44. Yuan, H. & Kelley, D. R. scBasset: sequence-based modeling of single-cell ATAC-

471 seq using convolutional neural networks. *Nat Methods* **19**, 1088–1096 (2022).

472 45. Vaishnav, E. D. *et al.* The evolution, evolvability and engineering of gene regulatory

473 DNA. *Nature* 1–9 (2022) doi:10.1038/s41586-022-04506-6.

474 46. Zrimec, J. *et al.* Controlling gene expression with deep generative design of

475 regulatory DNA. *Nat Commun* **13**, 5099 (2022).

476 47. Killoran, N., Lee, L. J., Delong, A., Duvenaud, D. & Frey, B. J. Generating and

477 designing DNA with deep generative models. Preprint at

478 https://doi.org/10.48550/arXiv.1712.06148 (2017).

479 48. Alipanahi, B., Delong, A., Weirauch, M. T. & Frey, B. J. Predicting the sequence

480 specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol* **33**, 831–

481 838 (2015).

482 49. Preger-Ben Noon, E. *et al.* Comprehensive Analysis of a cis-Regulatory Region

483     Reveals Pleiotropy in Enhancer Function. *Cell Rep* **22**, 3021–3031 (2018).

484     50.     Brennan, K. J. *et al.* Chromatin accessibility in the Drosophila embryo is determined

485     by transcription factor pioneering and enhancer activation. *Developmental Cell* **58**, 1898-

486     1916.e9 (2023).

487     51.     Vincent, B. J., Estrada, J. & DePace, A. H. The appeasement of Doug: a synthetic

488     approach to enhancer biology. *Integrative Biology* **8**, 475–484 (2016).

489     52.     Swanson, C. I., Schwimmer, D. B. & Barolo, S. Rapid Evolutionary Rewiring of a

490     Structurally Constrained Eye Enhancer. *Current Biology* **21**, 1186–1196 (2011).

491     53.     Koo, P. K. & Ploenzke, M. Improving representations of genomic sequence motifs in

492     convolutional networks with exponential activations. *Nat Mach Intell* **3**, 258–266 (2021).

493     54.     King, D. M. *et al.* Synthetic and genomic regulatory elements reveal aspects of cis-

494     regulatory grammar in mouse embryonic stem cells. *eLife* **9**, e41279 (2020).

495     55.     Davis, J. E. *et al.* Dissection of c-AMP Response Element Architecture by Using

496     Genomic and Episomal Massively Parallel Reporter Assays. *Cell Systems* (2020)

497     doi:10.1016/j.cels.2020.05.011.

498     56.     Tsai, A., Alves, M. R. & Crocker, J. Multi-enhancer transcriptional hubs confer

499     phenotypic robustness. *eLife* **8**, e45325 (2019).

500     57.     Fuqua, T. *et al.* Dense and pleiotropic regulatory information in a developmental

501     enhancer. *Nature* **587**, 235–239 (2020).

502     58.     Imrichova, H. & Aerts, S. ChIP-seq meta-analysis yields high quality training sets for

503     enhancer classification. 388934 Preprint at https://doi.org/10.1101/388934 (2018).

504
505
506
507
508

## Methods

**Data reporting**

No statistical methods were used to predetermine sample size. The number of synthetic enhancers that were tested using transgenic flies was determined to be minimally 6 per cell type and it was bounded by the feasibility of the transgenic animal generation experiments. In total, 68 transgenic flies were generated. The number of synthetic enhancers that were used with luciferase assays was determined to be minimally 10 per different category (in silico evolution, motif embedding, GAN, repressors, mutational steps). In total, 97 sequences were tested using luciferase assay. The initial random sequences (used for sequence evolution and motif implantation) were sampled from the sequence space that matches the GC content of the genomic sequences. Flies fitting the gender (equal amount of male and female) and age (<10 days) criteria were selected randomly for all experiments. In this study, we didn't perform experiments that needed to be allocated into different groups. The investigators were blinded when performing cloning, transfection, antibody staining, and luciferase experiments by using enhancer IDs.

**Statistics and reproducibility**

Statistics were calculated using Scipy (v1.6.0)[59]. The results here and throughout the manuscript were visualised using matplotlib (v3.1.1)[60]. The deep learning models were run in a conda environment where python (v3.7), tensorflow-gpu (v1.15)[61], numpy (v1.19.5)[62], ipykernel (v5.1.2), and h5py (v2.10.0) packages were installed. The same results were obtained from different replication experiments. Multiple brains (at least 10) were stained and imaged for the fly experiments. Three biological replicates were performed for the main luciferase experiments. Two biological replicates were performed for the negative control luciferase experiments. No biological replicates performed for ATAC-seq or ChIP-seq experiments.

**In silico saturation mutagenesis**

To measure the effect of each possible single mutation on a given DNA sequence, we performed in silico saturation mutagenesis, as described earlier[9,48,63]. We first generated the sequences of all single mutations for a given 500 bp sequence (3 possible mutations for each nucleotide, making 1500 sequences in total). We scored these sequences and the initial sequence with the deep learning models. For a chosen class, we calculated the delta prediction score by subtracting the score of the initial sequence from the score of the mutated sequence for each mutation.

**Random sequence generation**

We generated random 500 bp sequences to use as a prior set for the in silico sequence evolution and motif implantation by using the *numpy.random.choice(["A","C","G","T"])* command. For each position, instead of using 25% probability for each nucleotide to be chosen, we used the frequency of the nucleotides from fly or human genomic regions for each position. In these genomic regions, the GC-content was higher in the center of the regions on average relative to the flankings. We used 6,126 KC regions for fly and 3,885 MEL regions for human that we identified in our previous publications[37,39].

**In silico sequence evolution**

By using the saturation mutagenesis scores mentioned above, we performed in silico sequence evolution. For the in silico evolution from random sequences, we calculated saturation mutagenesis

554     scores for a random sequence. Then, we selected the mutation that had the highest positive delta
555     prediction score for the selected class (for γ-KC, class no. 35 in DeepFlyBrain; for PNG, class no. 34 in
556     DeepFlyBrain; for MEL, class no. 16 in DeepMEL2). For the selected sequence with one mutation, we
557     re-calculated the saturation mutagenesis scores for each nucleotide and again selected the mutation
558     with the highest delta score and repeated this procedure until the initial random sequence
559     accumulated 20 mutations.
560     Even though we used a simple objective function to direct the sequence evolution towards a single
561     cell type, without explicitly penalising off-target cell types, the generated sequences were mostly
562     active only in the targeted cell type. We believe this is due to the type of enhancer models we are
563     using, which were trained on cell-type specific accessible regions. When more general models are
564     used, for example trained on entire ATAC-seq tracks, adapted objective functions can be used and are
565     available in our code. The cell type specific activity of our synthetic enhancers suggests that: (1)
566     activator binding sites were not created for other cell types; and (2) repressor sites, which are present
567     in random sequences by chance, were not destroyed for other cell types. For example, in Kenyon Cells
568     we observed that activator binding sites are usually longer than repressor sites (18 bp and 10 bp versus
569     5 bp and 6 bp for Ey, Mef2, Mamo, and CAATTA respectively). This implies that a random sequence is
570     more likely to have multiple repressor binding sites by chance compared to activator sites (Extended
571     Data Fig. 1f). Indeed, the average prediction scores of our initial 6,000 random sequences were close
572     to zero for all classes. This may at least in part explain why earlier enhancer design efforts may have
573     failed.
574     We used 6,000 initial random sequences for KC and PNG and 4,000 for MEL. For the generation of KC
575     enhancers from genomic regions, we performed 6 iterative mutations. For the multiple cell-type code
576     enhancers, we started from optic lobe enhancers and in each iteration we manually selected the
577     mutations that increased the γ-KC prediction score while maintaining the optic lobe prediction scores
578     high. For the pruning experiment of a multiple cell type code enhancer into only KC code, we manually
579     selected the mutations that maintain the γ-KC prediction score high while decreasing the optic lobe
580     prediction scores. The DeepFlyBrain class numbers used for optic lobe neurons are 23 for T1, 20 for
581     T2, and 2 for T4 neurons.
582     To rescue the designed enhancers that were weak or negative, we performed 5 additional mutations
583     on both from-scratch and from-genomic sequences.
584     To repress the sequences with the creation of repressor binding sites, we selected single or double
585     mutations manually, by going over in silico saturation mutagenesis plots calculated on the evolved
586     sequences.
587     To explore the alternative in silico sequence evolution paths besides choosing the best mutation
588     (greedy algorithm), we chose the top 20 mutations on each sequence for every incremental step
589     starting from a random sequence. We followed this procedure for 5 incremental mutational steps.
590     Starting from the random sequence used to generate enhancer KC EFS-4, we obtained 3.2 million
591     paths/sequences at the end.
592

593     **Nucleotide contribution scores**
594     We used a network explaining tool, called DeepExplainer (SHAP package[64,65]), to calculate the
595     contribution of each nucleotide to the final prediction of the deep learning model for the chosen class.
596     We used randomly selected 250 genomic regions to initialize the explainer.
597     DeepFlyBrain model takes a single strand as an input. For a given 500 bp, we multiplied the explainer's
598     output by the one-hot encoded DNA sequence and visualized it as the height of the nucleotide letters.

599 DeepMEL2 model takes forward and reverse strands separately as an input. In this case, the explainer
600 results in contribution scores for each strand. We first took the average contribution score for each
601 nucleotide and then multiplied it by the one-hot encoded DNA sequence to visualize.
602
603 **Motif annotation**
604 To identify TF binding sites during the in silico evolution of designed sequences, we used TF-Modisco
605 (v0.5.5.4)[66] and Cluster-Buster[67]. Firstly, we calculated the nucleotide contribution scores on every
606 mutational step including random sequences. Then, we ran TF-Modisco on each mutational step
607 separately to identify which patterns are appearing/disappearing. The TF-Modisco parameters we
608 used were num_to_samp=5000, sliding_window_size=15, flank_size=5, target_seqlet_fdr=0.15,
609 trim_to_window_size=15, initial_flank_to_add=5, final_flank_to_add=5, final_min_cluster_size=60.
610 After investigating the TF-Modisco patterns that were identified on each mutational steps, we used
611 mutational step 1 for KC and mutational step 4 for MEL to collect the identified patterns, since they
612 contained all the activator and repressor patterns (Earlier steps didn't have good representation of
613 activators since they are close to random sequences. Later steps didn't have good representation of
614 repressors since they were destroyed during the mutational steps). We trimmed the patterns based
615 on information content (threshold=0.1) and saved them as a .cb file to be used by the Cluster-Buster.
616 By using the TF-Modisco patterns, we ran ClusterBuster (with -c 0 and -m 3 options) to identify motifs
617 on each mutational step, including random sequences. We selected only the motif instances from
618 Cluster-Buster results and merged (by using BEDTools v2.30.0[68]) the overlapping hits of the motifs into
619 a single hit. We calculated mean+std on the hit scores coming from random sequences for each motif
620 separately and used these thresholds to get the significant hits.
621 Identification of TF binding sites similar to TF-Modisco patterns was performed using Tomtom[69] using
622 the cisTarget motif collection[70].
623
624 **Scoring the fly genome**
625 To identify the regions that have high prediction scores for γ-KC but have less accessibility in γ-KC, we
626 scored the whole fly genome. We used the *bedtools makewindows -g dm6.chromsize -w 500 -s 50*
627 command[68] to create the coordinates of the binned fly genome with a 500 bp window and 50 bp
628 stride. We removed the regions that are not exactly 500 bp. This resulted in 2,750,893 regions to be
629 scored with the DeepFlyBrain model. We used the *stats* function of deeptools/pyBigWig package[71] to
630 calculate mean γ-KC accessibility values for each bin.
631
632 **Motif implanting**
633 To implant binding sites into 500 bp sequences, we started from a random sequence. We implanted a
634 binding site into every possible location on the random sequence one-by-one by replacing the
635 nucleotides on the random sequences with the binding site. Then, we scored these sequences with
636 the model. We selected the binding site position that gives the highest prediction score and implanted
637 the motif on that position. Then, starting from this sequence with one binding site implanted, we
638 implanted the next binding sites one-by-one by using the same procedure. The sequence of binding
639 sites that maximize the TF-Modisco pattern score were selected to implant and they are as follows;
640 Ey: TGCTCACTCAAGCGTAA, Mef2: CTATTTATAG, Onecut: ATCGAT, Sr: CCACCC, SOX10:
641 AACAATGGGCCCATTGTT, MITF: GTCACGTGAC, and TFAP2: GCCTGAGGC. We used 2,000 initial random
642 sequences for KC and 2,000 for MEL. The weaker binding sites taken from the *IRF4* enhancer are as

643 follows: SOX10_1: GTGAATGACAGCTTTGTT, SOX10_2: TACAAGTATCTCCATTGT, MITF_1:
644 ATCATGTGAA, MITF_2: GCCATATGAC, TFAP2_1: TCTTCAGGC, and TFAP2_2: CCCTGTGGT.
645 When TF motifs are implanted at random positions in a random sequence, prediction scores are very
646 low, likely because repressor sites remain present. Likewise, to be able to generate a functional
647 enhancer through random sequence generation, many sequences need to be generated (i.e., 100
648 million and 1 billion[38,72]).
649 To measure if there is a preference for a flanking sequence when performing motif implanting, we
650 aggregated all the sequences aligned by the location of the implanted motif. Then, we calculated the
651 position probability matrix and visualised it by subtracting 0.25 from each position.
652 To measure the effect of different background sequences on the minimal KC enhancer, we generated
653 1 million random sequences with the size of 20 bp. Then, we replaced the 20 bp spanning the position
654 where Ey, Mef2, and Onecut binding sites implanted that occupied the 6 bp flankings on both sides
655 and 8 bp inter-motif space. Then, we scored the sequences with the model and measured the effect
656 of different backgrounds around the motif implantation area.
657
658 **Generative Adversarial Network**
659 To train a GAN model, we used Wasserstein GAN architecture with gradient penalty[73] similar to earlier
660 work[47]. The model consists of two parts: generator and discriminator. Generator takes noise as input
661 (size is 128), followed by a dense layer with 64,000 (500 * 128) units with ELU activation, a reshape
662 layer (500, 128), a convolution tower of 5 convolution blocks with skip connections, a 1D convolution
663 layer with 4 filters with kernel width 1, and finally a SOFTMAX activation layer. The output of the
664 generator is a 500 × 4 matrix, which represents one-hot encoded DNA sequence. Discriminator takes
665 500 bp one-hot encoded DNA sequence as input (real or fake), followed by a 1D convolution layer with
666 128 filters with kernel width 1, a convolution tower of 5 convolution blocks with skip connections, a
667 flatten layer, and finally a dense layer with 1 unit.
668 Each block in the convolution tower consists of a RELU activation layer followed by 1D convolution
669 with 128 filters with kernel width 5. The noise is generated by the *numpy.random.normal(0, 1,*
670 *(batch_size, 128))* command. We used a batch size of 128. For every *train_on_batch* iteration of the
671 generator, we performed 10 *train_on_batch* iteration for the discriminator. We used Adam optimizer
672 with learning_rate of 0.0001, beta_1 of 0.5, and beta_2 of 0.9. We trained the models for around
673 260,000 batch training iteration for KC and around 160,000 batch training iteration for MEL.
674 We used 6,126 KC regions for the fly model and 3,885 MEL regions for the human model, which we
675 identified in our previous publications, as real genomic sequences to train the models. After the
676 training, we sampled 6,144 (48 * batch size) sequences for KC and 3,968 (31 * batch size) sequences
677 for MEL by using the generator for every 10,000 batch training iteration. The sampled synthetic
678 sequences were generated by calculating predictions on noise and then the *numpy.argmax()*
679 command was used to convert the predictions into one-hot encoded representations.
680
681 **Background model**
682 To compare against the GAN-generated sequences, we generated random sequences in different
683 orders by using the *CreateBackgroundModel* function from the INCLUSive package[74] based on the
684 same genomic regions that we used to train GANs.
685
686 **Training ChromBPNet models**

687 For training ChromBPNet models we used a pre-released version (v1.3-pre-release) from the
688 ChromBPNet GitHub repository (https://github.com/kundajelab/chrombpnet/tree/v1.3-pre-release).
689 We followed all the preprocessing and training steps as described in the tutorial: from the aligned
690 ATAC reads in the MM001 BAM file, we made a BigWig of Tn5 insertion sites, trained a bias model
691 that predict Tn5 binding sites in non-peak regions which is then used in the ChromBPNet model to
692 filter out Tn5 bias. ChromBPNet uses 2,114 bp DNA sequence as input and predicts both the ATAC
693 track and the natural log count of the aligned reads for the central 1000 bp. To be able to score 500
694 bp DNA sequences (*IRF4* enhancer and synthetic enhancers), we used the flanking sequences of the
695 cloned/integrated enhancer sequences surrounded by the integrated cassette. Both scalar and track
696 prediction were plotted. Flanking sequences are provided in the Supplementary Code.
697
698 **Using the Enformer model**
699 We used the Enformer model to do in silico CRISPR experiments. We took the *IRF4* locus
700 (Chr6:339,010:453,698) centred by the *IRF4* enhancer (Chr6:396,104:396,604). We replaced the
701 endogenous *IRF4* enhancer with the random / evolved / repressed designed sequences and calculated
702 the prediction scores for the related cell types. The prediction scores were plotted as showing the
703 whole locus. For DNase and ChIP-Histone:H3K27ac tracks, the mean values were calculated using the
704 middle 3 bins or 1 bin spanning the enhancer location. For CAGE tracks, the mean values are calculated
705 using 1 bin spanning the TSS of *IRF4*. The index of the tracks that we used to get the prediction scores
706 are as follow; 4832: CAGE/melanoma cell line:G-361, 162: DNase/SK-MEL-5, 2162: ChIP-
707 Histone:H3K27ac/foreskin melanocyte male newborn.
708 To measure the locational effect of the designed enhancers on gene expression, chromatin
709 accessibility, and histone modification, we moved the synthetic enhancer around the *IRF4* locus; (1)
710 to 10 kb upstream, (2) 5 kb upstream (which is next to the promoter of the *IRF4* gene), and (3) 17.5 kb
711 downstream of the original location.
712
713 **Cloning of synthetic *Drosophila* enhancers**
714 Synthetic sequences were ordered from Twist Bioscience, pre-cloned in the pTwist ENTR vector. The
715 motif-implantation and double-coded sequences were synthesized with an additional 5' CACC
716 sequence as double-stranded DNA (gBlocks Gene Fragments) by IDT. 49 bp motif-implantation
717 sequence was ordered from IDT as forward and reverse single-stranded DNA oligos, which were then
718 annealed for 5 min at 95°C and cooling down to RT over one hour. The double-stranded DNA
719 sequences were then cloned into the pENTR/D-TOPO plasmid (Invitrogen).
720 All sequences were introduced in a modified pH-Stinger vector[75], containing nuclear GFP, Hsp70
721 promoter, gypsy insulators, and attB site for phiC31 integration, via Gateway LR recombination
722 reaction (Invitrogen). 2 µl of the reaction was transformed into 25 µl of Stellar chemically competent
723 bacteria (Takara). Plasmid minipreps were performed using the NucleoSpin Plasmid Transfection-
724 grade Mini kit (Macherey-Nagel) and sequenced with Sanger sequencing to confirm the correct
725 insertion of the regions in the destination plasmid. After confirmation of the sequence, plasmid
726 midipreps were performed using the NucleoBond Xtra endotoxin-free Midi kit (Macherey-Nagel).
727 Next, the plasmids were sent to FlyORF (CH) for injection in *Drosophila* embryos (21F site on
728 chromosome 2l) and positive transformants were selected based on eye colour.
729 *Drosophila* flies were raised on a yeast-based medium at 25°C under a 12 h-12 h day-night light cycle.
730
731 **Immunohistochemistry analysis of *Drosophila* brains**

732 Brains of adult flies (*Drosophila melanogaster*, <10 days old, equally mixed sex) were dissected in PBS
733 and transferred to a tube for fixation in 4% formaldehyde in PBS for 20 min. All incubations were done
734 at room temperature, unless otherwise indicated. Brains were washed in PBS with 0.3% Triton-X
735 (PBST) three times for 10 min each, then they were placed in blocking solution (5% normal goat serum
736 (Abcam) in PBST) for 3 hours. We incubated the brains overnight at 4°C in primary antibodies diluted
737 in blocking solution (rabbit anti-GFP, IgG (Invitrogen), 1:1000 and mouse anti-Dachshund, mAB dac1-
738 1 (DSHB), 1:250). The brains were then washed in PBST three times for 10 min each and incubated
739 with the fluorochrome-conjugated secondary antibodies diluted in blocking solution for 2 hours (Alexa
740 Fluor 488 donkey anti-rabbit IgG (Invitrogen), 1:500 and Alexa Fluor 647 goat anti-mouse IgG
741 (Invitrogen), 1:500). Next, brains were washed in PBS three times for 10 min each. Finally, samples
742 were mounted onto microscope slides with Prolong Glass Antifade Mountant (Invitrogen).
743 For image acquisition, a Zeiss LSM900 microscope equipped with Airyscan2 in combination with a 20x
744 objective (Plan Apo 0,80 Air) was used. The setup was controlled by ZEN blue (version 3.4.91, Carl Zeiss
745 Microscopy GmbH). GFP was excited with a blue diode 100mW at 488 nm and tiled images were
746 collected with emission filter BP450-490/BS495/BP500-550.
747

748 **Cloning of synthetic human enhancers**
749 500 bp synthetic sequences were ordered from Twist Bioscience, pre-cloned in the pTwist ENTR
750 vector. 500 bp regions were introduced in the pGL4.23-GW luciferase reporter vector (Promega) via
751 Gateway LR recombination reaction (Invitrogen) and 2 µl of the reaction was transformed into 25 µl
752 of Stellar chemically competent bacteria (Takara).
753 Synthetic sequences shorter than 150 bp were ordered as gBlocks from IDT (Integrated DNA
754 Technologies) with 5' (cccgtcgacgaattctgcagatatcacaagtttgtacaaaaaagcaggct) and 3'
755 (acccagctttcttgtacaaagtggtgataaacccgctgatcag) adaptors. The pGL4.23-GW luciferase reporter vector
756 was linearized via inverse PCR with primers Lin_pSA335_short_ME_For (gtggtgataaacccgctgatcag) and
757 Lin_pSA335_short_ME_Rev (tctgcagaattcgtcgacggg). The short sequences and the linearized vector
758 were combined in an NEBuilder reaction (New England Biolabs, Ipswich, MA) and 2 µl of the reaction
759 was transformed into 25 µl of Stellar chemically competent bacteria.
760 For all cloning procedures, plasmid minipreps were performed using the NucleoSpin Plasmid
761 Transfection-grade Mini kit (Macherey-Nagel) and sequenced with Sanger sequencing to confirm the
762 correct insertion of the regions in the destination plasmid.
763 To generate stable cell lines with synthetic enhancers, the synthetic sequences were cloned into the
764 pSA351_SCP1_intron_eGFP vector (Addgene #206906). The vector was linearized via inverse PCR with
765 primers Lin_pSA351_For (ctgagctccctagggtact) and Lin_pSA351_Rev (cgactcgaggctagtctc). The
766 synthetic sequences were PCR-amplified from their respective pGL.23-GW vector with their respective
767 primer pairs: MM_EFS_1_For (gagactagcctcgagtcgctgattgtttgaaccattgttacgatttgg) and
768 MM_EFS_1_Rev (agtaccctagggagctcagcaattttgtttttttgcgcgtgac) for MM-EFS-1 sequences;
769 MM_EFS_4_For (gagactagcctcgagtcgtgatatgtattcacccatgccctca) and MM_EFS_4_Rev
770 (agtaccctagggagctcaagggtttgtatatgtatgctcctttatacga) for MM-EFS-4 sequences; MM_EFS_8_For
771 (gagactagcctcgagtcgatacgcacgacaaagcctcat) and MM_EFS_8_Rev
772 (agtaccctagggagctcacactgtacaaggcatcccgc) for MM-EFS-8 sequences; IRF_4_For
773 (gagactagcctcgagtcggctgccattggtgtggattttaag) and IRF_4_Rev (agtaccctagggagctcaactggcatcgagacggg)
774 for IRF-4 sequences. The PCR amplicons and the linearized vector were combined in an NEBuilder
775 reaction and 2 µl of the reaction was transformed into 25 µl of Stellar chemically competent bacteria.
776 Plasmid minipreps were performed using the NucleoSpin Plasmid Transfection-grade Mini kit

777 (Macherey-Nagel) and sequenced with Sanger sequencing to confirm the correct insertion of the
778 regions in the vector. After confirmation of the sequence, plasmid maxipreps were performed using
779 the NucleoBond Xtra endotoxin-free Maxi kit (Macherey-Nagel).
780
781 **Transfection and luciferase assay**
782 MM001 and MM047 were seeded in 24-well plates and transfected with 400 ng pGL4.23-enhancer
783 vector + 40 ng pRL-TK *Renilla* vector (Promega) with Lipofectamine 2000 (Thermo Fisher Scientific). As
784 positive controls, the previously published enhancers MLANA_5-I, IRF4_4-I and TYR_-9-D or
785 ABCC3_11-I and GPR39_23-I were used for MM001 and MM047 respectively[76]. One day after
786 transfection, luciferase activity was measured via the Dual-Luciferase Reporter Assay System
787 (Promega) by following the manufacturer's protocol. Briefly, cells were lysed with 100 µl of Passive
788 Lysis Buffer for 15 min at 500 rpm. 20 µl of the lysate was transferred in duplicate in a well of an
789 OptiPlate-96 HB (PerkinElmer, Waltham, MA) and 100 µl of Luciferase Assay Reagent II was added in
790 each well. Luciferase-generated luminescence was measured on a Victor X luminometer
791 (PerkinElmer). 100 µl of the Stop & Glo Reagent was added to each well, and the luminescence was
792 measured again to record *Renilla* activity. Luciferase activity was estimated by calculating the ratio
793 luciferase/*Renilla*; This value was normalized by the ratio calculated on blank wells containing only
794 reagents. Three biological replicates were done per condition for MM001 and two biological replicates
795 for MM047.
796
797 **Production of lentivirus**
798 The lentivirus plasmids were transfected in HEK 293T cells by use of the Lipofectamine 3000 reagent
799 (Thermo Fisher Scientific). 30 µg of pooled plasmid DNA was combined with 20 µg of a Pax2 plasmid
800 (Addgene #12260) and 10 µg of the MD2.G plasmid (Addgene #12259). 48 hours post-transfection,
801 medium was collected and refreshed. 72 hours post-transfection, medium was collected a second
802 time. Both medium collections were combined and spun down for 5 min at 1,500 rpm. Supernatants
803 was carefully collected with a blunt needle and a syringe and filtered through a 45 µm syringe disc
804 filter (Millex-HV Millipore) into an Ultra-15 MWCO100 centrifugal filter (Amicon). The concentrator
805 tube containing the supernatants was spun down at 4,000 rpm for approximately 45 min until the
806 desired volume of 250 µl was reached. The virus suspension was aliquoted and stored at -80°C.
807
808 **Transduction of melanoma cells**
809 The MM001 cells were seeded into a 6-well plate at a density of 250,000 cells per well. Transduction
810 was performed by adding 5-40 µl of lentivirus and Polybrene at 8 µg/ml. Cells were incubated for 24h
811 before washing away the Polybrene with PBS and with growth medium. After 3 days the cells were
812 split and expanded further.
813
814 **OmniATAC-seq**
815 Omni-assay for transposase-accessible chromatin using sequencing (OmniATAC-seq) was performed
816 as described previously[77]. Briefly, 50,000 MM001 cells transduced with the enhancer pools were
817 resuspended in 50 µL of cold ATAC-seq resuspension buffer (RSB; 10 mM TrisHCl pH 7.4, 10 mM NaCl,
818 and 3 mM MgCl2 in water) containing 0.1% NP40, 0.1% Tween-20, and 0.01% digitonin by pipetting
819 up and down three times. This cell lysis reaction was incubated on ice for 3 min. After lysis, 1 mL of
820 ATAC-seq RSB containing 0.1% Tween-20 was added, and the tubes were inverted to mix. Nuclei were
821 then centrifuged for 10 min at 500 g in a pre-chilled (4°C) fixed-angle centrifuge. Supernatant was

822    removed and nuclei were resuspended in 50 μL of transposition mix (25 μL 2x TD buffer, 2.5 μL
823    transposase (Nextera Tn5 transposase, Illumina), 16.5 μL PBS, 0.5 μL 1% digitonin, 0.5 μL 10% Tween-
824    20, and 5 μL water) by pipetting up and down six times. Transposition reactions were incubated at
825    37°C for 30 min in a thermoblock. Reactions were cleaned-up by MinElute (Qiagen). Transposed DNA
826    was            amplified          (10         cycles)         with          primers         i5_Indexing_For
827    (aatgatacggcgaccaccgagatctacacnnnnnnnnntcgtcggcagcgtcagatgtg)              and          i7_Indexing_Rev
828    (caagcagaagacggcatacgagatnnnnnnngtctcgtgggctcggagatgt).  All  libraries  were  sequenced  on  a
829    NextSeq2000 instrument (Illumina).
830    Reads         were         demultiplexed          using         bcl2fastq         (v2.20;
831    https://emea.support.illumina.com/sequencing/sequencing_software/bcl2fastq-conversion-
832    software.html).       Adapters       were       trimmed       by       trimgalore       (v0.6.7;
833    https://github.com/FelixKrueger/TrimGalore). Reads were mapped to a custom hg38 genome, which
834    contains integrated sequences as additional chromosomes, using bwa-mem2 (v2.2.1)[78]. By using
835    SAMtools (v 1.16.1)[79], reads were sorted and deduplicated, and reads from the blacklisted regions
836    (https://www.encodeproject.org/files/ENCFF356LFX/)  were  cleaned.  Bigwig  files  with  RPGC
837    normalisation were generated by using deepTools (v3.5.0) bamCoverage[71].
838
839    **ChIP-seq**
840    ChIP-seq was performed by following the Myers Lab ChIP-seq Protocol v011014 on 2x10⁷ MM001 cells.
841    5 μg of rabbit anti-ZEB2 antibody (1 mg/ml; Bethyl A302-473A) was used for ChIP. 15 ng of
842    immunoprecipitated DNA was used to perform library preparation according to the Illumina TruSeq
843    DNA Sample preparation guide. Briefly, the immunoprecipitated DNA was end-repaired, A-tailed, and
844    ligated to diluted sequencing adapters (1/100). After PCR amplification with i5_Indexing_For and
845    i7_Indexing_rev (18 cycles) and bead purification (Agencourt AmpureXP, Analis), the libraries with
846    fragment size of 300-500 bp were sequenced using the NextSeq2000 instrument (Illumina).
847    Reads were demultiplexed using bcl2fastq (v2.20). Adapters were trimmed by trimgalore (v0.6.7).
848    Reads were mapped to hg38 using bwa-mem2 (v2.2.1)[78]. By using SAMtools (v 1.16.1)[79], reads were
849    sorted    and    deduplicated,    and    reads    from    the    blacklisted    regions
850    (https://www.encodeproject.org/files/ENCFF356LFX/)  were  cleaned.  Bigwig  files  with  RPGC
851    normalisation were generated by using deepTools (v3.5.0) bamCoverage[71]. Peaks were called using
852    MACS2 (v2.1.2.1) callpeak[80].
853
854    **Cell lines**
855    MM001, MM047, and MM099 were obtained from Prof. Dr. Ghanem Ghanem and were cultured in
856    Ham's F-10 Nutrient Mix (Invitrogen) + 10% FBS (Invitrogen). We authenticated the cell lines by
857    checking their genomic, transcriptomic, and epigenomic profiles[8,81,82]. HEK293T used for lentivirus
858    production was obtained from ATCC (CAT# CRL-3216) and were cultured in DMEM (Invitrogen) + 10%
859    FBS (Invitrogen). Cell lines were tested for mycoplasma contamination prior to experiments, and were
860    found negative.
861
862    **Code availability**
863    Code used to load deep learning models, create random sequences, perform sequence evolution,
864    perform motif implantation, and train GAN models together with the IPython Notebooks that
865    reproduces all the figures were provided as Supplementary Code. The data to run the scripts, the

models, and the intermediate files can be found together with the code here 10.5281/zenodo.10184648.

**Data availability**

Cloned *Drosophila* and human sequences were provided as Supplementary Tables. DeepMEL, DeepMEL2, and DeepFlyBrain deep learning model files were obtained from Kipoi[83] (http://kipoi.org/models/DeepMEL, https://kipoi.org/models/DeepFlyBrain) with Zenodo record ids 3592129, 4590308, and 5153337. The fasta files used to train GAN models and the trained GAN models are available on Zenodo at https://doi.org/10.5281/zenodo.6701504. Custom genomes (hg38 and dm6) generated in this study are available on Zenodo at https://doi.org/10.5281/zenodo.10184648. Chromatin accessibility values in Kenyon Cells in adult *Drosophila* brains were obtained from GSE163697[39]. In vitro saturation mutagenesis on *IRF4* data was obtained from https://kircherlab.bihealth.org/satMutMPRA/ [84]. Chromatin accessibility of *Drosophila* and transduced melanoma lines and ZEB2 ChIP-seq data generated for this study have been submitted to the NCBI Gene Expression Omnibus (GEO, https://www.ncbi.nlm.nih.gov/geo/) under accession number GSE240003.

**Additional references**

59. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* **17**, 261–272 (2020).

60. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering* **9**, 90–95 (2007).

61. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. (2015).

62. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (2020).

63. Kelley, D. R., Snoek, J. & Rinn, J. L. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* **26**, 990–999 (2016).

64. Shrikumar, A., Greenside, P. & Kundaje, A. Learning Important Features Through Propagating Activation Differences. Preprint at https://doi.org/10.48550/arXiv.1704.02685 (2019).

65. Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. in *Proceedings of the 31st International Conference on Neural Information Processing Systems* 4768–4777 (Curran Associates Inc., 2017).

900   66.   Shrikumar, A. *et al.* Technical Note on Transcription Factor Motif Discovery from

901       Importance Scores (TF-MoDISco) version 0.5.6.5. Preprint at

902       https://doi.org/10.48550/arXiv.1811.00416 (2020).

903   67.   Frith, M. C., Li, M. C. & Weng, Z. Cluster-Buster: finding dense clusters of motifs in

904       DNA sequences. *Nucleic Acids Res* **31**, 3666–3668 (2003).

905   68.   Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing

906       genomic features. *Bioinformatics* **26**, 841–842 (2010).

907   69.   Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. S. Quantifying

908       similarity between motifs. *Genome Biology* **8**, R24 (2007).

909   70.   Bravo González-Blas, C. *et al.* SCENIC+: single-cell multiomic inference of

910       enhancers and gene regulatory networks. *Nat Methods* **20**, 1355–1367 (2023).

911   71.   Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing

912       data analysis. *Nucleic Acids Research* **44**, W160–W165 (2016).

913   72.   Sahu, B. *et al.* Sequence determinants of human gene regulatory elements. *Nat

914       Genet* 1–12 (2022) doi:10.1038/s41588-021-01009-4.

915   73.   Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V. & Courville, A. Improved Training

916       of Wasserstein GANs. Preprint at https://doi.org/10.48550/arXiv.1704.00028 (2017).

917   74.   Thijs, G. *et al.* INCLUSive: INtegrated Clustering, Upstream sequence retrieval and

918       motif Sampling. *Bioinformatics* **18**, 331–332 (2002).

919   75.   Aerts, S. *et al.* Robust Target Gene Discovery through Transcriptome Perturbations

920       and Genome-Wide Enhancer Predictions in Drosophila Uncovers a Regulatory Basis for

921       Sensory Specification. *PLOS Biology* **8**, e1000435 (2010).

922   76.   Mauduit, D. *et al.* Analysis of long and short enhancers in melanoma cell states. *eLife*

923       **10**, e71735 (2021).

924   77.   Corces, M. R. *et al.* An improved ATAC-seq protocol reduces background and

925       enables interrogation of frozen tissues. *Nat Methods* **14**, 959–962 (2017).

926   78.   Vasimuddin, Md., Misra, S., Li, H. & Aluru, S. Efficient Architecture-Aware

927       Acceleration of BWA-MEM for Multicore Systems. in *2019 IEEE International Parallel and*

928     *Distributed Processing Symposium (IPDPS)* 314–324 (2019).

929     doi:10.1109/IPDPS.2019.00041.

930  79.     Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,

931     2078–2079 (2009).

932  80.     Gaspar, J. Improved peak-calling with MACS. Preprint at

933     https://doi.org/10.1101/496521 (2018).

934  81.     Verfaillie, A. *et al.* Decoding the regulatory landscape of melanoma reveals TEADS

935     as regulators of the invasive cell state. *Nature Communications* **6**, 6683 (2015).

936  82.     Wouters, J. *et al.* Robust gene expression programs underlie recurrent cell states

937     and phenotype switching in melanoma. *Nature Cell Biology* **22**, 986–998 (2020).

938  83.     Avsec, Ž. *et al.* The Kipoi repository accelerates community exchange and reuse of

939     predictive models for genomics. *Nat Biotechnol* **37**, 592–600 (2019).

940  84.     Kircher, M. *et al.* Saturation mutagenesis of twenty disease-associated regulatory

941     elements at single base-pair resolution. *Nat Commun* **10**, 3583 (2019).

942

## Acknowledgements

## Author contributions

I.I.T. and S.A. conceived the study; I.I.T. performed all computational analyses and designed synthetic enhancers; V.C. performed enhancer cloning with assistance from K.I.S and D.M.; V.C. performed luciferase assays with assistance from D.M.; K.I.S. performed antibody staining and visualization with assistance from I.I.T., H.D., and J.I.; R.V. performed lentivirus production and cell line transduction; V.C. performed ATAC-seq and ChIP-seq experiments with assistance from H.D., K.T., and A.P.; G.H. performed ATAC-seq and ChIP-seq data preprocessing; N.K. trained ChromBPNet models with assistance from E.C.E.; I.I.T. and S.A. wrote the manuscript with assistance from D.M..

## Competing interest

963    The authors declare no competing interests.

964

965    **Figure legends**

966    **Figure 1: Deep learning based enhancer design**
967    Overview of enhancer design strategies and activity measurements of designed enhancers in *Drosophila* brains and human
968    cell lines.

969

970    **Figure 2: In silico sequence evolution towards functional enhancers**
971    **a**, Prediction score distribution of the sequences for the γ-KC class (⬚ = 6,000 sequences) after each mutation. The box plots
972    show the median (center line), interquartile range (box limits), and 5th and 95th percentile range (whiskers) for KC-directed
973    (blue) or random drift (orange) mutations. **b**, Nucleotide contribution scores for the γ-KC class of a selected random sequence
974    in its initial form (top) and after in silico evolution (bottom). **c**, Prediction scores of 13 selected sequences at each mutational
975    step. Dashed line indicates the selected iteration (10$^{th}$ or 15$^{th}$ mutation). **d**, In vivo enhancer activity of the cloned KC
976    sequences with positive enhancer activity. **e**, Nucleotide contribution scores for the PNG class of the same selected random
977    sequence as in panel **b** (top) and after PNG-directed mutations (bottom). **f**, In vivo enhancer activity of the cloned PNG
978    sequences. Top-middle: initial random sequence, top-left: random sequence after 10 mutations toward KC evolution, top-
979    right: random sequence after 15 mutations toward PNG evolution, bottom: Three other random sequences after mutations
980    toward PNG evolution. **g**, In vivo enhancer activity of the cloned genomic sequences with 6 mutations (11 for FP3). **h**,
981    Nucleotide contribution scores of a selected genomic sequence in its initial form (top) and after 6 iterations (bottom). In
982    panels **b**, **e**, **h**, dashed line shows the position of the mutations, the mutational order and type of nucleotide substitutions
983    are written in between top and bottom plots, and motif annotation is indicated with strong (s) or weak (w) motif instances.
984    In panels **d**, **f**, **g**, the expected location of KC is shown with dashed circles. Scale bars, 100 μm.

985

986    **Figure 3: Spatial expansion and restriction of enhancer activity**
987    **a**, Nucleotide contribution score and delta prediction score for in silico saturation mutagenesis of the EFS-4 enhancer after
988    10 mutations (first and second row) and after adding repressors (third and fourth row). Dashed line shows the position of
989    the mutations. Black circles: selected mutations to generate repressor sites. Motif annotation is indicated with strong (s) or
990    weak (w) motif instances. **b**, In vivo enhancer activity of enhancers before (top-left) and after adding repressor sites. **c**,
991    Chromatin accessibility profile near the *amon* gene. **d**, In vivo enhancer activity of the *amon* enhancer. **e**, *amon* enhancer
992    prediction scores for each cell-type. **f**, Prediction scores for the γ-KC and T4 classes after each mutational step. **g**, In vivo
993    enhancer activity of the *amon* enhancer after 13 mutations. The *amon* enhancer conserved exactly the same pattern of
994    activity for T4 following incorporation of the KC code. **h**, Number of regions that score high (>0.3) for multiple cell-types. **i**,
995    Comparison between γ-KC and T1 prediction score for the accessible regions in fly brain (⬚ = 95,931). The selected region
996    with high γ-KC and T1 prediction is highlighted with a blue dot. **j**, Chromatin accessibility profile of this region (*Pkc53e*) in
997    multiple cell-types. **k**, In vivo enhancer activity of the *Pkc53e* enhancer. **l**, *Pkc53e* enhancer prediction scores for each cell-
998    type. **m**, Prediction scores for the γ-KC, T1, and T2 classes after each mutational step. **n**, In vivo enhancer activity of the multi
999    cell-type enhancer after 9 mutations. In panels **b**, **d**, **g**, **k**, **n** dashed circles show the expected location of KC. Scale bars,
1000   100 μm. In panels **c**, **e**, **j**, **l**, AST: astrocytes; CTX: cortex glia; ENS: ensheathing glia; PNG: perineurial glia; SUB: subperineurial
1001   glia; T1-T5: T1-T5 neurons; α/β: α/β-Kenyon cells; α'/β': α'/β'-Kenyon cells; γ: γ-Kenyon cells.

1002

1003   **Figure 4: Motif implantation towards minimal enhancer design**
1004   **a**, Prediction score distribution of the sequences for the γ-KC class (⬚ = 2,000 sequences) after each motif implantation at
1005   best location (blue), random location (orange), and after 15 mutations (Nuc-15). The box plots show the median (center line),
1006   interquartile range (box limits), and 5th and 95th percentile range (whiskers). **b**, Distribution of Mef2 locations relative to Ey
1007   (⬚ = 2,000). **c**, Distribution of Onecut locations relative to Ey (⬚ = 2,000). **d**, Prediction scores for motif implanted sequence
1008   (ME-1) after each motif implanting and generation of repressor sites. **e**, Nucleotide contribution scores of ME-1 in its initial
1009   form (first track) and after Ey, Mef2, and Onecut implantations (second track). Dashed lines show the position of the motifs.
1010   Delta prediction score for in silico saturation mutagenesis (third track). Black circles: selected mutations to generate
1011   repressor sites. Nucleotide contribution scores after generation of repressor sites (fourth track). Dashed lines show the
1012   position of the mutations. **f-g**, In vivo enhancer activity of the cloned 500 bp sequence with Ey, Mef2, and Onecut
1013   implantations (**f**) and after generation of repressor sites (**g**). **h**, Zoom into the selected 49 bp part of the 500 bp sequence
1014   from **e**. The size of the motifs, the spaces between motifs, and the flankings are shown at the bottom. **i**, In vivo enhancer

1015 activity of the cloned 49 bp sequence with Ey, Mef2, and Onecut implantations. In panels **f, g, i**, the expected location of γ-
1016 KC is shown with dashed circles. Scale bars, 100 μm. Abbreviations in **a**: Ey (E), Mef2 (M), Onecut (O), and Sr (S).

1017
1018 **Figure 5: Human enhancer design**
1019 **a-b,** Prediction score distribution (MEL class, ▯=4,000 sequences (**a**) and 10 selected sequences (**b**)) after each mutation. **c,**
1020 Nucleotide contribution scores of a synthetic sequence pre (top) and post (bottom) 15 mutations, with binding site names,
1021 mutation positions (dashed lines) and orders (between top and bottom plots). **d,** Mean luciferase signal (log$_2$ fold-change
1022 over *Renilla*) of synthetic sequences from in silico sequence evolution and genomic enhancers. **e,** MM001 ATAC-seq profile
1023 of 3 integrated EFS reporters: initial, evolved and evolved with repressor sites. Red lines: enhancer boundaries. **f,** DeepMEL2
1024 prediction scores (left), luciferase activity (middle) and their correlation (right) for EFS-4 sequences. **g,** MM001 ATAC-seq ,
1025 SOX10 and ZEB2 ChIP-seq tracks for *IRF4* gene; enhancer location in red. **h,** ZEB2 ChIP-seq signal (x-axis), SOX10 ChIP-seq
1026 signal (y-axis), and ATAC-seq signal (color) for top ZEB2 regions in MM001. **i,** in vitro and in silico saturation mutagenesis
1027 values of the *IRF4* enhancer. **j,** Enformer predictions for EFS-4 sequences replacing *IRF4* enhancer: initial score and score
1028 changes post-mutations. **k**, Enformer predictions per mutation step and after repressor addition for MEL EFS sequences. **l,**
1029 Prediction scores for top 50 DNase tracks for EFS-4 sequences. Four first tracks are foreskin melanocyte male newborn and
1030 SK-MEL-5 tracks. **m,** ChromBPNet ATAC MM001 (MEL) and MM047 (MES) prediction scores for EFS sequences, across
1031 mutations and post-repressor addition. **n,** Prediction score distribution for MEL class (▯=2,000 sequences) after motif
1032 implantation. **o,** Relative TF locations distribution (▯=2,000). **p,** Luciferase signal (log$_2$ fold-change over *Renilla*) comparison
1033 of motif-implanted sequences and genomic enhancers. In **a, n**, box plots show the median (center line), interquartile range
1034 (box limits), and 5th and 95th percentile range (whiskers). Error bars in **d, f, p** denote mean standard error (▯=3 biological
1035 replicates). In **n, p**, S:SOX10, M:MITF, T:TFAP2.

1036
1037 # Extended data figure legends

1038
1039 **Extended Data Figure 1: In silico sequence evolution from random sequences**
1040 **a,** Distribution of GC-content in GC-adjusted random sequences (green) and fly genomic regions (red). **b,** Prediction score
1041 distribution of the sequences (▯ = 6,000 sequences) for all classes after 10 mutations. The KC specific classes and their class
1042 number are indicated. In **b, c**, the box plots show the median (center line), interquartile range (box limits), and 5th and 95th
1043 percentile range (whiskers). **c,** Prediction score distribution of the sequences that do not reach 0.5 prediction score threshold
1044 after 15 mutations for the γ-KC class (▯ = 180 sequences) after each mutation. **d,** Distribution of distances (▯ = 6,000) between
1045 farthest mutations on each sequence after 10 iterative mutations. The orange line shows the median. **e,** Location of the
1046 generated mutations across the random sequences (▯ = 6,000 sequences). **f,** Average number of motif hits at each mutational
1047 step compared to genomic enhancers. **g,** Delta number of motifs in each mutational step. The TF-Modisco patterns and the
1048 most similar position weight matrices from the cisTarget motif database are shown at the top of each plot. The patterns that
1049 are upside-down are the ones contributing negatively to the model's prediction and they are destroyed by the model on
1050 each step. **h,** Top panel: Dachshund staining (red) highlights KC location in the fly brain. Bottom panel: colocation of the
1051 Dachshund (red) and GFP (green) staining from enhancer EFS-13. **i,** In vivo enhancer activity of the cloned sequences with no
1052 or weak enhancer activity. **j,** Prediction scores, at each mutational step, of 4 sequences with no enhancer activity after 10
1053 mutations. The selected iterations (10$^{th}$ and 15$^{th}$ mutations) are indicated with a dashed line. **k,** Dachshund (red) and GFP
1054 (green) staining for three negative enhancers. **l,** *Drosophila* adult brain bulk-ATAC-seq profile of 6 transgenic flies that have
1055 the designed enhancers integrated. The chromatin accessibility profile of the integrated enhancers (left) and two control
1056 regions gish enhancer (middle) and Appl enhancer (right) are shown. **m,** Prediction scores, at each mutational step, of 3 EFS
1057 sequences. The selected iterations to study intermediate mutational steps (0, 2, 4, 6, 8, 10 mutations) are indicated with a
1058 dashed line. **n,** In vivo enhancer activity of fly lines with subsequent mutational steps. After 8 mutations of a random
1059 sequence, the enhancer becomes active in all three lines (EFS-3, 4, and 7) marked by GFP expression. In panels **h, i, k, n**, the
1060 expected location of γ-KC is shown with dashed circles. Scale bars, 100 μm.

1061
1062 **Extended Data Figure 2: State space optimization, design of perineurial glia enhancers and modification of genomic
1063 sequences toward KC enhancers**
1064 **a,** Prediction score distribution for 3 million sequences generated by selecting the top 20 best mutations for 5 incremental
1065 mutational steps. Blue line represents the path that was taken by the greedy algorithm. **b,** Zoomed-in version of panel **a** to
1066 the sequences that have higher prediction score than 0.25. **c,** Prediction score of evolved sequences by greedy algorithm
1067 (EFS-4) vs the best of 3 million sequences on each mutational step. **d,** Nucleotide contribution score of the original and
1068 evolved sequences as well as delta prediction score of in silico saturation mutagenesis for EFS-4 (top) and the top scoring
1069 sequence (bottom) **e,** Prediction scores of 6 selected PNG sequences at each mutational step for PNG model (left) and KC

1070   model (right). The selected iteration (15th mutation) is indicated with a dashed line. **f**, In vivo enhancer activity of the cloned
1071   PNG sequences with no enhancer activity. **g**, Comparison between γ-KC prediction score and mean γ-KC accessibility for the
1072   binned fly genome regions. The selected regions with high prediction and low accessibility are highlighted with blue, orange,
1073   green, and red dots. **h**, γ-KC ATAC-seq profile of the four selected regions. The exact location of the regions is indicated with
1074   dashed lines. **i**, Prediction scores of 4 selected KC near-enhancer sequences at each mutational step for KC model. The
1075   selected iteration (6th mutation) is indicated with a dashed line. After the 6th mutation, 4 more mutations are performed in
1076   FP-3 to improve prediction score while 7 or 8 mutations are performed in the three other sequences to generate repressor
1077   sites. **j**, In vivo enhancer activity of the cloned WT genomic "near-enhancer" sequences with no enhancer activity. The
1078   expected location of KC is shown with dashed circles. Scale bars, 100 μm.
1079
1080   **Extended Data Figure 3: Enhancer design towards multiple cell type codes**
1081   **a**, Chromatin accessibility profile near *CG15117* gene. **b**, In vivo enhancer activity of the wild-type ( WT) *CG15117* enhancer.
1082   **c**, *CG15117* enhancer prediction scores for each cell type (top) and prediction scores for the γ-KC and T1 classes after each
1083   mutational step. **d**, Nucleotide contribution scores of WT *CG15117* enhancer sequence and after 14 mutations for T1 (top)
1084   and γ-KC (bottom). **e**, Nucleotide contribution scores of WT *amon* enhancer sequence and after 13 mutations for T4 (top)
1085   and γ-KC (bottom). **f**, In vivo enhancer activity of the WT *CG15117* enhancer after 14 mutations. The CG15117 enhancer
1086   displayed a slightly altered T1 pattern following incorporation of the KC code. **g**, Nucleotide contribution scores of WT *Pkc53e*
1087   enhancer sequence and after 9 mutations for T2 (top), T1 (middle) and γ-KC (bottom). In panels **b** and **e**, the expected location
1088   of KC is shown with dashed circles. Scale bars, 100 μm. In panels **d**, **f**, **g**, the position of the mutations is shown with dashed
1089   lines, the mutational order is written in-between top and bottom plots, and motif annotation is indicated with strong (s) or
1090   weak (w) motif instances.
1091
1092   **Extended Data Figure 4: Enhancer design by motif implanting**
1093   **a**, Preferred nucleotides flanking implanted motifs (⬚ = 2,000). Dashed lines indicate the boundaries of the motifs. **b**,
1094   Distribution of Onecut locations relative to Mef2, Sr to Ey, Sr to Mef2, and Sr to Onecut, respectively (⬚ = 2,000). **c**,
1095   Distribution of Mef2 locations relative to Ey when both are on same strand, Ey is on the negative strand, Mef2 is on the
1096   negative strand, and both are on the negative strand, respectively (⬚ = 2,000). **d**, Distribution of Onecut locations relative to
1097   Ey when Ey is on the positive strand and when Ey is on the negative strand, respectively (⬚ = 2,000). **e**, DeepFlyBrain KC
1098   prediction score of the ME-2 sequence after consecutive motif implanting. **f**, In vivo enhancer activity of ME-2 enhancer. The
1099   expected location of KC is shown with dashed circles. Scale bar, 100 μm. **g**, Nucleotide contribution scores of the ME-2 motif
1100   implanting sequence (top) and in silico saturation mutagenesis assays (bottom). Each dot on the saturation mutagenesis plot
1101   represents a single mutation and its effect on the prediction score (⬚ axis).
1102
1103   **Extended Data Figure 5: Human enhancer design by in silico evolution**
1104   **a**, Distribution of GC-content in GC-adjusted random sequences (green) and human genomic regions (red). **b**, Average
1105   number of motif hits at each mutational step compared to genomic enhancers. **c**, Delta number of motifs in each mutational
1106   step. The TF-Modisco patterns and the most similar position weight matrices from the cisTarget motif database are shown
1107   at the top of each plot. The patterns that are upside-down are the ones contributing negatively to the model's prediction
1108   and they are destroyed by the model on each step. **d**, Bar plot showing the mean luciferase signal (log₂ fold-change over
1109   *Renilla*) in a MES melanoma line (MM047) of the synthetic MEL enhancers (generated by in silico sequence evolution),
1110   showing no activity compared to positive control genomic MES enhancers. The bar shows the mean (⬚ = 2 biological
1111   replicates). **e**, MM001 (left) and MM099 (right) ATAC-seq profiles of all integrated lentiviral EFS reporters. Red dashed lines
1112   indicate boundaries of the enhancer. **f**, MM001 ATAC-seq profile of 3 integrated EFS reporters: initial (top), evolved (middle)
1113   and post-evolution with repressive sites (bottom). Red lines mark enhancer boundaries. **g**, DeepMEL2 prediction score (left),
1114   luciferase activity levels in MM001 (middle) and correlation between prediction score and activity (right) for EFS-1 (top) and
1115   EFS-8 (bottom) sequences after incremental mutation steps. In **g**, the error bars show the standard error of the mean (⬚ = 3
1116   biological replicates)
1117
1118   **Extended Data Figure 6: Intermediate steps of in silico evolution and generation of repressor sites in human generated**
1119   **enhancers**
1120   Nucleotide contribution scores of EFS-4 at different mutational steps; 0 (random sequence), 3, 4, 7, 8, 12, 15, 15+Repressors.
1121   ZEB2 motif annotation is indicated with strong (s) or weak (w) motif instances.
1122
1123   **Extended Data Figure 7: Human enhancer design by in silico evolution**

1124 **a**, Prediction scores for the top 50 DNAse tracks for MEL EFS sequences. The four first DNAse tracks are: foreskin melanocyte
1125 male newborn, SK-MEL-5, foreskin melanocyte male newborn, SK-MEL-5. **b**, Enformer prediction tracks for three classes and
1126 ChromBPNet MM001 ATAC prediction tracks (right) for melanoma EFS-1 (top) and EFS-8 (bottom) sequences added in place
1127 of the *IRF4* enhancer. Top track: random sequence prediction score, other tracks: delta of mutated sequence prediction score
1128 vs random sequence prediction score. **c**, Enformer prediction tracks for three classes for melanoma EFS-4 sequences added
1129 10 kb upstream, 5 kb upstream or 17.5 kb downstream of the *IRF4* enhancer. Top track: random sequence prediction score,
1130 other tracks: delta of mutated sequence prediction score vs random sequence prediction score.
1131
1132 **Extended Data Figure 8: ZEB2 repression of in silico evolved MEL enhancers**
1133 **a**, Prediction scores for each mutational step and after the addition of repressor sites for 3 EFS sequences. **b**, Nucleotide
1134 contribution scores (DeepMEL2 MEL class) showing the creation of single or multiple repressor binding sites by single or
1135 double mutations in the EFS-4 sequence. **c-e**, In vivo enhancer activity of EFS-4 (**c**), EFS-1 (**d**), and EFS-8 (**e**) after the
1136 generation of repressor binding sites. ZEB2 motif annotation is indicated with strong (s) or weak (w) motif instances. The
1137 error bars in **c-e**, show the standard error of the mean ($n$ = 3 biological replicates).
1138
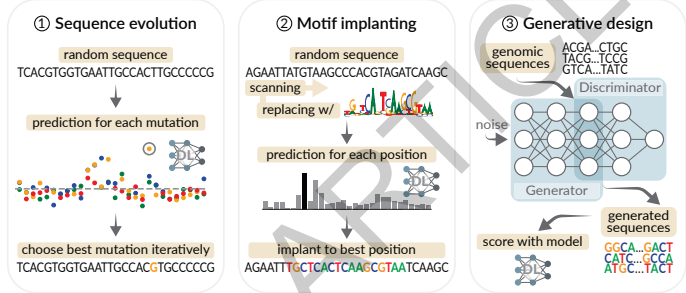1139 **Extended Data Figure 9: Human enhancer rescue**
1140 In the fly brain, we applied in silico sequence evolution to create enhancers from genomic regions with high scores that did
1141 not show chromatin accessibility and could consequently be considered as 'near-enhancer' sequences. We extended this
1142 approach to MEL enhancers. We started from a human sequence that has no MEL enhancer activity, but its homologous
1143 sequence in the dog genome is accessible and active as MEL enhancer. We used DeepMEL to introduce 4 mutations that
1144 restored the activator binding sites in the human sequence, resulting in a rescue of the activity, as measured by luciferase
1145 activity. **a**, Dot plot showing the mean luciferase signal (log$_2$ fold-change over *Renilla*) versus prediction score for the MEL
1146 class of the WT human and dog genomic sequences and the rescued human sequences. **b**, Nucleotide contribution scores of
1147 the dog, human-rescued, and human-WT sequences (top 3 rows) and in silico saturation mutagenesis assay of human-WT
1148 sequence (bottom). **c,** As a variation of this approach, we introduced two mutations in a weak MEL enhancer which resulted
1149 in a 10-fold increase in enhancer activity. Dot plot showing the mean luciferase signal (log$_2$ fold-change over *Renilla*) versus
1150 prediction score for the MEL class of the wild-type and enhanced enhancers. **d**, Nucleotide contribution scores of the wild-
1151 type (middle) and enhanced (top) enhancers and in silico saturation mutagenesis assay of wild-type enhancer (bottom). In **a**
1152 and **c**, the error bars show the standard error of the mean ($n$ = 3 biological replicates). In **a**, **c**, abbreviations are used for
1153 SOX10 (S), MITF (M), and TFAP2 (T). In **b**, **d**, each dot on the saturation mutagenesis plot represents a single mutation and its
1154 effect on the prediction score ($y$ axis). In **b**, **d**, the position of the mutations is shown with dashed lines and circles.
1155
1156 **Extended Data Figure 10: Human enhancer design by motif implantation**
1157 **a-b**, Bar plots show the mean luciferase signal (log$_2$ fold-change over *Renilla*) of the synthetic sequences, which were
1158 generated by motif implantation, tested in MM001 (**a**, MEL melanoma cell line, $n$ = 3 biological replicates) and MM047 (**b**,
1159 MES melanoma cell line, $n$ = 2 biological replicates). Values of 2 previously validated MES regions are displayed for MM047.
1160 The error bars in **a,** show the standard error of the mean. The bars in **b**, show the mean. **c**, Nucleotide contribution scores of
1161 the selected synthetic sequences in their initial form (first row), after adding SOX10, MITF, and TFAP2 motifs once (second
1162 row), after adding SOX10, MITF, and TFAP2 motifs twice (third row), weaker-motif version of the third row after replacing
1163 implanted motifs with weaker sites (fourth row), cut version of the second row where only the part with the binding sites
1164 were taken (fifth row, left), and minimal version of the second row where MITF and TFAP2 placed as close as possible to
1165 SOX10 (fifth row, right). The names of the motifs and their implantation order are indicated at the top. The position of the
1166 motifs is shown with dashed lines.

Enhancer Design Strategies

① Sequence evolution

random sequence
TCACGTGGTGAATTGCCACTTGCCCCCG

prediction for each mutation

choose best mutation iteratively
TCACGTGGTGAATTGCCAC**G**TGCCCCCG

② Motif implanting

random sequence
AGAATTATGTAAGCCCACGTAGATCAAGC
scanning

replacing w/

prediction for each position

implant to best position
AGAATTT**GCT**CACTCAA**GCGT**AATCAAGC

③ Generative design

genomic sequences
ACGA...CTGC
TACG...TCCG
GTCA...TATC

Discriminator

noise

Generator

score with model

generated sequences
GGCA....GACT
CATC...GCCA
ATGC....TACT

Enhancer Activity Measurements

① Genomic integration (transgenic flies)

| Designed enhancer | minP | nGFP (reporter gene) |

inject fly embryos

stain (α-GFP)
adult fly brains

② Episomal transfection (luciferase assays)

| Designed enhancer | minP | LUC (reporter gene) |

transfect cultured cells

measure
luciferase activity

Strong
Weak
Off

Luc. act.

**a**

Prediction score

Number of mutations

— Directed evolution
— Random drift

**c**

EFS-1
EFS-2
EFS-3
EFS-4
EFS-5
EFS-6

EFS-7
EFS-8
EFS-9
EFS-10
EFS-11
EFS-12
EFS-13

Number of mutations

**d**

EFS-1-15M — Weak
EFS-2-15M — Pos.
EFS-3-10M — Pos.
EFS-4-10M — Pos.
EFS-5-15M — Pos.
EFS-6-15M — Pos.
EFS-7-10M — Pos.
EFS-11-15M — Pos.
EFS-12-15M — Pos.
EFS-13-15M — Pos.

**b**

Random sequence

Repressor | Mamo (s) | AAGA (s) | Repressor

1:A>G   2:G>C 5:A>G 4:G>A 3:A>T 6:T>G   9:A>C 10:G>C 8:T>A 7:C>A

KC-directed mutations (EFS-4)

Ey (s)   Onecut (w)   Onecut (s)   Ey (w)   Mef2 (s)

**e**

Random sequence

Repressor   Repressor   Repressor   Repressor

12:A>C 14:A>T 15:G>T 9:T>C 10:G>T 2:G>T   11:T>C 8:C>A 7:G>A 6:G>T 13:T>C   1:A>T 5:A>G 3:T>G 4:T>A

PNG-directed mutations (PNG-2)

Activator   Activator   Activator Activator   Activator   Activator   Activator

**f**

Random Sequence — Negative / PNG mutations / KC mutations
KC-directed evolution (EFS-4) — KC Positive
PNG-directed evolution (PNG-2) — PNG Positive
PNG-1 — Pos.
PNG-4 — Pos.
PNG-5 — Pos.

**g**

Genomic sequence (FP1) — Neg. / KC mutations
KC-directed evolution (FP1-6M) — Pos.
FP2-6M — Pos.
FP3-10M — Pos.
FP4-6M — Pos.

**h** FP-4 (ChrX:9786300-9786800)

WT

GATC (s)   TAATTA (s)

6:T>A   1:G>A 4:A>C 2:C>A 5:T>C   3:A>G

6 mutations

Ey (s)   Onecut (s)   Sr (s)   Sr (s)   Mef2 (s) Ey (w) Ey (w)

**a** 0.015

**KC-directed 10 mutations (EFS-4)**

Ey (s)    Onecut (w)    Onecut (s)    Ey (w)    Mef2 (s)

○ Selected mutations

**Adding repressors by mutations (EFS-4-Repressed)**

TTTGGG (s) CAATTA (s)    AAGA (s)    GATC (s)    GATC (s)    AAGA (s)

● A
● C
● G
● T

130    150    170    190    210    230    250    270    290    310    330    350    370    390

**b**

EFS-4-10M    EFS-4-Repressed
Repressor - creating mutations    Pos.    Neg.

EFS-3-Repressed    EFS-7-Repressed
Neg.    Neg.

FP2-Repressed    FP4-Repressed
Neg.    Neg.

FP10-Repressed
Neg.

**c** AST CTX ENS PNG SUB T1 T2 T2a T3 T4 T5 α/β α'/β' γ

*amon*

**d** *amon* enhancer (WT)
KC-creating mutations    T Pos.

**g** *amon* enhancer (13 Mut.)
KC + T Pos.

**e** Pred. score
AST CTX ENS PNG SUB T1 T2 T2a T3 T4 T5 α/β α'/β' γ

**f** Pred. score
● γ-KC
● T4
Number of mutations

**h** x1000
# of regions
# of cell-types

**i** T1 pred. score
γ-KC pred. score

**j** AST CTX ENS PNG SUB T1 T2 T2a T3 T4 T5 α/β α'/β' γ

**k** *Pkc53e* enhancer (WT)
T1/2-destroying mutations    KC + T Pos.

**n** *Pkc53e* enhancer (9 Mut.)
KC Pos.

**l** Pred. score
AST CTX ENS PNG SUB T1 T2 T2a T3 T4 T5 α/β α'/β' γ

**m** Pred. score
● γ-KC
● T1
● T2
Number of mutations

**a** Pred. score
Best location
Random location
Random Sr Onecut Mef2 Ey M.E E.M E.M,O E.M,O,S Nuc-15

**b** # of sequences
5bp space
Location of Mef2 relative to Ey

**c** # of sequences
3bp space
Location of Onecut relative to Ey

**d** Pred. score
Random E E E E M M M O O R.

**e**
x0.001
Repressor          Repressor     Repressor        Repressor
**Random sequence**
60                                                                    400

Mef2   Ey   Onecut
**Ey, Mef2, Onecut added**

A
C
G
T
○ Selected mutations

**Repressed**
GATC (s)   GATC (s)   Repressor   GATC (s)   TCTT (s)
60  70  80  90  100  110  120  130  140  150  160  170  180  190  200  210  220

**f** Ey, Mef2, Onecut (500bp)
ME-1          Pos.

**g** Ey, Mef2, Onecut + Repressors
ME-1-repr.    Neg.

**h** Mef2   Ey   Onecut
4   10   5   17   3   6   4
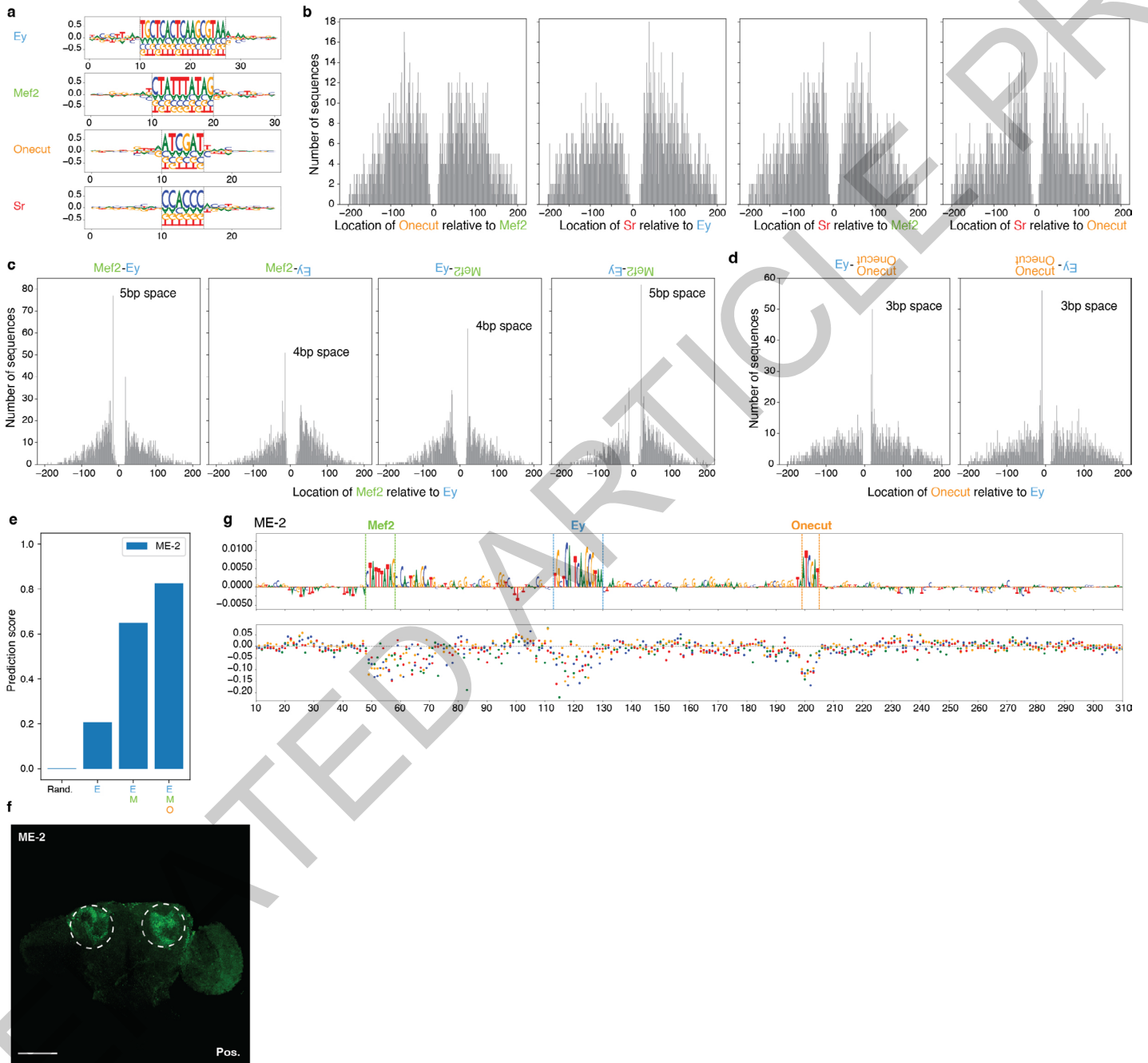49bp

**i** Ey, Mef2, Onecut (49bp)
ME-1-short    Pos.
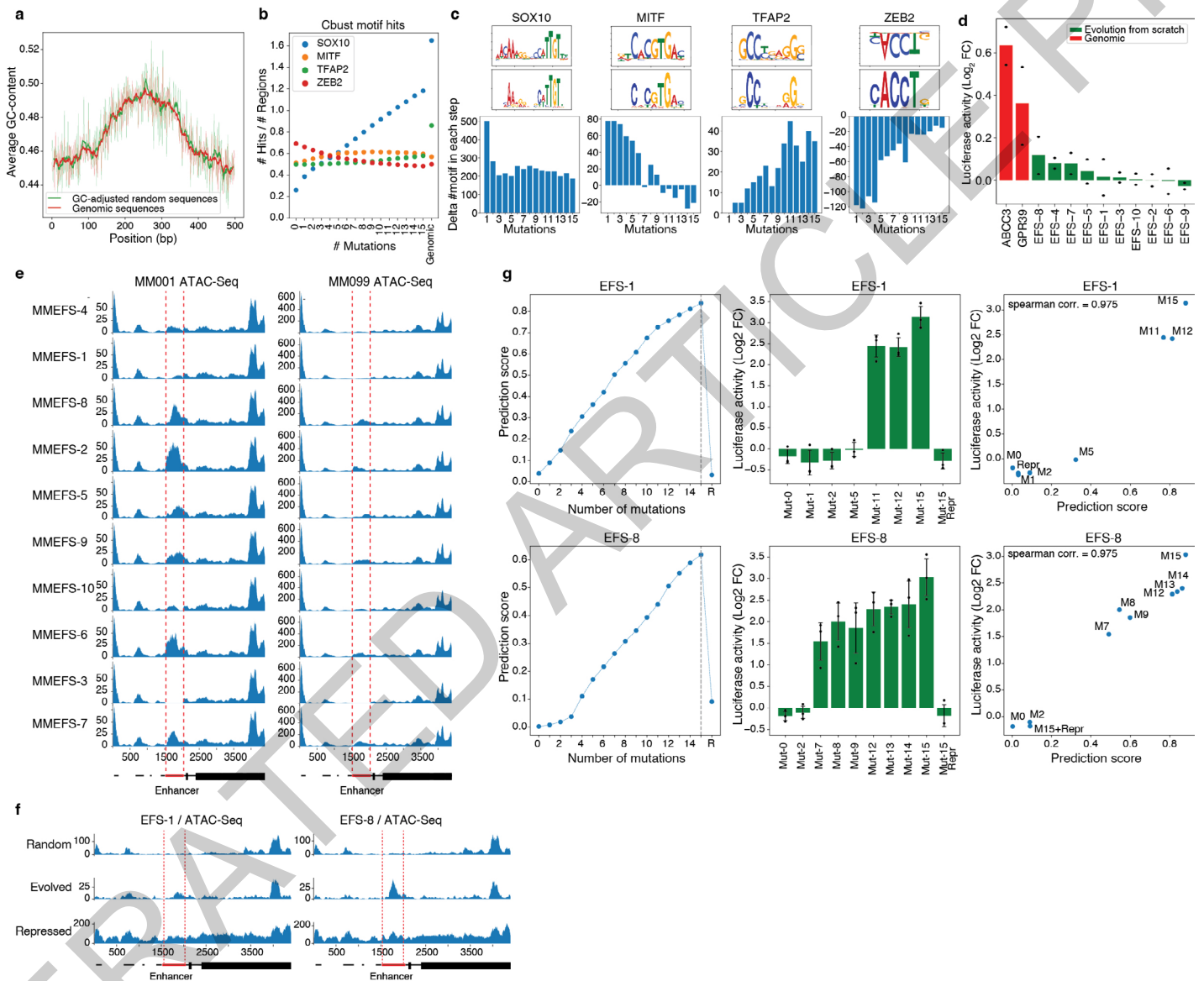
**Extended Data Fig. 1**

**Extended Data Fig. 2**

**Extended Data Fig. 3**

**Extended Data Fig. 4**

**Extended Data Fig. 5**

Mutation 0: Random sequence

Mutation 3: SOX

Mutation 4: ZEB2 -> MITF

Mutation 7: Weak SOX

Mutation 8: TFAP2

Mutation 12: SOX

Mutation 15: Weak SOX

Selected mutations

Mutation 15 + Repressors

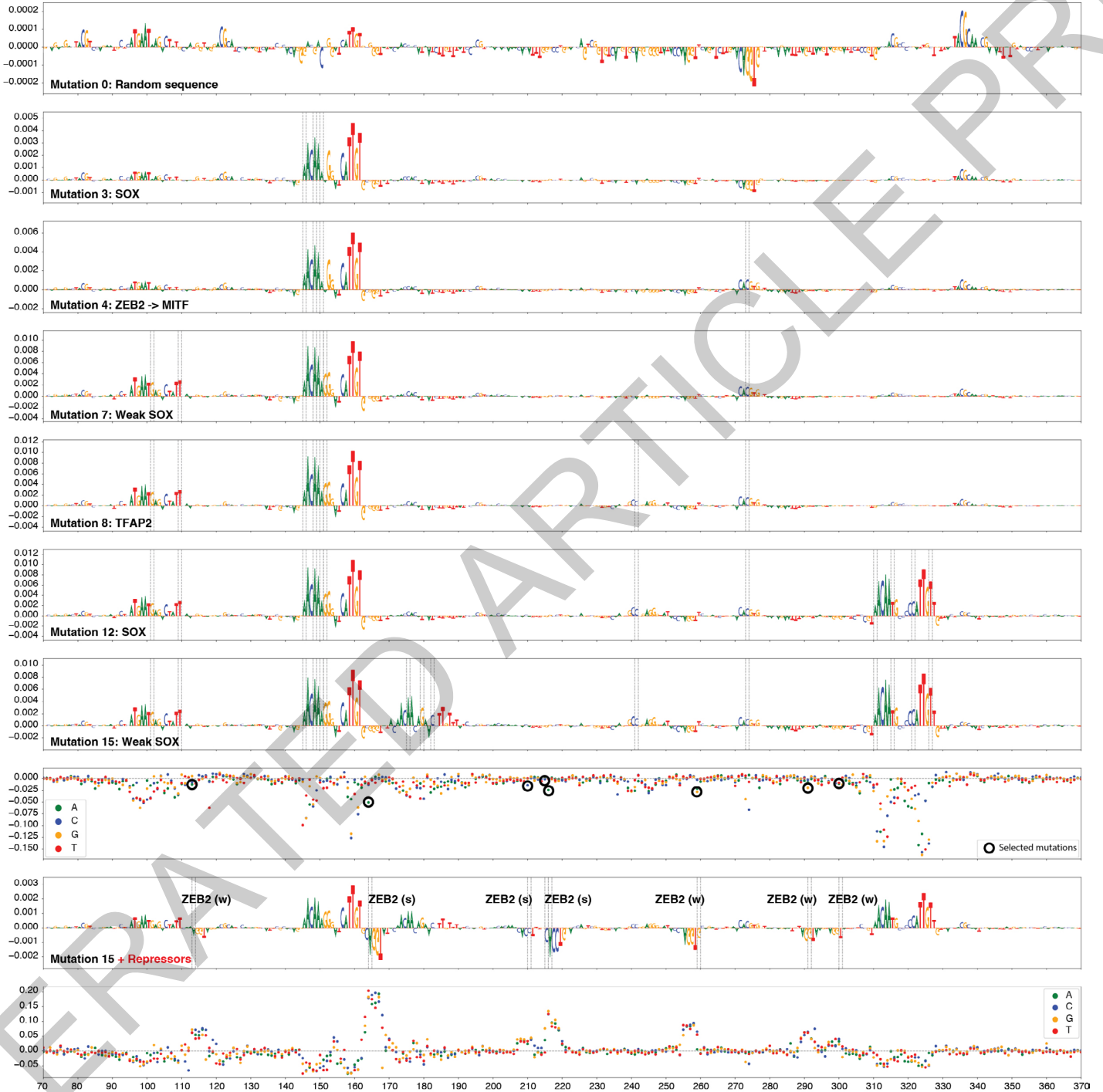ZEB2 (w)    ZEB2 (s)    ZEB2 (s)    ZEB2 (s)    ZEB2 (w)    ZEB2 (w)    ZEB2 (w)

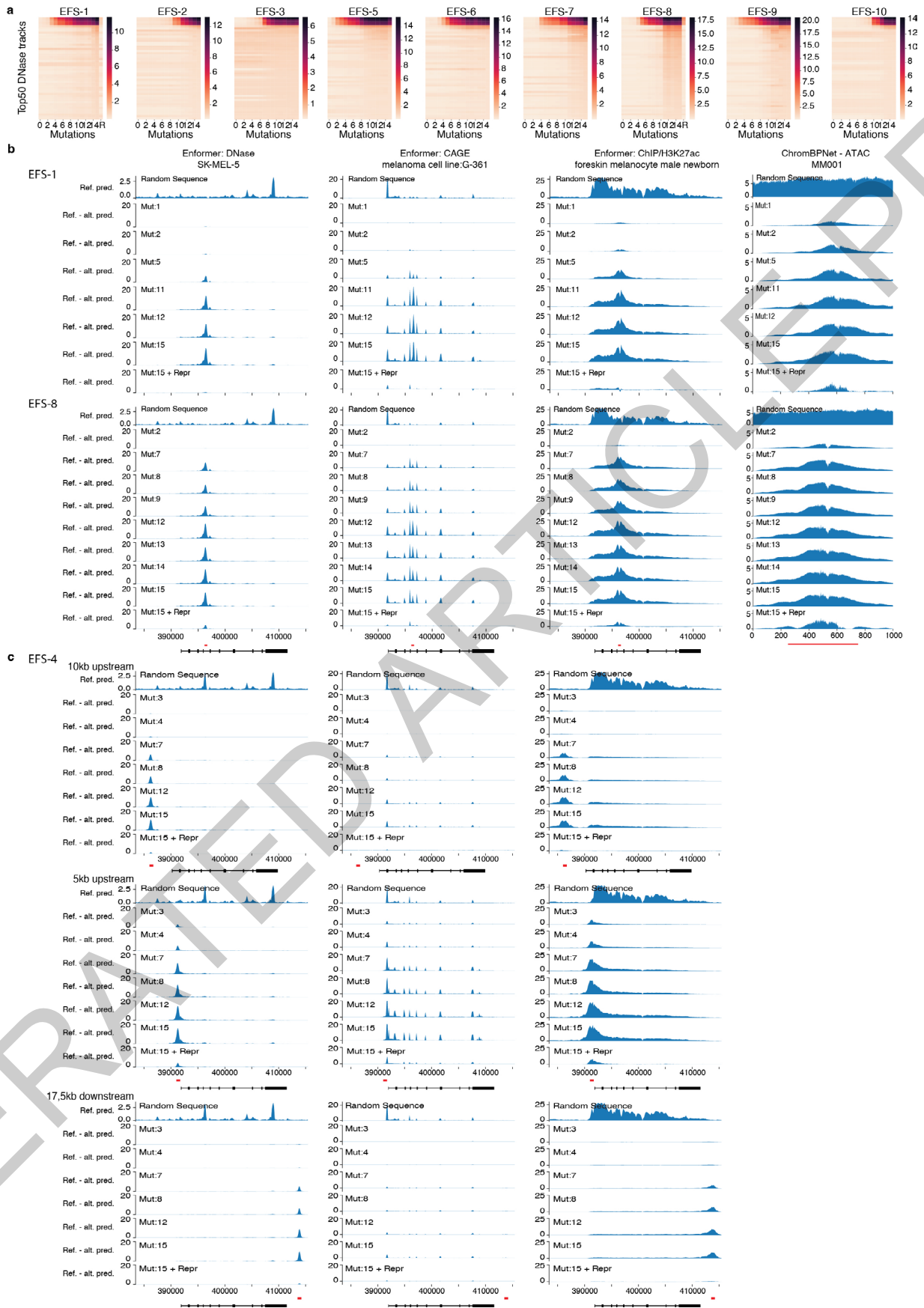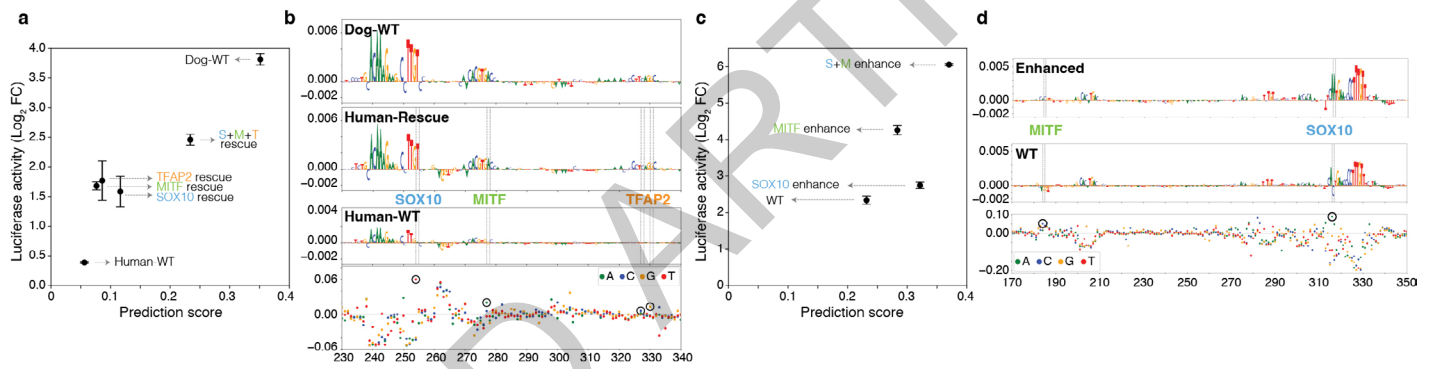**Extended Data Fig. 6**

Extended Data Fig. 7

**Extended Data Fig. 8**

**Extended Data Fig. 9**

Extended Data Fig. 10

# nature portfolio

| | |
|---|---|
| Corresponding author(s): | Stein Aerts |
| Last updated by author(s): | Nov 22, 2023 |

# Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our Editorial Policies and the Editorial Policy Checklist.

## Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

| n/a | Confirmed | |
|---|---|---|
| ☐ | ☒ | The exact sample size (*n*) for each experimental group/condition, given as a discrete number and unit of measurement |
| ☒ | ☐ | A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| ☒ | ☐ | The statistical test(s) used AND whether they are one- or two-sided<br>*Only common tests should be described solely by name; describe more complex techniques in the Methods section.* |
| ☒ | ☐ | A description of all covariates tested |
| ☒ | ☐ | A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| ☐ | ☒ | A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| ☒ | ☐ | For null hypothesis testing, the test statistic (e.g. *F*, *t*, *r*) with confidence intervals, effect sizes, degrees of freedom and *P* value noted<br>*Give P values as exact values whenever suitable.* |
| ☒ | ☐ | For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| ☒ | ☐ | For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| ☐ | ☒ | Estimates of effect sizes (e.g. Cohen's *d*, Pearson's *r*), indicating how they were calculated |

*Our web collection on statistics for biologists contains articles on many of the points above.*

## Software and code

Policy information about availability of computer code

| Data collection | Confocal images: ZEN blue 3.4.91 |
|---|---|
| Data analysis | Custom codes: https://doi.org/10.5281/zenodo.10184648<br><br>DL Python environment to use DeepMEL 1.0, DeepMEL2 1.0, and DeepFlyBrain 1.0:<br>python=3.7  tensorflow-gpu=1.15 numpy=1.19.5 matplotlib=3.1.1 shap=0.29.3 ipykernel=5.1.2 h5py=2.10.0<br><br>DL Python environment to train GAN models:<br>python=3.6  tensorflow-gpu=1.14.0 keras-gpu=2.2.4 numpy=1.16.2 matplotlib=3.1.1 shap=0.29.3 ipykernel=5.1.2<br><br>To perform motif analysis: TF-Modisco 0.5.5.4, Tomtom (MEME 5.5.1),  ClusterBuster 2022-04-21, BEDTools 2.30.0<br>To create higher-order background sequences: INCLUSive 3.2<br><br>To calculate statics: Scipy 1.6.0<br><br>To train ChromBPNet model: ChromBPNet 1.3-pre-release<br><br>ATAC-seq and ChIP-seq data analysis:<br>Demultiplexing with bcl2fastq 2.20<br>Adapter trimming with trimgalore 0.6.7<br>Mapping with bwa-mem2 2.2.1 |

Sorting with SAMtools 1.16.1
Deduplicating with SAMtools 1.16.1
Removing blacklist regions with SAMtools 1.16.1
Generating bigwig with deepTools 3.5.0
Peak calling with MACS2 2.1.2.1

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio guidelines for submitting code & software for further information.

## Data

Policy information about availability of data

All manuscripts must include a data availability statement. This statement should provide the following information, where applicable:
- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our policy

Cloned Drosophila and human sequences were provided as Supplementary Tables. DeepMEL, DeepMEL2, and DeepFlyBrain deep learning model files were obtained from Kipoi (http://kipoi.org/models/DeepMEL, https://kipoi.org/models/DeepFlyBrain) with Zenodo record ids 3592129, 4590308, and 5153337. The fasta files used to train GAN models and the trained GAN models are available on Zenodo at https://doi.org/10.5281/zenodo.6701504. Custom genomes (hg38 and dm6) generated in this study are available on Zenodo at https://doi.org/10.5281/zenodo.10184648. Chromatin accessibility values in Kenyon Cells in adult Drosophila brains were obtained from GSE16369739. In vitro saturation mutagenesis on IRF4 data was obtained from https://kircherlab.bihealth.org/satMutMPRA/. Chromatin accessibility of Drosophila and transduced melanoma lines and ZEB2 ChIP-seq data generated for this study have been submitted to the NCBI Gene Expression Omnibus (GEO, https://www.ncbi.nlm.nih.gov/geo/) under accession number GSE240003.

## Research involving human participants, their data, or biological material

Policy information about studies with human participants or human data. See also policy information about sex, gender (identity/presentation), and sexual orientation and race, ethnicity and racism.

| | |
|---|---|
| Reporting on sex and gender | No human research participants involved |
| Reporting on race, ethnicity, or other socially relevant groupings | No human research participants involved |
| Population characteristics | No human research participants involved |
| Recruitment | No human research participants involved |
| Ethics oversight | No human research participants involved |

Note that full information on the approval of the study protocol must also be provided in the manuscript.

# Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences  ☐ Behavioural & social sciences  ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

# Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

| | |
|---|---|
| Sample size | The number of synthetic enhancers that were tested using transgenic flies was determined to be minimally 6 per cell type and it was bounded by the feasibility of the transgenic animal generation experiments. In total, 68 transgenic flies were generated. The number of synthetic enhancers that were used with luciferase assays is determined to be minimally 10 per different category (in silico evolution, motif embedding, GAN, repressors, mutational steps). In total, 97 sequences were tested using luciferase assay. |
| Data exclusions | No data was excluded. |
| Replication | The same results were obtained from different replication experiments. Multiple brains (at least 10) were stained and imaged for the fly experiments. 3 biological replicates were performed for the main luciferase experiments. 2 biological replicates were performed for the negative control luciferase experiments. No biological replicates performed on ATAC-seq or ChIP-seq experiments. |
| Randomization | The initial random sequences (used for sequence evolution and motif implantation) were sampled from the sequence space that matches the |

| Randomization | GC content of the genomic sequences.<br>Flies fitting the gender(equal amount of male and female) and age (<10days) criteria were selected randomly for all experiments.<br>In this study, we didn't perform experiments that needed to be allocated into different groups. |
|---|---|
| Blinding | The investigators were blinded when performing cloning, transfection, antibody staining, and luciferase experiments by using enhancer IDs. |

# Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

## Materials & experimental systems

| n/a | Involved in the study |
|---|---|
| ☐ | ☒ Antibodies |
| ☐ | ☒ Eukaryotic cell lines |
| ☒ | ☐ Palaeontology and archaeology |
| ☐ | ☒ Animals and other organisms |
| ☒ | ☐ Clinical data |
| ☒ | ☐ Dual use research of concern |
| ☒ | ☐ Plants |

## Methods

| n/a | Involved in the study |
|---|---|
| ☐ | ☒ ChIP-seq |
| ☒ | ☐ Flow cytometry |
| ☒ | ☐ MRI-based neuroimaging |

## Antibodies

| Antibodies used | 1 - Rabbit polyclonal anti-GFP (1:1000 dilution); Life Technologies CAT# A-6455; RRID: AB_221570<br>2 - Donkey polyclonal anti-rabbit Alexa Fluor 488 (1:500 dilution); Life Technologies CAT# A-21206; RRID: AB_2535792<br>3 - Rabbit anti-ZEB2; Bethyl CAT# A302-473A (1mg/ml and we used 5 micrograms for ChIP)<br>4 - Mouse anti-Dachshund (1:250 dilution); DSHB; CAT# dac1-1<br>5 - Alexa Fluor 647 goat anti-mouse IgG (1:500 dilution); Invitrogen, CAT# A-21235 |
|---|---|
| Validation | 1- References provided, statement on manufacturer's website: "This Antibody was verified by Relative expression to ensure that the antibody binds to the antigen stated.". Selected references out of 238: PMID 36067320, 35142344, 34908527, 34644579, 33932333, 33846330, 33463521, 33174166, 33112231, 32640222.<br>2- References provided, no statement on manufacturer's website. Selected references out of 6277: PMID 36067320, 35142344, 34908527, 34644579, 33932333, 33846330, 33463521, 33174166, 33112231, 32640222.<br>3- Testing and references provided, we performed ChIP-seq using ZEB2 antibody and the most enriched motif was the ZEB2 motif. No statement on the manufacturer's website. References: PMID 33614228, 20515682<br>4- References provided, statement on manufacturer's website: "The antibody reproduces the pattern observed by in situ hybridization with a dac cDNA probe (unpublished observations) and an enhancer trap insert in dac.". References: PMID 7821215, 17868668, 32781577, 18430931, 25670791, 8756723, 9845371, 24142104, 22874913, 34409041, 34322481, 33982759, 32738261, 32781577, 32184260, 31453329.<br>5- References provided, statement on the manufacturer's website: "The antibody "was used with a concentration of 2μg/mL.". Selected references out of 1448: PMID 35297981, 35017509, 33570489, 32878938, 32649914, 33659324, 32579612, 32317641, 37332603, 36879821, 36355348, 36649336, 34459871, 34605405, 33689682. |

## Eukaryotic cell lines

Policy information about cell lines and Sex and Gender in Research

| Cell line source(s) | MM001, MM047, and MM099 were obtained from Prof. Dr. Ghanem Ghanem with a Material Transfer Agreement.<br>HEK293T was obtained from ATCC (CAT# CRL-3216). |
|---|---|
| Authentication | We have used MM001, MM047, and MM099 in previous studies (Verfaillie et al., Nature Communications, 2015; Wouters et al., Nature Cell Biology 2020; Minnoye et al., Genome Research, 2020; Kalender-Atak et al., Genome Research 2021). We authenticated the cell lines by tracking their morphology overtime and by checking their genomic profile and mutations (Verfaillie et al., Kalender-Atak et al.), transcriptomic profile (Wouters et al.), and epigenomic profile (Verfaillie et al., Wouters et al., Kalender-Atak et al.).<br>HEK293T cells were only used for lentivirus production in this study, and the final products were tested and confirmed by sequencing. No authentication was needed for this cell line. |
| Mycoplasma contamination | Cell lines were tested for mycoplasma contamination, and were found negative. |
| Commonly misidentified lines<br>(See ICLAC register) | No commonly misidentified cell lines were used in this study. |

# Animals and other research organisms

Policy information about [studies involving animals](); [ARRIVE guidelines]() recommended for reporting animal research, and [Sex and Gender in Research]()

| | |
|---|---|
| Laboratory animals | Transgenic Drosophila melanogaster strains were used in this study. Young adult flies (<10-days-old) were used when performing antibody stainings. |
| Wild animals | No wild animals were used. |
| Reporting on sex | Sexes were equally mixed when performing antibody staining on adult Drosophila melanogaster brains. |
| Field-collected samples | No field-collected samples were used. |
| Ethics oversight | No approval required. |

Note that full information on the approval of the study protocol must also be provided in the manuscript.

# Plants

| | |
|---|---|
| Seed stocks | No plants were used. |
| Novel plant genotypes | No plants were used. |
| Authentication | No plants were used. |

# ChIP-seq

## Data deposition

☒ Confirm that both raw and final processed data have been deposited in a public database such as [GEO]().

☒ Confirm that you have deposited or provided access to graph files (e.g. BED files) for the called peaks.

| | |
|---|---|
| Data access links *May remain private before publication.* | GSE240003 |
| Files in database submission | MM001_ZEB2_ChIP-seq<br>MM001_input_ChIP-seq |
| Genome browser session (e.g. [UCSC]()) | no longer applicable |

## Methodology

| | |
|---|---|
| Replicates | n=1 |
| Sequencing depth | ZEB2_ChIP-seq: 83410868<br>input_ChIP-seq: 168512695 |
| Antibodies | Rabbit anti-ZEB2; Bethyl CAT# A302-473A |
| Peak calling parameters | macs2 callpeak default parameters |
| Data quality | 31866 peaks are called with 5% FDR |
| Software | Demultiplexing with bcl2fastq 2.20<br>adapter trimming with trimgalore 0.6.7<br>mapping with bwa-mem2 2.2.1<br>sorting with SAMtools 1.16.1<br>deduplicating with SAMtools 1.16.1<br>removing blacklist regions with SAMtools 1.16.1<br>generating bigwig with deepTools 3.5.0 |

peak calling with MACS2 2.1.2.1