1

# Current sequence-based models capture gene expression determinants in promoters but mostly ignore distal enhancers

Alexander Karollus[1], Thomas Mauermeier[1], Julien Gagneur[1,2,3,4]

[1]Department of Informatics, Technical University of Munich, Garching, Germany, [2]Institute of Human Genetics, Technical University of Munich, Munich, Germany, [3]Institute of Computational Biology, Helmholtz Zentrum München, Neuherberg, Germany, [4]Munich Data Science Institute, Technical University of Munich, Garching, Germany

## Abstract

**Background: The largest sequence-based models of transcription control to date have been obtained by predicting genome-wide gene regulatory assays across the human genome. This setting is fundamentally correlative, as those models are exposed during training solely to the sequence variation between human genes that arose through evolution, questioning the extent to which those models capture genuine causal signals.**

**Results: Here we confront predictions of state-of-the-art models of transcription regulation against data from two large-scale observational studies and five deep perturbation assays. The most advanced of these sequence-based models, Enformer, by and large captures causal determinants of human promoters. However, models fail to capture the causal effects of enhancers on expression, notably in medium to long distances and particularly for highly expressed promoters. More generally, the predicted impact of distal elements on gene expression predictions is small and the ability to correctly integrate long-range information is significantly more limited than the receptive fields of the models suggest. This is likely caused by the escalating class imbalance between actual and candidate regulatory elements as distance increases.**

**Conclusions: Our results suggest that sequence-based models have advanced to the point that in-silico study of promoter regions and promoter variants can provide meaningful insights and we provide practical guidance on how to use them. Moreover, we foresee that it will require significantly more and particularly new kinds of data to train models accurately accounting for distal elements.**

Correspondence to: AK alexander.karollus@tum.de or JG gagneur@in.tum.de

Keywords: Deep Learning, Transcription, Variant Effect, Gene Expression, Promoter, Enhancer

2

# Introduction

Regulatory regions in the genome encode instructions determining gene product abundance in response to developmental and environmental cues encode instructions determining gene product abundance in response to developmental and environmental cues. Inherited or acquired genetic alterations in these regulatory regions can result in the dysregulation of gene expression, which ultimately can cause a variety of diseases. Accordingly, models which can reliably predict gene expression directly from sequence would not only be of scientific interest but could potentially find many uses in the design of personalized diagnoses and treatments.

Currently, no sequence-based model is capable of holistically accounting for all stages of gene expression from transcription initiation to protein degradation and can thus predict the abundance of each processed protein isoform in any given cellular context. However, deep learning models have recently been proposed which - at least in theory - can predict measures of RNA abundance directly from arbitrary input sequences centered on a gene of interest for large sets of human cell types and tissues. The focus of this study is to understand how successful these models are. For the sake of simplicity, we will follow the convention in the literature and use RNA abundance and gene expression synonymously, even if the former is only an imperfect proxy of the latter [1,2].

We will study Xpresso [3], a cell-type agnostic model which predicts gene expression from a small sequence window around the transcription start site (TSS), Basenji1 [4] and 2 [5], deep convolutional models which were trained on many genome-wide assays and cell lines of the ENCODE project [6,7] and use about 40 kilobases (kb) of context, and Expecto [8], a linear model trained on top of the deep convolutional model DeepSea [9], itself trained to predict ENCODE genome-wide assays. Moreover, we will investigate the performance of the largest model trained to date, Enformer [10], a deep transformer [11] with ~250 million trainable parameters - an order of magnitude more than preceding models. As input, Enformer gets a 196 kb long sequence and predicts the value of 5,313 different ENCODE tracks in bins of 128bp. These tracks include chromatin-immunoprecipitation signal for hundreds of transcription factors, DNase footprinting (DNase), which measures genome accessibility, and cap-analysis of gene expression (CAGE) measurements for hundreds of cellular contexts (defined as combinations of cell lines and treatments). For this study, the CAGE predictions are the most relevant, as these provide the sought-after measure of RNA abundance.

As Avsec et al. [10] convincingly show in their paper, Enformer provides unparalleled performance when predicting CAGE signal of held-out test genes. At least in some cell types, the model is close to experimental accuracy. However, such aggregate measures of performance carry significant caveats. The main issue is that all the models we named were trained using the sole genetic diversity available across the human and mouse genomes. This setting is fundamentally correlative, as genomic sequences of any organism are not a random sample from the space of possible sequences but rather have been selected and co-evolved over millions of years of evolution. Thus, it is unclear to what extent the models have learned and employ causal principles, rather than mere correlations, to make their predictions. If so, predictions may become very misleading when applied in a medical diagnostic context to interpret rare germline or somatic variants, or if the model is used to generate new mechanistic hypotheses.

Additionally, improved performance of a model on aggregate measures such as total explained variance does not provide insights into the reasons for the improvement. Enformer has more parameters than previous models, but also a much wider receptive field - i.e. the

2

3

89 size of the sequence window it can integrate over to predict expression at a particular
90 location. Is the wider receptive field the deciding factor for its improved performance?
91 Understanding the actual source of performance improvements will help to design better
92 models.

93 To address these questions, we conducted in-silico reproductions of two large-scale
94 observational assays which measured gene expression in different tissues and stages of
95 development and five deep perturbation assays, including designed massively parallel
96 reporter assays for promoters and enhancers, CRISPRi enhancer-knockdown and saturation
97 mutagenesis experiments. Moreover, we conducted two in-silico perturbation studies. In this
98 way, we can probe in a targeted fashion to what extent a model actually makes use of
99 particular regulatory elements in its expression predictions and whether it does so in a way
100 consistent with experiments.

4

# Results

## Evaluating Deep Models through in-silico reproduction of experiments

| Variable Sequence | Fixed Sequence | In-silico Experimental Design | In-vivo Data | Cell Type |
|---|---|---|---|---|
| All | None | 196 kb | GTEx RNA-seq (47 tissues) | Human Tissues (Matched with linear Model) |
| | | | Cardoso-Moreira et al. (RNA-Seq, 7 tissues, 23 development stages) | |
| Promoter | Enhancers, General Sequence Context | | Designed Reporter Assays (Weingarten-Gabbay et al. Bergman et al.) | K562 |
| Single Nucleotide Variant | Everything Else | vs. | Kircher et al. Saturation Mutagenesis in plasmid for selected loci | Different Cell Lines (manually matched) |
| Enhancer | Promoter, General Sequence Context | | Designed Reporter Assay (Bergman et al.) | K562 |
| | | vs. | CRISPR Enhancer Knockdown (Fulco et al., Gasperini et al.) | |
| eQTL Variant | Everything Else (up to variants in linkage) | vs. | GTEx eQTL (finemapped with SuSiE) | Human Tissues (Matched with linear Model) |
| Enhancers, General Sequence Context | Promoter | | Trip-seq (Hong et al.) | K562 |

**Legend:** Promoter · Enhancer — Sequence Context · Variant · Location of Prediction · NNNNN Masking

***Figure 1: Overview of our in-silico experiments.*** To assess the generalization power of the models we performed analyses (rows), in which certain sequences varied (first column) while others were kept fixed (second column).

The overview figure (Fig 1) summarizes the different datasets we used in our study, which regulatory element(s) each dataset focuses on, and how we reproduced the experiment in-silico. Generally, replicating an experiment with a sequence-based model is not straightforward and requires three preparatory steps:

1. As many experiments involve some modification of the endogenous genome, we must construct in silico the correct sequences.
2. Most datasets do not report CAGE tracks but alternative measures of gene

4

114 expression including gene-level RNA-sequencing read counts and reporter
115 fluorescence. Hence, with the exception of Xpresso, which gives gene-level
116 predictions, we need to decide for which transcription start site predictions should be
117 made, i.e. from which bins to record the CAGE predictions.
118 3. Many experiments are conducted in cell types and tissues that do not exactly match
119 those of ENCODE. Some matching or mapping between those cellular contexts must
120 therefore be done.

121 The first step is the most intricate one. Several experiments we analyzed integrated
122 sequences into the endogenous genome. Usually, these sequences consist of the regulatory
123 element of interest, a reporter, post-transcriptional elements (e.g. chimeric introns and viral
124 polyadenylation sites), and technical elements which facilitate the genomic integration (e.g.
125 retrotransposon long terminal repeats) and sequencing. Many of these technical elements
126 are highly artificial and thus unlike anything the deep models we consider will have seen in
127 their training data. For this reason, we performed each replication twice, once with a faithful
128 reproduction of the full insert and once with a minimal insert, consisting only of the regulatory
129 element of interest and the reporter. Interestingly, almost always this minimal insert led to
130 predictions that better correlated with the experiment. Therefore, We decided to report only
131 the results from these minimal inserts. When possible, we avoided plasmid-based assays,
132 as there is no way to represent a circular chromosome in the models we consider. For lack
133 of alternatives, we made two exceptions, namely the Bergman et al. [12] promoter-enhancer
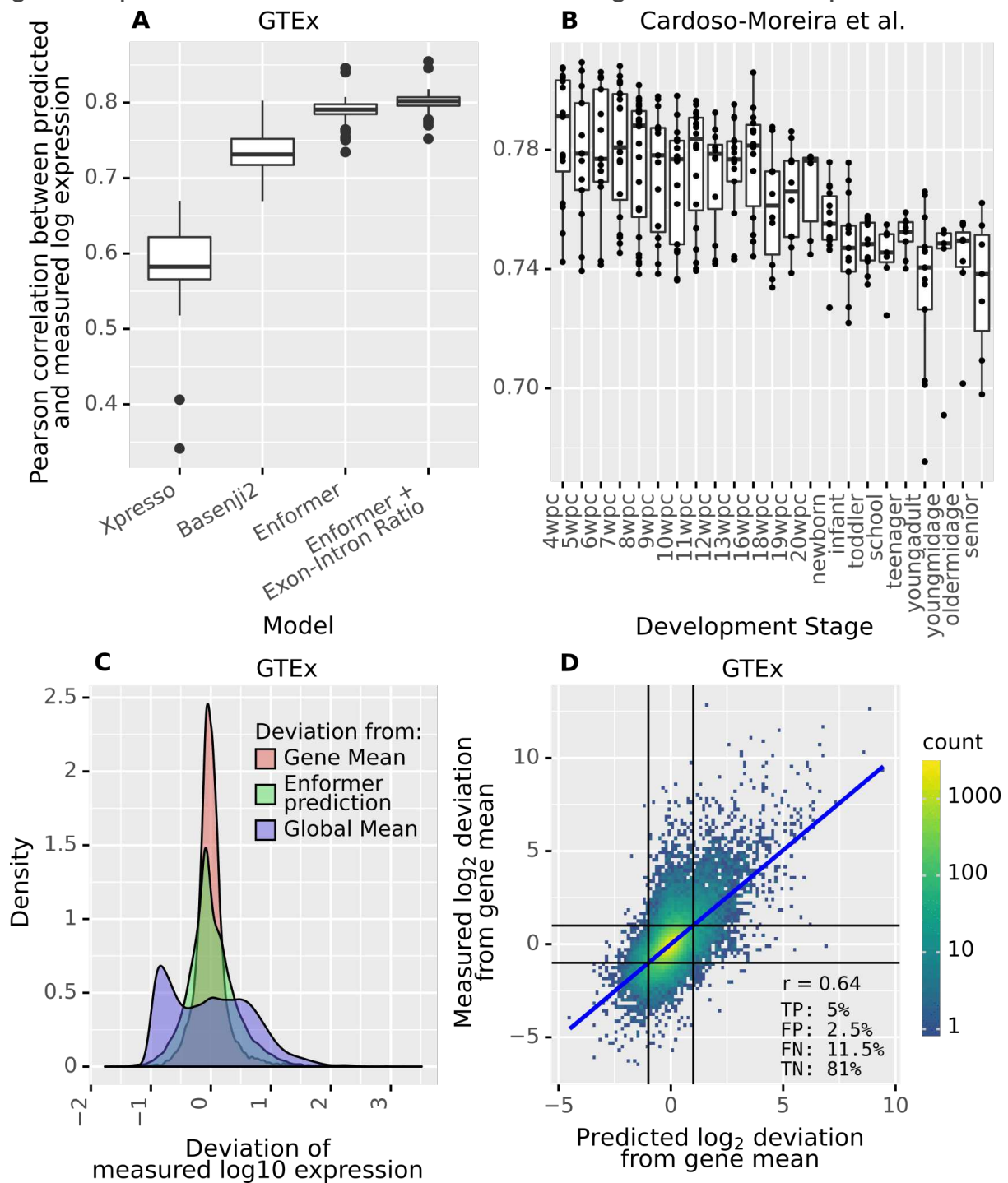134 compatibility study and the Kircher et al. [13] saturation mutagenesis study.

135 The second step is somewhat easier. For RNA-seq datasets (GTEx [14], Cardoso-Moreira et
136 al. [15], GTEx eQTL [16]), we considered CAGE prediction for the TSS of the Ensembl [17]
137 canonical transcript. This might not be the correct transcript in some tissues, but
138 considerably reduces the number of predictions to be made. The only exception to this rule
139 is the CRISPRi enhancer-knockout study, where we used the same TSS sites as Avsec et
140 al. [10] If we did not know the location of the TSS, but we knew where the core promoter was
141 located (e.g. because the experiment involves integrating promoters at arbitrary locations,
142 then we took the prediction around the promoter midpoint. For the Kircher et al. saturation
143 mutagenesis study, we took the prediction at the variant as done previously ([10]).

144 For the third step, all but three of the experiments we considered were done in K562 cells. In
145 these cases, we simply took the K562 CAGE track. For the studies done using RNA-seq of
146 human tissue biosamples, we fitted for each tissue a simple linear model (ridge regression)
147 predicting RNA abundance from all ENCODE CAGE tracks. We held out the
148 Enformer/Basenji2 test-set genes while fitting these regressions and evaluated only on this
149 held-out data. These fitted regressions were also used to analyze GTEx eQTL. For the
150 Kircher et al. saturation mutagenesis data, we used a manual matching procedure.

151 Once these preparatory steps are completed, we can run the sequence-based model on the
152 constructed sequences and collect the relevant predictions. Here it is important to note that
153 many of these models are sensitive to small changes in the input (e.g. small shifts),
154 particularly if regulatory elements fall directly on bin boundaries. To mitigate this, we
155 computed each prediction six times - for both strands and with small offsets respectively -
156 and took the average. We also always summed predictions over three neighboring bins. The
157 same technique was used by Avsec et al. [10], presumably for the same reason.

158

6

## Deep Models, particularly Enformer, provide very accurate predictions of gene expression in human tissues and during human development



**Figure 2: Enformer provides effective gene expression prediction for endogenous genes. A)** *Pearson correlation between predicted and measured log-transformed expression on GTEx tissues for different models. Enformer can predict endogenous RNA abundance, as measured in adult tissues (GTEx [14]), better than previous models. Adding the exon-intron ratio, a (weak) proxy of RNA half-life, as an additional predictor slightly improves performance.* **B)** *Same as A) for Enformer predictions on developmental samples (Cardoso-Moreira et al. [15] dataset). Enformer predicts endogenous gene expression very well overall*

6

169 *yet somewhat worse for later stages of development.* **C)** *Distribution of deviations of GTEx*
170 *measured log expression values from (1) the global mean (across genes and tissues, blue),*
171 *(2) the gene mean (across tissues, red), and (3) the Enformer prediction (green). The first*
172 *indicates overall variation in expression, the second indicates between-tissue variation and*
173 *the third indicates the magnitude of errors of Enformer. Enformer accuracy is sufficient to*
174 *explain much of the between-gene variation but not for the variation of genes between*
175 *tissues.* **D)** *Measured between-tissue deviations of gene expression against prediction.*
176 *Enformer predicts large between-tissue changes in expression reasonably well on average,*
177 *but there is significant room for improvement. The numbers indicate the percentages of true*
178 *positives (TP), false positives (FP), false negatives (FN), and true negatives (TN) when*
179 *predicting 2-fold changes (black lines).*

180 Avsec et al. [10] already provide ample evidence that Enformer can predict the expression of
181 endogenous genes very well. However, these validations were mostly done on ENCODE,
182 which is highly enriched for cancer cell lines that may provide an imperfect proxy of in-vivo
183 expression. Thus, to provide a slightly more complete picture, we benchmarked the model
184 on two additional RNA-Seq datasets: GTEx [14], which measured gene expression in 49
185 human tissues, and Cardoso-Moreira et al. [15] which measured gene expression in 7
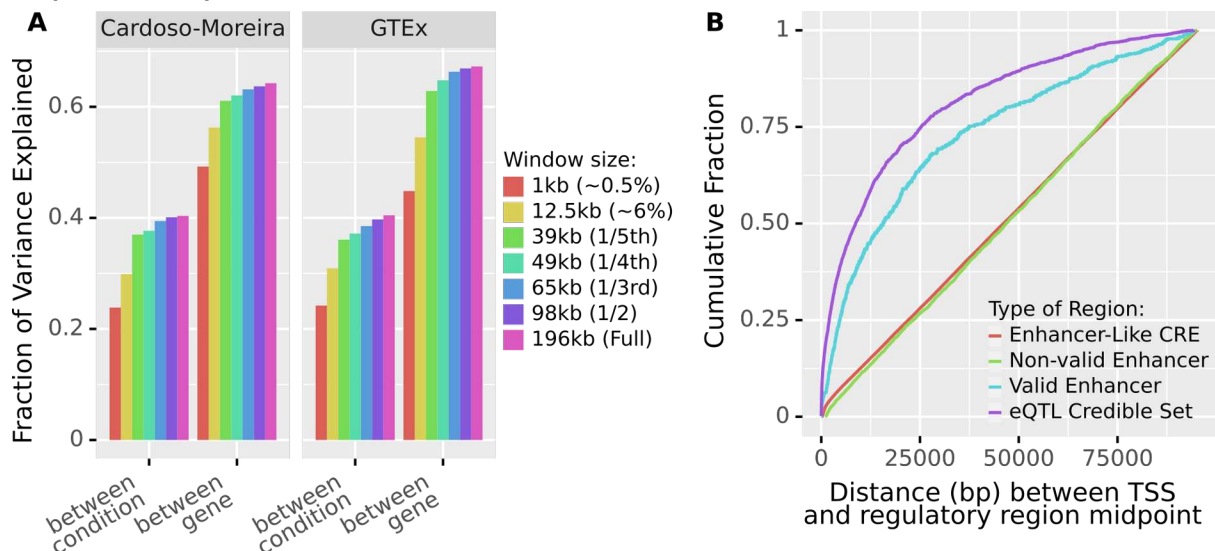186 tissues for 23 stages of development - from 4 weeks post-conception to senescence.

187 In both datasets, Enformer performs very well at predicting the between-gene variation (Fig
188 2A). For the median GTEx tissue, predictions and measurements correlate with r = 0.79.
189 Moreover, Enformer consistently outperforms both Xpresso and Basenji2. Adding the tissue-
190 specific exon-intron ratio of each gene - a (weak) proxy of RNA stability - as an additional
191 predictor improves the performance on GTEx. This suggests that Enformer predictions could
192 be augmented using a model that better captures post-transcriptional regulation. The median
193 correlation in the Cardoso-Moreira et al. dataset is very similar (r = 0.77). Notably,
194 performance is uneven across tissues and also degrades in later stages of development (Fig
195 2B, Fig S1). This might be due to varying data quality, the influence of environmental factors,
196 because some tissues and stages of development intrinsically feature more complex
197 regulation (e.g. testis [18]) or because they are not properly covered by the ENCODE cell
198 lines Enformer was trained on.

199 The above-mentioned correlations are computed across all genes, which span a large
200 dynamic range. Specifically, the mean absolute deviation from the grand mean of log RNA
201 abundance (across expressed genes and tissues) in the GTEx data is ~4-fold (Fig 2C). In
202 contrast, the mean absolute deviation to the mean per gene across tissue is only ~1.5-fold.
203 As Enformer has a mean-absolute error of about 2-fold (Fig 2C), it naturally struggles to
204 predict smaller differences between tissues. Nevertheless, the correlation between
205 measured and predicted log fold changes of genes between GTEx tissues is still remarkable
206 (r = 0.64, Fig 2D, Basenji2: r = 0.54). This performance translates to a decent precision of
207 66% for a recall of 30% at predicting 2-fold changes between GTEx tissues.

208 Overall, we find that Enformer can predict endogenous RNA abundance very well and
209 consistently outperforms previous models. This being said, when we compare the
210 expression of different genes, we are comparing highly dissimilar regulatory sequences.
211 These genes will generally have different promoters, different GC-content, different
212 enhancers, and will be located in different chromosomal contexts. Thus, these aggregated
213 results do not tell us which features of the sequence the model uses to make its predictions.

214 ## Most of the receptive field has a very minor impact on Enformer gene
215 ## expression predictions



216

***Figure 3: Enformer has very similar predictive power even if we severely restrict its input window, partially because most strong regulators are proximal. A)** Fraction of variance in log-transformed expression, both between conditions and between genes, which Enformer can explain given varying amounts of sequence context. Values computed on Enformer held-out data. Most of the signal comes from the sequence immediately around the TSS, with the distal two-thirds of the input window contributing very little. **B)** Distribution of the distance within 98 kb of TSS of bona fide regulatory elements (eQTL, purple, and CRISPRi validated enhancers, blue) and candidate elements (ENCODE CRE with enhancer-like signal, red) and CRISPRi tested but not validated enhancers (green). Most bona fide regulatory elements lie close to their target gene whereas candidate elements are uniformly distributed. We only consider elements within 98kb of a TSS, i.e. within the Enformer receptive field.*

229 Due to its wide receptive field (196kb), Enformer can theoretically account for the impact of
230 regulatory elements up to a distance of 98 kb on either side of a TSS. Given our
231 observations in the previous section, we wondered to what extent Enformer relies on these
232 distal elements to correctly predict gene expression.

233 To this end, we created sequence windows of varying sizes centered on the TSS of genes
234 by masking distal parts of the endogenous sequence (with "N" nucleotides). We then let
235 Enformer predict the CAGE gene expression at these TSS for each window size. In this way,
236 we can evaluate how the predictive power of Enformer on the GTEx [14] and the Cardoso-
237 Moreira et al. [15] data changes when it can no longer use distal elements to inform its
238 predictions.

239 Reassuringly, expanding the sequence window consistently improves the gene expression
240 predictions (Fig 3A). However, Enformer already achieves substantial explained variance
241 (about two-thirds of what is achieved with the full sequence) with only a tiny sequence
242 window of 1001bp (~0.5% of the total receptive field) around the TSS. Moreover, we face
243 strong diminishing returns when adding additional sequence context. While expanding from
244 1kb to ~40kb yields substantial improvements, the entire distal two-thirds of the sequence
245 only added ~1% of the total explained variation. Importantly, the improvement of Enformer
246 against its predecessor Basenji2 appears consistent across all distances (Fig S2), with

247  substantially better predictions even when Enformer can only access the same amount of
248  sequence than Basenji2 (~40kb), particularly for between-condition predictions.

249  In conclusion, Enformer extracts most of the signal to predict gene expression from promoter
250  and promoter-proximal sequences, with sequences further than 30kb from the TSS having a
251  negligible impact on its overall explained variance. Surprisingly, distal sequences also have
252  little impact on predicting between-condition variation, where we would have expected
253  relatively more influence from distal elements. What is not clear from this analysis is whether
254  the seeming irrelevance of distal elements reflects biological reality.

## Most known strong regulators are located close to their target genes, leading to an extreme class imbalance at higher distances

257  Our previous results indicate that Enformer considers most distal sequences to have a
258  negligible impact on gene expression at the TSS. We sought to examine the biological
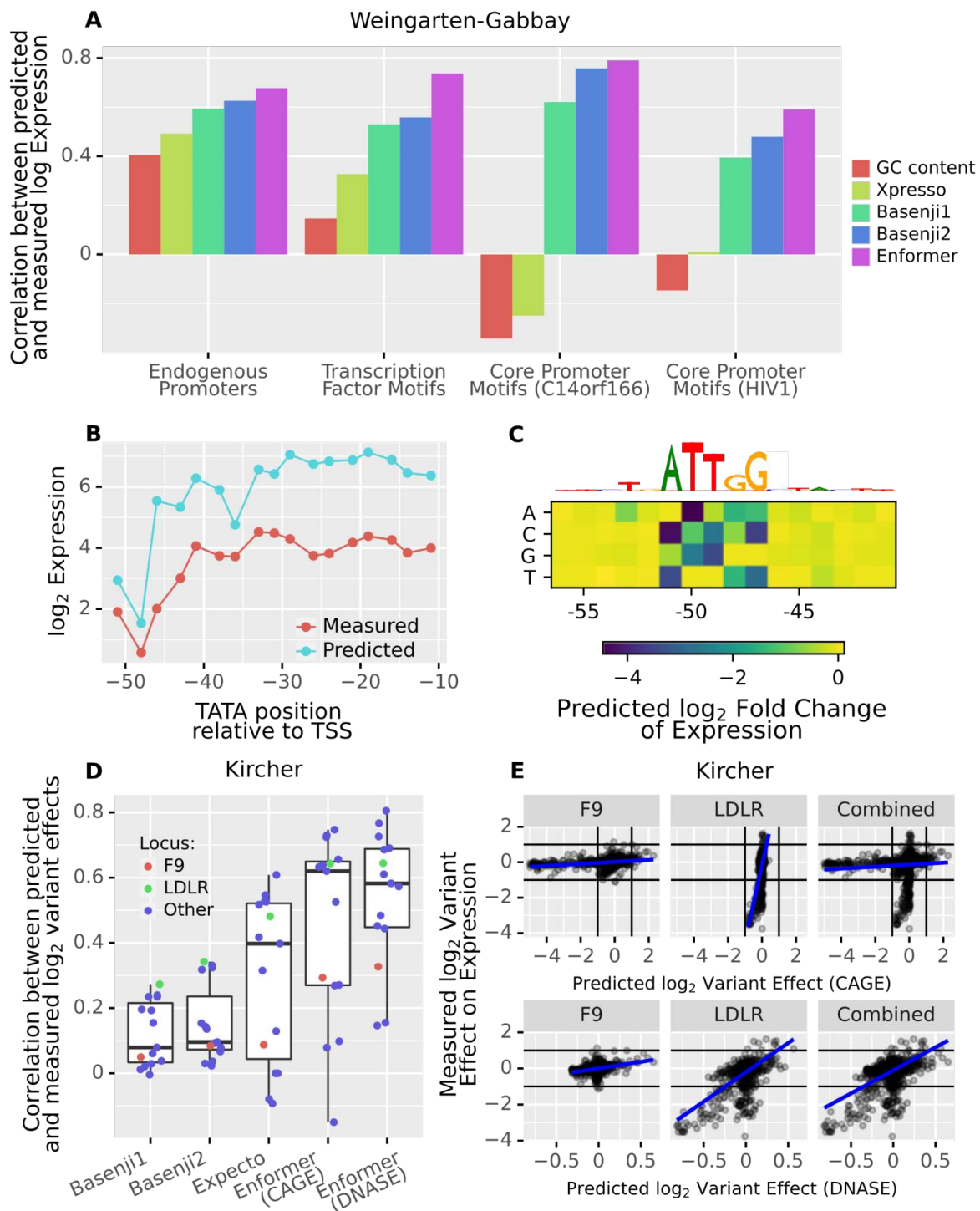259  plausibility of this.

260  For this purpose, we computed the distance distribution of *bona fide* regulatory elements,
261  specifically, eQTL [16] (which pass a series of filters, Methods) and CRISPRi validated
262  enhancers [19,20]. Most of these *bona fide* regulatory elements are located close to the TSS
263  of their target gene, in striking contrast to CRISPRi-tested but not validated candidate
264  enhancers (Fig 3B). Because of the limited power of the underlying assays, we cannot
265  necessarily conclude from this that most regulators of a gene will be proximal, but we can
266  conclude that the majority of strong regulators (i.e. those with large, individual effects) will be
267  proximal. These observations, consistent with previous studies [19,21,22], suggest that there
268  is a biological basis for Enformer to attribute more importance to local sequences.

269  However, the distribution of bona fide regulatory elements seems to not be as imbalanced as
270  the one implied by Enformer explained variance. Specifically, we see that the distal two-
271  thirds of the sequence still contains ~20% of known eQTL and ~25% of validated enhancers.
272  Interestingly, candidate regulatory elements from ENCODE (ENCODE SCREEN cis-
273  Regulatory Elements with enhancer-like signature [23]) show a uniform distribution within
274  100 kb of TSS (Fig 3B). Thus, for a typical gene, we can usually find similar amounts of
275  „enhancer-like" sequences at every distance to the TSS. The result is that, as we move
276  further from the TSS, the ratio of relevant regulatory elements to candidates (i.e. sequences
277  which look reasonably similar to regulatory elements) will necessarily become very
278  unfavorable very quickly. Thus, at higher distances, every long-range model will face an
279  extreme class imbalance. Perhaps this class imbalance and the difficulty to distinguish
280  enhancers targeting a given TSS from other enhancer-like elements is the cause of the
281  apparent under-usage of distal regulatory sequences.

282  We next examined perturbation experiments in detail to determine whether individual
283  regulatory elements contribute to Enformer predictions in a causal manner.

284 Deep Models, particularly Enformer, correctly predict promoter strength
285 and the impact of many promoter modifications



286

287 *Figure 4: Enformer accurately predicts genetic perturbations of promoters. A)*
288 *Correlation between model predictions and measurements across synthetic promoters of the*
289 *Weingarten-Gabbay et al. [24] parallel reporter assay. Enformer outperforms preceding*
290 *models, notably in targeted perturbation experiments (Transcription factor motif ablation and*
291 *core promoter motif perturbations). B) Enformer can detect and often correctly interpret the*

292 *impact on expression of subtle locational shifts of the TATA box in the RPLP0 promoter*
293 *background.* **C**) *In-silico mutagenesis analysis suggests that the large drop of expression*
294 *around position -50 likely from panel B is due to the disruption of a CAT-box at this location,*
295 *rather than the positional preference of the TATA box.* **D**) *Pearson correlation between*
296 *predicted and observed variant effects for the different loci tested by Kircher et al. [13]*
297 *Enformer outperforms other models for most loci.* **E**) *Measured vs. predicted variant effects*
298 *for two loci (F9, LDLR) individually and then both loci combined for CAGE (top) and DNase*
299 *(bottom). The CAGE predictions appear to be miscalibrated across those two loci.*

300 Most of the signal used by Enformer to predict gene expression seems to derive from the
301 promoter and promoter-proximal elements. Thus we first sought to assess to what extent
302 Enformer - and previous models - can predict the causal determinants of human promoters.

303 For this purpose, we performed an in-silico replication of the MPRA study by Weingarten-
304 Gabbay et al. [24] who measured the impact of different endogenous and synthetic
305 promoters on the expression of a reporter in K562 cells. In this study, 2274 164bp long
306 sequence fragments corresponding to known human promoters and pre-initiation complex
307 binding sites were inserted together with a fluorescence reporter at the AAVS1 viral
308 integration site.

309 Generally, Enformer predicts the relative effects of those promoters very well (the Pearson
310 correlation between measured and predicted log expression values is 0.68, Fig 4A).
311 Moreover, Enformer outperforms all other models on this task.

312 Weingarten-Gabbay et al. additionally constructed synthetic sequences to measure how
313 individual promoter elements affect its overall strength. In one experiment, they inserted 133
314 different TF binding sites in two different backgrounds. We find that Enformer predicts the
315 expression impact of these different TF motifs very well in both backgrounds (r = 0.75 and
316 0.73 respectively) and outperforms all other models (Fig 4A). We conclude that Enformer not
317 only recognizes these motifs but also generally correctly determines whether they act as
318 repressors, activators, or neither in K562 cells.

319 In yet another experiment, Weingarten-Gabbay et al. tested many combinations of six known
320 core promoter motifs (including the TATA box and the initiator) in five backgrounds.
321 Enformer predicts the impact of these perturbations in two backgrounds (r = 0.78 and 0.59
322 respectively, Fig S3), but no significant correlation between predictions and measured
323 values was found for the other three backgrounds (but this is also true for the other models
324 tested). On these two backgrounds, Enformer once again outperforms the competition (Fig
325 4A).

326 Additionally, Weingarten-Gabbay et al. measured the effect of shifting the position of a TATA
327 box in four different backgrounds. In two out of four backgrounds, Enformer predictions show
328 significant correlations with the measured effects (Fig S4). In the RPLP0 background, in
329 particular, the correlation is almost perfect (r = 0.9, Fig 4B). Given that Enformer pools the
330 sequence into 128bp bins, it is quite impressive that it is sensitive to such small shifts. Note,
331 however, that Basenji1 and Basenji2 - but not Xpresso - deliver similar predictions on this
332 task.

333 We wondered whether the large drop in both measured and predicted expression when the
334 TATA was placed around position -50 was a result of the position preference of the TATA or
335 due to another factor. As Weingarten-Gabbay et al. do not discuss this, we sought to use
336 Enformer to answer this question. For this purpose, we first shifted a neutral k-mer
337 (NNNNNNNN) through the sequence instead of the TATA. We find that this had little impact,
338 except again around position -50 (Fig S5). We then performed an in-silico mutagenesis

339  around this position and found that the expression prediction was most sensitive to the
340  bases ATTGG (Fig 4C). This Enformer-based analysis provides a biologically plausible
341  hypothesis, namely that the drop in expression is caused by a disruption of a CAT-box rather
342  than due to a positional preference of the TATA box.

343  Lastly, we reanalyzed the Kircher et al. [13] saturation mutagenesis data which was already
344  used as validation in the Avsec et al. [10] paper and additionally provided a comparison to
345  Basenji1, Basenji2, and Expecto. In this experiment, short stretches of sequence from 15
346  loci (mostly promoters) were selected to serve as regulatory elements for a reporter gene in
347  a plasmid. The authors then introduced almost every possible single nucleotide variant in
348  these sequences and recorded their respective impacts on the expression of the reporter.

349  Consistent with Avsec et al., we found that observed and predicted relative variant effects
350  (computed as $\log_2$ fold change of predicted expression at the variant, Methods) correlated
351  well for most loci. Moreover, Enformer outperformed all other models (Fig 4D).

352  However, Enformer predictions appeared to be miscalibrated between loci (Fig S6). While a
353  linear correlation is often present, the value of the slope varied drastically between loci. The
354  problem is striking when comparing the F9 and LDLR loci, which were both measured and
355  predicted in HepG2 (Fig 4E). As a consequence, while the individual correlations are
356  significant for both loci (r = 0.3 and r = 0.64 respectively), the correlation drops substantially
357  (r = 0.1) when pooling the data of the two loci. As those correlations are computed on log-
358  transformed abundance, these observations imply that Enformer variant effect predictions
359  only accurately reflect the experimental values up to a locus-specific exponential factor. This
360  being said, it is unclear whether this discrepancy is due to Enformer or an artifact of the
361  plasmid-based assay.

362  Interestingly, if we use accessibility (DNase) predictions as a proxy for RNA abundance
363  predictions, then the effects are on average too small (a predicted $\log_2$ fold change of 0.65
364  corresponds roughly to a measured effect of 1), but this scaling factor is by and large
365  consistent between loci (Fig 4E, Fig S7). Thus, for use cases such as genome-wide variant
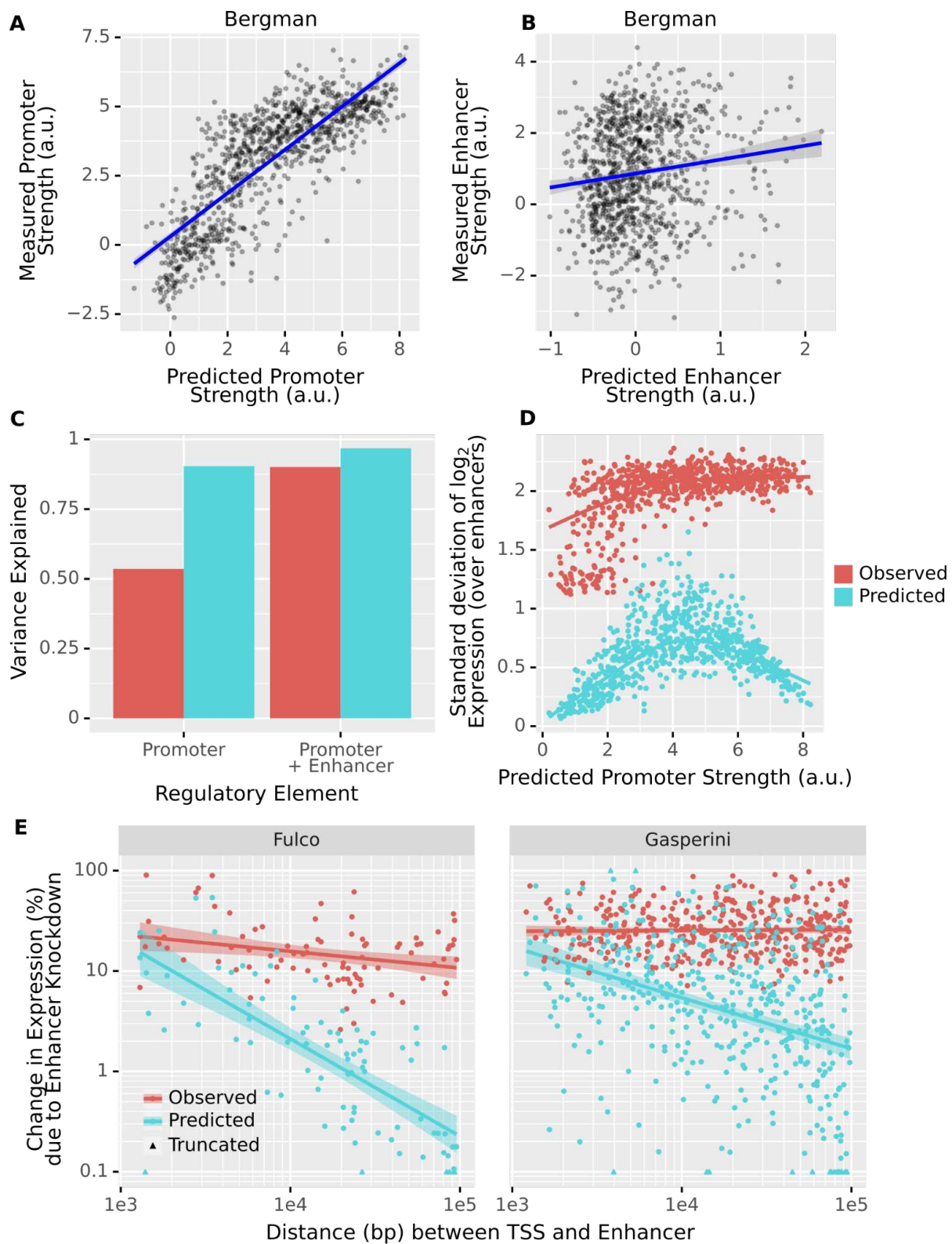366  prioritization, DNase prediction score will likely be more useful.

367  Overall, we see that current sequence-based models, particularly Enformer, can match
368  experimental data measuring the expression of different promoters in a controlled synthetic
369  context very well. Moreover, it can predict the impact of different TF-binding sites, single
370  nucleotide variants, and core promoter motifs, at least in some backgrounds. These
371  analyses thus show that Enformer significantly captures genuine causal regulatory elements
372  of the promoter.

13

### 373 Enformer attributes considerably less importance to enhancers than
### 374 experiments suggest

375



376

13

377 ***Figure 5: The predicted impact of enhancers, particularly distal enhancers, is***
378 ***significantly smaller than experiments suggest. A)*** *Predicted promoter strength vs*
379 *measured promoter strength. These strengths are determined by fitting linear models to the*
380 *data/predictions **B)** Predicted enhancer strength vs measured enhancer strength. **C)** The*
381 *promoter can explain 54% of the variation in the measurements of Bergmann et al. [12], with*
382 *enhancers explaining another 36%. However, 90% of the variation in Enformer predictions*
383 *for the same sequences is driven by the promoter alone. Thus, at least in a plasmid context,*
384 *Enformer strongly underestimates the importance of the enhancer for determining gene*
385 *expression. **D)** In Enformer, the predicted variation of expression induced by the enhancer*
386 *also heavily depends on the promoter. Only promoters of intermediate predicted strength are*
387 *sensitive to the choice of Enhancer. In the experimental data, strong and intermediate*
388 *promoters show similar sensitivity. **C)** The measured and predicted changes in gene*
389 *expression (expressed as an unsigned percentage) due to enhancer knockout as a function*
390 *of the distance between the gene and the enhancer. Values < 0.1% and > 100% are*
391 *truncated. Shown are only validated enhancer-gene pairs from Fulco et al. [20] and*
392 *Gasperini et al. [19]. Enformer attributes significantly less effect to most validated enhancers*
393 *than the experiments suggest. The effect is particularly strong for distal enhancers.*

394 Having established that Enformer captures causal elements located in promoters, we next
395 study to what extent it can correctly predict the effect of enhancers. Bergman et al. [12]
396 assayed all combinations of one thousand 264 bp human promoter and one thousand 264
397 bp human enhancer fragments using a plasmid-based MPRA. A major result of this study
398 was that transcriptional output appeared to be well modeled as the product of promoter
399 strength and enhancer strength.

400 In agreement with our above-mentioned analysis of promoters, Enformer predicted the
401 promoter strength of the Bergman et al. assay remarkably well (r = 0.81, Fig. 5A). In
402 contrast, the predicted enhancer effect correlates poorly with the reported enhancer strength
403 (r = 0.137, Fig. 5B). Moreover, promoters alone explained most of the Enformer predictions
404 (90% predicted variation in log expression) but only half of the experimentally observed
405 variation (54%, Fig. 5C).

406 This is not to say that the enhancers never matter in the predictions. If we plot, for each
407 promoter, the standard deviation in predicted log expression induced by the different
408 enhancers, we notice that the choice of enhancer does sometimes matter, but only for
409 promoters of intermediate predicted strength (Fig 5D). Promoters that Enformer considers
410 very strong seem to basically override the enhancer. In the experimental data, we do not
411 observe such a pattern.

412 Overall, this analysis indicates that Enformer inadequately accounts for the impact of
413 enhancers on gene expression and cannot predict their measured average effects. However,
414 one limitation of the Bergman et al. study is that it is based on a plasmid construct, which
415 Enformer was not trained to handle and may not reflect endogenous gene regulation.
416 Moreover, in this assay, enhancers are placed very close to the promoter and are part of the
417 transcript, which may interfere with co- and post-transcriptional mechanisms.

418 Given these concerns, we wondered whether the results described above generalize when
419 we apply Enformer to enhancer-promoter pairs in the endogenous genome. To find such
420 pairs, we use the enhancer knockdown screens performed by Gasperini et al. [19] and Fulco
421 et al. [20] In these screens, CRISPR interference (CRISPRi) is used to perturb enhancers
422 and then the corresponding impact on gene expression is measured.

423 To replicate the CRISPRi knockdown experiments in the model, we used an in-silico
424 mutagenesis (ISM) approach. Specifically, we computed the change in predicted expression

425   at the gene TSS upon shuffling the enhancer sequence (Methods), assuming that a shuffled
426   enhancer sequence should generally be non-functional. To account for random noise
427   introduced by shuffling, we repeated the procedure 25 times and recorded the average
428   impact.

429   We find that the model sometimes attributes plausibly large effects to enhancers, particularly
430   if they are promoter-proximal. However, as the distance between the TSS and the enhancer
431   increases, the model very quickly begins to excessively discount the importance of the
432   enhancer for gene expression (Fig 5E). This discounting scales proportionally to the inverse
433   of distance, which is not what is observed for validated enhancers. This is a strong decay,
434   which is already substantial at a distance of only 10 kb (2% for Fulco et al, 5% for Gasperini
435   et al., Fig 5E).

436   In the respective assays, knocking down the median validated enhancer (n = 522) has a
437   measured impact on expression of ~20%. According to Enformer, this median effect is only
438   ~4%. Moreover, for ~60% of genes with a validated enhancer (n = 385), none of the tested
439   candidates (whether validated or not, n = 2070) has an impact on expression that exceeds
440   10% (Fig S8).

441   As for the Bergman et al. [12] data, Enformer predicts smaller effects of enhancer
442   knockdown for genes with high predicted basal expression (i.e. the predicted expression
443   after knockdown, Fig S9). Whether or not a similar relationship exists in the experimental
444   data is difficult to say, because the underlying assays have higher power to detect smaller
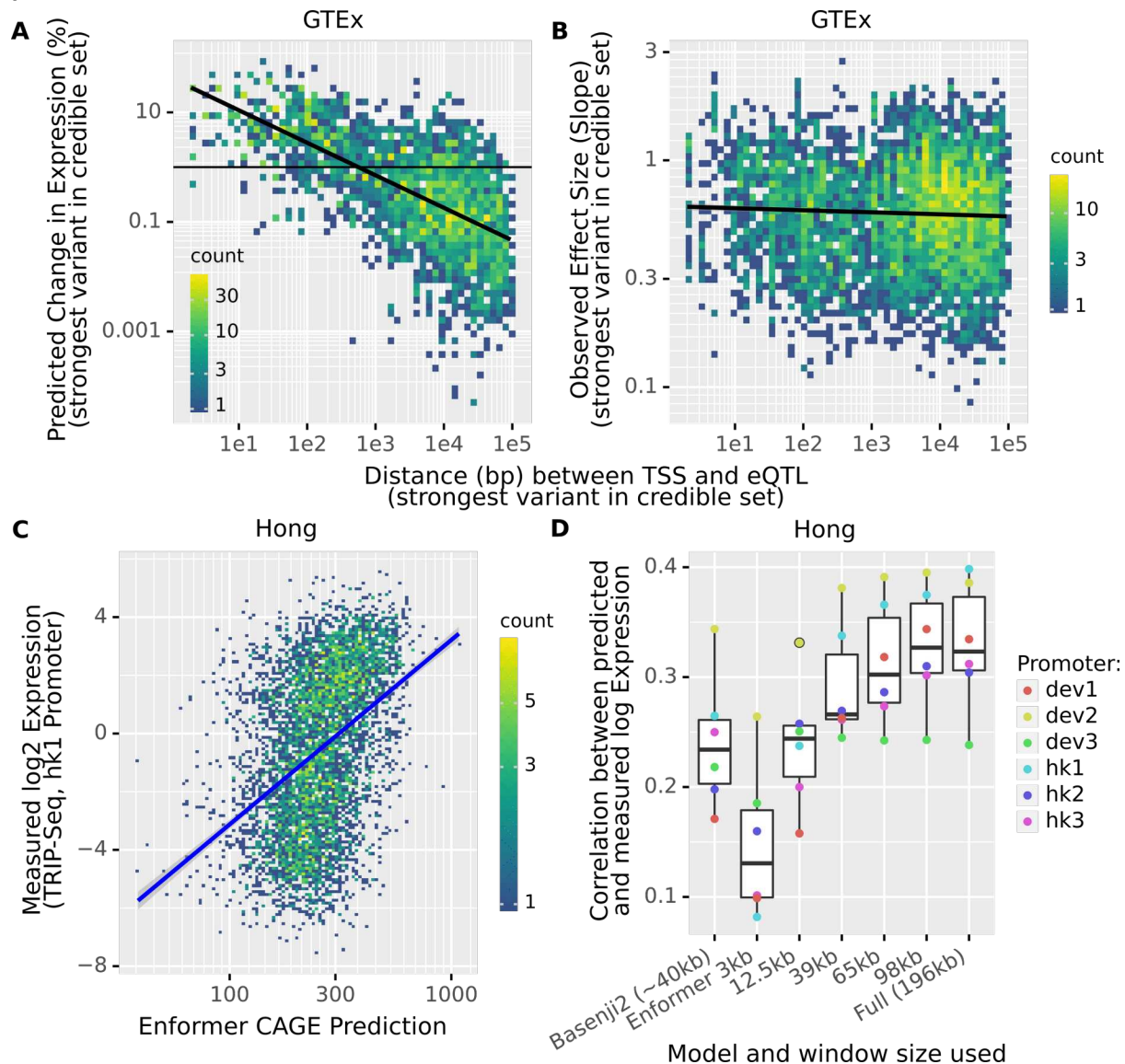445   effects for more highly expressed genes.

446   Our results do not contradict those of Avsec et al. [10], who showed with the same data that
447   Enformer performs reasonably well at prioritizing validated enhancer-promoter pairs. In fact,
448   on average, Enformer does attribute somewhat larger effects to validated than to non-
449   validated enhancers, after controlling for distance and basal expression of the promoter (Fig
450   S10). However, even when it correctly links enhancer and gene, this impact is usually small.
451   If we rank enhancers by their effect, then the threshold for ~50% recall of validated
452   enhancers is a ~3% predicted effect (Fig S11). Moreover, we here defined the enhancer
453   impact in relative terms, i.e. as percentage change, as we consider this to be the biologically
454   plausible metric. However, it is possible to achieve a slightly better classification in enhancer
455   prioritization using the predicted absolute change as in the original study (Fig S12). This is
456   because, in contrast to the predicted percentage change, which as we noted above,
457   decreases with promoter strength, the absolute change increases proportionally to the
458   predicted basal expression (Fig S13). Since predicted expression correlates well with actual
459   expression, this metric thus privileges highly expressed genes. This likely adds to the
460   classification performance because highly expressed genes mechanically will tend to have
461   more validated regulatory elements, due to the limited power of the underlying assay (Fig
462   S14). We, therefore, believe that the predicted absolute change delivers somewhat inflated
463   performance.

464   Overall, the model strongly underestimates the effect of known enhancers, particularly if they
465   are distal to the TSS, for gene expression regulation.

466

16

## Most distal eQTLs do not have a meaningful impact on expression predictions



**Figure 6: Enformer attributes no meaningful impact to distal eQTLs and performs poorly on tasks where long-range information is crucial. A)** *The measured and predicted changes in gene expression (expressed as an unsigned percentage) due to eQTL variants are plotted as a function of the distance between the (canonical) TSS and the eQTL. To account for linkage, we always take the maximal effect of all variants in the credible set. This predicted effect quickly decays with distance.* **B)** *The GTEx eQTL normalized effect size is plotted against the distance to the TSS. We observe no systematic decay with distance.* **C)** *Enformer struggles to predict the overall impact of the genomic environment on expression of the hk1 promoter, as measured by TRIP-seq. Note that this is the promoter for which the model performs best.* **D)** *Performance of Basenji2 and Enformer on the TRIP-Seq data. For Enformer, we computed predictions after restricting the input window, as previously. We find that most of the (limited) signal on this task once again derives from the proximal 20% of the input sequence.*

Given our results for enhancers, we next asked whether similar observations also apply for expression Quantitative Trait Loci (eQTLs).

16

485

486 To test this, we let Enformer predict the impact on gene expression of GTEx eQTLs [16]. We
487 applied a set of filters to exclude eQTLs that may act post-transcriptionally, have unclear
488 associated TSS, and for which Enformer cannot be straightforwardly applied (Methods). A
489 fine-mapping method was applied to associate each eQTL with a credible set of likely-causal
490 variants and account for linkage (SuSie [25,26]). We then computed the predicted impact on
491 expression at the canonical TSS of the target gene for all variants in the credible set. We
492 defined the predicted eQTL effect as the maximum of these individual impacts.

493

494 As for the enhancer analysis, Enformer predicted effects of eQTLs were large when close to
495 the TSS but decayed quickly and proportionally to the inverse of distance (Fig 6A). This
496 decay was extreme, showing an average predicted effect of less than 1% from 1kb on. We
497 do not see such decay in the measured normalized effect sizes of these eQTLs (Fig 6B).
498 Although the GTEx normalized effect sizes do not correspond to expression fold-changes
499 due to non-linear and gene-specific data preprocessing [16], a decreasing relationship with
500 distance would be expected if Enformer's predictions were correct. Moreover, Enformer
501 predicts for those eQTLs surprisingly tiny effects (median = 0.5% change). This gives further
502 evidence that the model strongly underweights the causal effect of distal regulatory elements
503 on gene expression.

## Current deep models mostly cannot predict the general impact of the genomic environment

506 To further probe the ability of Enformer to correctly account for distal sequence context in its
507 predictions, we used TRIP (thousands of reporters integrated in parallel) sequencing data
508 compiled by Hong et al. [27]. In this experiment, six different short promoter fragments were
509 integrated at thousands of locations across the genome and then their activity at each
510 location was measured. This assays the impact of the overall genomic environment on
511 promoter activity. We replicated this experiment in-silico by inserting the fragments at the
512 same locations and recording the predicted expression.

513 Consistent with the results described previously, Enformer perfectly ranked the promoters
514 according to their median expression (Fig S15). However, its ability to predict the variation in
515 expression of individual promoters when integrated at different genomic locations was quite
516 limited (r = 0.24-0.4 depending on the promoter, Fig 6C,D, Fig S16), albeit statistically
517 significant. Moreover, Enformer once again outperformed Basenji2 on this task, even when it
518 was limited to the same receptive field (Fig 6D).

519 Hong et al. also tested the impact of placing 676 different promoters in 4 different locations
520 in the genome (Patch-MPRA). When we reproduced this experiment in Enformer, we found
521 very coherent results: (1) Enformer predicts the average impact of promoters well (r = 0.73,
522 Fig S17), (2) the promoters alone explain more variation in Enformer predictions than in the
523 experimental data (69% vs 55%), and (3) only intermediate-strength promoters, but not
524 strong promoters, are predicted to be strongly affected by the genomic environment (Fig
525 S18). Of note, this is somewhat consistent with the experimental data, but there the effect is
526 far less pronounced (Fig S19).

527 Admittedly, the fact that both Basenji2 and more so Enformer have any predictive power on
528 the TRIP-seq data at all is impressive since it requires predicting the impact of a rather
529 extreme and unnatural modification of the genome. Nevertheless, the limited performance

17

18

530   on this data further indicates that more research is needed to achieve models of
531   transcriptional regulation which properly capture the effects of distal elements and genomic
532   context.

## Enformer's promoter-enhancer logic is (mostly) multiplicative

534   A remarkable result, obtained independently by Bergman et al. [12] and Hong et al. [27], is
535   that promoters and enhancers generally act multiplicatively on transcriptional input. We
536   asked whether Enformer qualitatively captured the same rule. To this end, we designed an
537   in-silico experiment. This is not trivial to test because, as discussed previously, the promoter
538   usually dominates all other sequence elements in Enformer. Therefore, we selected
539   endogenous triplets of promoter, enhancer and sequence background which were such that
540   the enhancer had a notable predicted impact on the promoter in that particular background
541   (Methods). We then predicted for every combination of these promoters (n = 89) and
542   enhancers (n = 115) in each background (n = 32). Remarkably, the Enformer predictions
543   could be well approximated as the product of promoter strength, enhancer strength, and
544   background strength (77% of the variance, Methods).

545   If we focus on individual backgrounds, we find that in some cases the enhancers can explain
546   up to 25% of the variation in predicted expression at the promoter (Fig S20). Moreover,
547   except for a few backgrounds where the variation in expression can be almost completely
548   explained by a background strength, we find that the enhancer strengths are highly
549   correlated across backgrounds (Fig S21). We find very similar results if we focus on each
550   individual promoter (Fig S22). Thus, it appears that the background and promoter determine
551   whether enhancers have any impact, but if they do, this impact is consistent.

552   Altogether, this analysis showed that Enformer qualitatively agrees with the recently reported
553   multiplicative model of promoter and enhancer effects, despite quantitatively underestimating
554   the effects of distant regulatory elements.

# Discussion

556   Here we have performed the most extensive benchmark of sequence-based models of gene
557   expression to unseen data. Specifically, we compared the predictions of the models against
558   two large-scale observational RNA-seq datasets of adult and developmental tissue, as well
559   as five deep perturbation assays, comprising designed reporter assays and CRISPRi
560   screens. With our approach, we specifically probed to what extent current sequence-based
561   models can account for the regulatory role of promoters, enhancers, eQTLs, and genomic
562   environment. This allowed us to evaluate the generalization power of these models. We
563   found that current sequence-based models show a remarkable ability to predict the
564   expression associated with human promoters. However, we observed that even Enformer,
565   with its very wide receptive field, can account for distal regulators only to a limited extent.

566   We will now discuss in more detail our findings and provide suggestions on how the field
567   could use and improve sequence-based models of gene expression.

## Do current sequence-based models learn causal effects?

569   Across a wide range of analyses, we repeatedly observed that current sequence-based
570   models, particularly Enformer, predict the impact of different promoters very well. It does not
571   matter if the strength of a promoter is measured in its proper genomic context, in a
572   completely different context, or even in a plasmid - Enformer predictions of expression

18

573   generally will show substantial correlations with these measurements.

574   Moreover, we observed that Enformer often correctly predicts the impact on expression of
575   diverse promoter modifications. These include single-nucleotide variants, addition or removal
576   of TF-binding sites, and in some cases mutagenesis or shifts of core-promoter motifs. This is
577   strong evidence that the model captures causal determinants of the promoter.

578   In contrast to the strong performance for promoters, the model underperforms when it comes
579   to accounting for the expression effect of enhancers, particularly if they are distal. However,
580   we note that, since our analysis was focused on expression, we did not analyze whether the
581   deep models can correctly account for the causal impact of enhancer variants or
582   perturbations on epigenetic marks at the enhancer itself. Performing a similar analysis as we
583   have done but focused on ATAC-Seq or CHIP-Seq-based MPRA or allele-specific binding
584   data [28] could shed light on this question.

585   Accordingly, our results should not be construed to mean that such models do not further our
586   understanding of enhancer biology or enhancer variants. We focused on showing that these
587   models attribute far too little importance to enhancers when predicting gene expression.
588   Indeed it has been repeatedly demonstrated in the literature that using the full range of
589   epigenetic predictions can provide real added value in enhancer prioritization or variant fine-
590   mapping [8–10,29].

## Do current sequence-based models learn long-range effects?

592   Overall, we found limited evidence that the sequence-based models we studied make use of
593   long-range information. For Enformer, the model with the largest receptive field (196kb), we
594   found that we can safely remove two-thirds of the input sequence, with minimal impact on
595   predictions. As we saw from the distribution of eQTLs and validated enhancers this does to
596   some extent reflect the underlying biology. Indeed, most strong regulatory elements do seem
597   to be located relatively close to their target gene. However, the datasets we analyzed clearly
598   showed that distal elements with large causal impacts on expression exist and that Enformer
599   will generally strongly underestimate their impact. As a result, the majority of validated
600   enhancers and eQTLs do not have a meaningful impact on Enformer predictions of gene
601   expression.

602   We hypothesize that a large part of the problem is the escalating class imbalance: The ratio
603   between actual regulatory elements among all candidate regulatory elements decreases with
604   increasing distance from a gene. Perhaps, the model responds to this worsening signal-to-
605   noise ratio by down weighting distal elements and thus, in some sense, distributing their
606   effect. One possible underlying reason is the model being unable to correctly identify which
607   distal regulatory elements will be functional in a certain cell type. Another non-exclusive,
608   possibility is that the model is unable to correctly link distal regulatory elements to their target
609   genes. In other words, is the source of difficulty the complexity of the enhancer code or is it
610   the complexity of the "folding code"? Disentangling these issues will help to determine which
611   additional training data and modeling assumptions are most likely to yield substantial
612   improvements.

## Do current sequence-based models learn between-condition regulation?

614   We showed that Enformer, and to a lesser extent also Basenji2, predictions do exhibit
615   significant correlations with measured between-condition differences in expression. Thus,
616   clearly, the models capture some signal on this task. However, the between-condition
617   predictions are less impressive, yet decent, than the between-gene ones.

20

618 We showed that this is partially because even Enformer has not yet reached the level of
619 precision necessary to predict between-condition variation, which is generally of a smaller
620 magnitude than between-gene variation. Thus, in principle, as models improve overall, they
621 should particularly improve on between-condition tasks.

622 Our between-condition analysis was performed using endogenous gene expression and not
623 using controlled perturbation experiments. For this reason, we cannot exclude the possibility
624 that some of the between-condition predictive power is correlative. Furthermore, we cannot
625 determine whether Enformer struggles at between-condition predictions precisely because of
626 its limited understanding of distal regulatory elements. Since distal enhancers appear to
627 mostly be a feature of multicellular organisms, one of their main roles could be to control
628 cell-type specific expression. Deep perturbation assays performed across multiple cell types
629 are needed to address these issues.

## Implications for in-silico investigations of gene expression

631 Our results indicate that sequence-based models can now be usefully employed to study
632 promoter mutations in silico. This could prove useful for fine mapping in GWAS studies, for
633 the detection of possibly pathogenic variants in rare disease diagnostics and oncology
634 settings, and to study the evolution of human promoters. In-silico experimentation and
635 interpretation techniques may even yield new candidate motifs or mechanistic hypotheses,
636 which can then be tested in future designed MPRA studies. Assuming sufficient GPU
637 resources, we recommend using Enformer for these tasks, as it clearly outperformed all
638 other methods. However, some limitations of Enformer should be considered. Because it
639 predicts in 128bp bins, Enformer cannot be used to predict the exact locations of TSS sites.
640 Also, the model cannot reliably predict the directionality of a promoter because it is trained to
641 predict similar values for both strands.

642 We found that variant effect predictions on gene expression from Enformer appear to be
643 miscalibrated when comparing between promoters. Thus, genome-wide rankings of
644 promoter variants made using Enformer may be misleading. We cannot fully exclude the
645 possibility that this is an artifact of the plasmid-based measurements, rather than a failure of
646 the model. Nevertheless, if a genome-wide ranking is the goal, using the predicted DNase
647 variant effects (or a mixture of CAGE and DNase) appears to be the safer choice.

648 Our results indicate that predicting the impact of a distal variant on expression at the TSS
649 rarely leads to meaningful results.

650 Finally, we note that, in addition to using sequence-based models to analyze variants or to
651 prioritize enhancers, they also have potential use as an in-silico experimental platform to
652 explore more abstract biological questions. A recent study, for example, used a deep model
653 to analyze the impact of helical periodicity on TF binding [30]. As models mature, open
654 questions in epigenomics, such as the pioneer factor hypothesis, the billboard vs.
655 enhanceosome discussion, and others may become amenable to in-silico analysis. This
656 could help to design more targeted experiments and also aid in the interpretation of
657 experimental results.

## Implications of our study for model development

659 The recent trend in deep supervised models of regulatory genomics has been towards
660 expanded receptive fields, through the use of dilated convolutions and attention. These
661 expansions have usually been associated with increased performance, suggesting that the
662 wider receptive field is the main cause of the improvement.

20

663 Paradoxically, however, we observed that Enformer, despite its huge receptive field, derives
664 most of its predictive power from a small fraction of the sequence. Moreover, Enformer
665 substantially outperformed Basenji2 even when it is restricted to the latter model's input
666 window and even on tasks where the receptive field size is irrelevant (such as plasmid-
667 based saturation mutagenesis). Therefore, it is unclear to what extent the large receptive
668 field of Enformer actually contributes meaningfully to its predictive power. Perhaps, similar to
669 results in natural language processing [11,31], the mere amount of parameters combined
670 with the transformer architecture, is the driving force behind the improvement.

671 How could a very good model of gene expression that does account for distal regulators be
672 built? Given that Enformer already achieves close to replicate predictions on CAGE in
673 ENCODE, it is likely that data complementary to ENCODE is required. Which data will help
674 the most depends on the exact nature of the problem. If the issue is that Enformer does not
675 correctly interpret the enhancer code, then adding more epigenetic, as well as expression,
676 data from more cell types and more species should be most beneficial. More cell types
677 provide more variation in enhancer activity, whereas more species introduce not just
678 variation in activity but also variation in functional enhancer sequence. Recent advances
679 have shown that, if done at a sufficient scale, perturbation experiments can also be used to
680 train models [32–34], although further progress will require modeling advances to
681 seamlessly integrate global and targeted assays.

682 If the main difficulty is folding, then adding the prediction of Hi-C data as an auxiliary task, or
683 integrating a model which can predict Hi-C, would be the natural solution. Promising models
684 predicting chromosomal contacts from sequence have been obtained [35–37], but
685 integration with gene expression prediction models remains to be explored.

686

# Methods

## Models used

689 Basenji1/2 have a very similar structure as Enformer. Accordingly, Basenji1/2 predictions
690 were always made in exactly the same way as Enformer predictions - the only difference
691 being a smaller input window. We used the pretrained Basenji2 model from:
692 https://storage.googleapis.com/basenji_barnyard/model_human.h5. For Basenji1, we used
693 the model from the Kipoi model repository [38].
694 Xpresso only predicts a single value per sequence. Moreover, it is strand-specific and it
695 expects the TSS to be at a particular position in the sequence. In cases where the TSS is
696 known, we thus made only one prediction with Xpresso, namely with the TSS at the correct
697 location and on the correct strand. In the Segal dataset, where the exact TSS is not known,
698 we tried a number of offsets and selected the best ones. We used the Xpresso model from
699 the Kipoi model repository, which does not account for RNA half-life covariates.
700 To get Expecto predictions of variant effects, we used the web interface at:
701 https://hb.flatironinstitute.org/expecto/.

## Endogenous Expression: GTEx and Cardoso-Moreira et al.

703 We collected gene expression measurements for GTEx tissues from the GTEx consortium
704 webpage and for different development stages from ArrayExpress [39] (E-MTAB-6814). Note
705 that this data was already normalized for sequencing depth and gene length. We log-

706   transformed the data, adding a pseudo count of one. We exclude mitochondrial genes and
707   the Y chromosome. We also exclude all genes which are never expressed in our data.

708   We extracted the reference genome sequence (hg38) around each gene, centered on the
709   canonical TSS (i.e. the TSS of the ENSEMBL [17] Canonical transcript), and computed the
710   predictions of expression at this TSS for each of the models.

711   Because Enformer/Basenji2 provide CAGE predictions only for certain ENCODE cell lines,
712   which do not permit a 1:1 matching to GTEx tissues, we instead perform the matching using
713   ridge regression. Specifically, we fit L2 regularized regressions for each tissue such that:

714
$$g_{it} = \mathbf{a}_t^T \mathbf{x}_i + \epsilon$$

715   Where $g_{it}$ is the (log-transformed) expression of gene i in tissue t, $\mathbf{a}_t$ is the vector of learned
716   weights for tissue t and $\mathbf{x}_i$ is the vector of (log-transformed) Enformer/Basenji2 CAGE
717   predictions for gene i.  Note that the intercept is zero by construction (see below).

718   To fit these regressions, we split the genes into train and test-set, whereby every gene which
719   is fully enclosed in an Enformer test-set region is included in the test set, and all genes
720   which never intersect any test-set region go to the train set. Genes which intersect a test-set
721   region, but are not contained by it, are excluded. In this way we prevent contamination of the
722   held out test-set.

723   To make between-condition comparisons meaningful, we compute the mean for each tissue
724   (or tissue-development-stage in the Cardoso-Moreira et al. [15] data) on the training set and
725   remove this mean from both the training and the test set. This ensures that our regressions
726   cannot learn tissue means, in a way that prevents leakage from the test set.

727   For Xpresso we follow the same procedure, but as this model is cell-type agnostic, the ridge
728   regression only rescales the predictions.

729   We do not show the results for Basenji1 on this data, as Basenji1 used a different train-test
730   split, thus the numbers are not comparable (nevertheless, it still performs worse than
731   Enformer). Note that the same is true for Xpresso, but we consider it unlikely that this model
732   overfits very much.

## Sequence context ablation study

734   We follow the same steps as we did previously for the GTEx and Cardoso-Moreira et al.
735   data. The only difference being that we now generate predictions with different sequence
736   window sizes around the TSS. We use the following window sizes: 1001, 3001, 12501,
737   34501, 39321, 49153, 65537, 98305 and 131073 bp. The last six correspond to a fifth, a
738   fourth, a third, half and two-thirds of the total receptive field respectively. These windows are
739   always centered on the TSS (so a window size of 1001bp means we extract the TSS site +/-
740   500bp). As the windows are smaller than Enformer's receptive field, we pad with "N"
741   nucleotides on the flanks.

742   We train separate regression models for each window size, using the same train-test split as
743   previously

## Class Imbalance

745   We downloaded ENCODE CREs [23] from https://screen.encodeproject.org/.

746  We used PyRanges [40] to define "Enformer-sized" (i.e. 196kb) windows around each gene
747  of interest (i.e. the genes with validated enhancers or eQTL). We then intersected our sets of
748  regulatory elements (eQTL, CRISPRi validated and non-validated enhancers and ENCODE
749  CRE) with these windows and for each hit we recorded the distance to the gene. Note that
750  we apply the same filters we applied to the eQTL data also to the ENCODE CRE, so as to
751  make these sets comparable (these filters are discussed in the eQTL methods section). The
752  validated and non-validated enhancers are comparable by design.

## Promoter determinants: Weingarten-Gabbay et al.

754  We collected the data from the supplementary materials of the Weingarten-Gabbay et al.
755  [24] manuscript and the construct sequences were kindly provided by the first author. We
756  followed the general procedures outlined previously to construct sequences and compute
757  predictions. As the exact TSS is, to our knowledge, unknown, we center predictions on the
758  midpoint of the promoter fragments. As this is a K562 experiment, we use the K562 CAGE
759  track as the predictor.

## Saturation Mutagenesis: Kircher et al.

761  Vikram Agarwal kindly provided the Kircher et al. [13] data.

762  To generate a variant effect prediction with Enformer (and other models), we create two
763  sequences: one centered on the reference allele and one on the alternative allele. We
764  predict for both (averaged over strands, offsets, and the neighboring bins as always), and
765  compute the log fold change in prediction centered on the variant (averaged over the
766  neighboring bins).

767  Kircher et al. used a number of different cell lines in their experiment, depending on the
768  locus. To match these cell lines to ENCODE tracks, we followed the same procedure as
769  Avsec et al. [10]. Specifically, we used tracks (CAGE/DNASE) whose ENCODE descriptions
770  contained substrings that correspond (more or less) to the cell line used for a particular
771  locus:

772  -  'HepG2' for *F9*, *LDLR*, and *SORT1*
773  -  'K562' for *GP1BB*, *HBB*, *HBG1*, and *PKLR*
774  -  'HEK293' for *HNF4A*, *MSMB*, *TERT* and MYCrs6983267
775  -  'pancreas' for *ZFAND3*
776  -  'glioblastoma' for *TERT*
777  -  'keratinocyte' for *IRF6*
778  -  'SK-MEL' for *IRF4*

779  If there was more than one matching track, we averaged predictions over the matching
780  tracks. This is different from Avsec et al., who instead extract principal components, but
781  ultimately this procedure yields very similar correlations.

782  In cases where there was no match at all, we averaged predictions over all tracks.

## Promoter x Enhancer: Bergman et al.

784  We collected the data and plasmid sequence from the supplementary materials of the
785  Bergman et al. [12] manuscript. We employ the same filtering strategy, keeping only data
786  points supported by at least 25 plasmids and at least 2 barcodes.

787  As this assay was conducted using a plasmid and has no clear analog in the endogenous

788 genome, we reproduced the plasmid sequences in-silico and added N-padding on the flanks
789 to adapt it to the Enformer input length. We placed the promoter and enhancer fragments
790 into their respective locations in the plasmid and centered predictions on the midpoint of the
791 promoter fragment, as the exact TSS is - to our knowledge - unknown.

792 Bergman et al. use their data to impute the intrinsic strengths of the promoter and enhancer
793 sequences. For this, they fit a Poisson model with promoter and enhancer indicators:

794
$$\mathrm{RNA} \sim \mathrm{Poisson}(\lambda = \exp(\beta \log(\mathrm{DNA}) + P + E))$$

795 where RNA is the measured RNA count, DNA is the amount of plasmid used, $\beta$ is a learned
796 weight, and P and E are the promoter and enhancer indicators ("strengths") respectively. We
797 reproduced this analysis using the package statsmodels [41]. Note that we also fitted a log-
798 linear OLS to this data, which gave very similar results, but we report the results of the
799 Poisson model to stay faithful to the source material.

800 To impute the *predicted* promoter and enhancer strengths, we use a similar strategy.
801 However, we do not fit a Poisson model, as Enformer predictions are not integer-valued.
802 Instead we fit a Gamma model:

803
$$\mathrm{CAGE}_{\mathrm{Enformer}} \sim \mathrm{Gamma}(\exp(P + E))$$

804 where CAGE is the Enformer CAGE prediction (for K562).

805 The Gamma distribution is often used to model the prior distribution of a Poisson lambda
806 and arguably the Enformer prediction in natural scale is a Poisson lambda, as Enformer is
807 trained using a Poisson loss function.

808 To ensure that our results are robust, we additionally fitted a linear regression to the log-
809 transformed Enformer prediction:

810
$$\log(\mathrm{CAGE}_{\mathrm{Enformer}} + 1) \sim \mathrm{Normal}(\mu = P + E, \sigma^2)$$

811 This gave very similar results.

812 Note that our calculations of explained variance refer to the variation in log expression. This
813 is why they slightly differ from the ones reported in Bergman et al.

## Enhancer knockdown

815 We collected the data from the supplements of the respective manuscripts. Additionally, the
816 sequences used for this benchmark in the Enformer paper were kindly provided by Ziga
817 Avsec. For each gene, Avsec et al. [10] determined the TSS site with the highest predicted
818 expression in K562. They then extracted sequences centered on these TSS, unless the
819 distance between the gene and the enhancer is such that this was not possible given the
820 receptive field size, in which case the TSS is shifted to accommodate the tested enhancers.
821 Predictions are made at the TSS sites, as per usual.

822 For each reference sequence (with the enhancer intact) we create 25 "knockout" sequences,
823 where we shuffle 2000 bp centered on the enhancer. In this way, we destroy the enhancer
824 without changing the nucleotide composition of the underlying sequence. The predicted
825 effect of enhancer knockout is then given by the average (over shuffles) change in predicted
826 gene expression at the TSS.

## eQTL

We downloaded SuSie credible sets for GTEx eQTL from the EBI eQTL catalog [26]. We additionally downloaded GTEx eQTL normalized effect sizes from the GTEx portal [14,16].

We apply the following filters:

- We only consider protein-coding eGenes.
- We exclude credible sets which span more than 5kb
- We demand that each variant in the credible set can be scored by Enformer when the canonical TSS is placed in the center of the sequence. In other words, the furthest variant of a credible set must be no further than 98kb from the canonical TSS of the eGene, otherwise the entire set is excluded.
- We demand that all variants in the credible set are upstream of the canonical TSS. This is to exclude post-transcriptionally acting variants (i.e. NMD variants, splice variants, etc). If a credible set contains even one downstream variant, we exclude it.
- We exclude all eGenes which have a GENCODE annotated protein-coding transcript upstream of the canonical one. In this way, we ensure that the canonical TSS is always the closest (protein-coding) TSS of the eGene to the variant.

6141 credible sets pass our filters, thus providing a total of 14139 variants to test.

For each variant in each credible set, we then compute predictions for all CAGE tracks at the canonical TSS of the eGene (using our usual strategy of summing over neighboring bins and averaging over strands and small offsets). We then use the previously fitted ridge regressions to match these CAGE predictions to GTEx tissues. Next, we compute, for each variant, the change in prediction vis-a-vis the predictions made with the reference sequence. We define the eQTL effect as the effect of the strongest variant (i.e. the one leading to the biggest change in predicted expression at the canonical TSS of the eGene in the tissue of interest) in the credible set. This follows the usual assumption in the literature that generally only one variant in a linkage block will be causal. Our strategy fails if the eQTL effect is the result of an epistatic interaction between several variants in linkage - however testing this possibility would require testing all combinations of variants in a credible set, which would be prohibitively expensive. Moreover, we would expect that such cases are rare anyways.

## Trip-Seq (Hong et al.)

We collected the data and relevant construct sequences from the supplements of the Hong et al. [27] manuscript. We followed the standard procedure to replicate this experiment.

To compute the explained variances, we fit log-linear OLS models, similarly as we did for the Bergman et al. data.

## In-silico Multiplicativity Assay

To identify triplets of promoter, enhancer, and background where Enformer attributes importance to the enhancer in determining expression at the promoter, we returned to our in-silico reproduction of the CRISPRi data. We first selected promoter-enhancer pairs where:

- the enhancer is further than 3kb from the TSS
- the enhancer is within 90kb of the TSS
- the in-silico enhancer knockdown has a predicted impact of at least 30% (note that this includes both repressive and activating enhancers)

26

869 To cover a slightly bigger range of possible enhancer effects, we identified for the same
870 genes some weak ( < 1% predicted effect, despite being located within 20kb of the TSS) and
871 intermediate strength enhancers (between 4% and 8%) predicted effect. We finally selected
872 a few "non-functional" promoter-enhancer pairs where the enhancer had no real effect at all (
873 < 0.1%).

874 This procedure yielded 89 unique promoters and 115 unique enhancers. To extract the
875 promoter sequences, as exact boundaries are unknown, we take a 1kb window around the
876 TSS. For the enhancer, we take a 2kb window around the enhancer midpoint.

877 Each promoter-enhancer pair comes with an endogenous sequence background, which also
878 determines the distance of the enhancer to the promoter. Ideally, we would have tested
879 every promoter-enhancer combination in every background, but the combinatorial explosion
880 makes this computationally expensive - specifically because we have to predict for each
881 sequence six times, i.e. for both strands and with small shifts, to limit the noise in the
882 prediction. Thus we selected 32 backgrounds. We took 6 backgrounds from our set of non-
883 functional triplets (to see if other enhancer-promoter pairs could be linked in such a
884 background) and we selected 6 backgrounds where the enhancer location is far from the
885 TSS location. The remaining backgrounds were selected from the strong triplets.

886 We then proceed as follows: for each background, we predict expression in K562 at the TSS
887 for every combination of promoter and enhancer. Hereby, we always replace the
888 endogenous promoter of this background with the promoter of interest and we replace the
889 endogenous enhancer with the enhancer of interest. Thus, for a given background, the
890 distance between promoter and enhancer is constant. We focus on K562 as most of the
891 experiments on this topic were also performed in this cell type.

892 We get a total of 327,520 combinations (promoter x enhancer x background).  We then fit
893 log-linear models to explain the variation in log expression in this data. We first do this for
894 the entire dataset, using an indicator variable for each promoter, each enhancer, and each
895 background (236 parameters). This model thus assumes that the log expression of a certain
896 promoter-enhancer pair in a certain background is determined by the innate strength of the
897 promoter which is scaled by the background and the enhancer.

898 We next fit log-linear models to explain the expression variation for each background, across
899 promoters and enhancers (32 backgrounds, with 10235 observations for each one). We use
900 promoter and enhancer indicators (204 parameters). We also fit log-linear models to explain
901 the expression variation for each promoter, across backgrounds and enhancers (89
902 promoters, with 3680 observations for each one). In this case we use background and
903 enhancer indicators (147 parameters). Lastly, we correlate the enhancer parameters we
904 estimated across promoters and across backgrounds.

# 905 Declarations

## 906 Ethics approval and consent to participate

907 Not applicable.

## 908 Consent for publication

909 Not applicable.

## Availability of data and materials

The datasets analyzed during the current study are available in the Zenodo repository, https://doi.org/10.5281/zenodo.7076228. All scripts used in the analysis can be found under https://github.com/Karollus/SequenceModelBenchmark.

## Competing Interests

The authors declare that they have no competing interests.

## Funding

## Author contributions

AK and JG conceived the study and wrote the manuscript. AK analyzed the data. TM and AK designed, implemented, and tested the computational pipeline. All authors read and approved the final manuscript.

## Acknowledgments

# References

1. Schwanhäusser B, Busse D, Li N, Dittmar G, Schuchhardt J, Wolf J, et al. Global quantification of mammalian gene expression control. Nature. Nature Publishing Group; 2011;473:337–42.

2. Eraslan B, Wang D, Gusic M, Prokisch H, Hallström BM, Uhlén M, et al. Quantification and discovery of sequence determinants of protein-per-mRNA amount in 29 human tissues. Mol Syst Biol. 2019;15:e8513.

3. Agarwal V, Shendure J. Predicting mRNA Abundance Directly from Genomic Sequence Using Deep Convolutional Neural Networks. Cell Rep. 2020;31:107663.

4. Kelley DR, Reshef YA, Bileschi M, Belanger D, McLean CY, Snoek J. Sequential regulatory activity prediction across chromosomes with convolutional neural networks. Genome Res. 2018;28:739–50.

5. Kelley DR. Cross-species regulatory sequence activity prediction. PLoS Comput Biol.

945    Public Library of Science; 2020;16:e1008050.

946    6. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human
947    genome. Nature. 2012;489:57–74.

948    7. Luo Y, Hitz BC, Gabdank I, Hilton JA, Kagda MS, Lam B, et al. New developments on the
949    Encyclopedia of DNA Elements (ENCODE) data portal. Nucleic Acids Res. 2020;48:D882–9.

950    8. Zhou J, Theesfeld CL, Yao K, Chen KM, Wong AK, Troyanskaya OG. Deep learning
951    sequence-based ab initio prediction of variant effects on expression and disease risk. Nat
952    Genet. 2018;50:1171–9.

953    9. Zhou J, Troyanskaya OG. Predicting effects of noncoding variants with deep learning-
954    based sequence model. Nat Methods. 2015;12:931–4.

955    10. Avsec Ž, Agarwal V, Visentin D, Ledsam JR, Grabska-Barwinska A, Taylor KR, et al.
956    Effective gene expression prediction from sequence by integrating long-range interactions.
957    Nat Methods. Nature Publishing Group; 2021;18:1196–203.

958    11. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention Is All
959    You Need. 2017 [cited 2022 Sep 13]; Available from:
960    http://dx.doi.org/10.48550/arXiv.1706.03762

961    12. Bergman DT, Jones TR, Liu V, Ray J, Jagoda E, Siraj L, et al. Compatibility rules of
962    human enhancer and promoter sequences. Nature. Nature Publishing Group;
963    2022;607:176–84.

964    13. Kircher M, Xiong C, Martin B, Schubach M, Inoue F, Bell RJA, et al. Saturation
965    mutagenesis of twenty disease-associated regulatory elements at single base-pair
966    resolution. Nat Commun. Nature Publishing Group; 2019;10:1–15.

967    14. Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, et al. The Genotype-
968    Tissue Expression (GTEx) project. Nat Genet. Nature Publishing Group; 2013;45:580–5.

969    15. Cardoso-Moreira M, Halbert J, Valloton D, Velten B, Chen C, Shao Y, et al. Gene
970    expression across mammalian organ development. Nature. 2019;571:505–9.

971    16. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across
972    human tissues. Science. 2020;369:1318–30.

973    17. Cunningham F, Allen JE, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, et al.
974    Ensembl 2022. Nucleic Acids Res. 2022;50:D988–95.

975    18. Xia B, Yan Y, Baron M, Wagner F, Barkley D, Chiodin M, et al. Widespread
976    Transcriptional Scanning in the Testis Modulates Gene Evolution Rates. Cell.
977    2020;180:248–62.e21.

978    19. Gasperini M, Hill AJ, McFaline-Figueroa JL, Martin B, Kim S, Zhang MD, et al. A
979    Genome-wide Framework for Mapping Gene Regulation via Cellular Genetic Screens. Cell.
980    2019;176:1516.

981    20. Fulco CP, Nasser J, Jones TR, Munson G, Bergman DT, Subramanian V, et al. Activity-
982    by-contact model of enhancer-promoter regulation from thousands of CRISPR perturbations.

983 Nat Genet. 2019;51:1664–9.

984 21. Võsa U, Claringbould A, Westra H-J, Bonder MJ, Deelen P, Zeng B, et al. Large-scale
985 cis- and trans-eQTL analyses identify thousands of genetic loci and polygenic scores that
986 regulate blood gene expression. Nat Genet. 2021;53:1300–10.

987 22. Ferraro NM, Strober BJ, Einson J, Abell NS, Aguet F, Barbeira AN, et al. Transcriptomic
988 signatures across human tissues identify functional rare genetic variation. Science [Internet].
989 2020;369. Available from: http://dx.doi.org/10.1126/science.aaz5900

990 23. ENCODE Project Consortium, Moore JE, Purcaro MJ, Pratt HE, Epstein CB, Shoresh N,
991 et al. Expanded encyclopaedias of DNA elements in the human and mouse genomes.
992 Nature. 2020;583:699–710.

993 24. Weingarten-Gabbay S, Nir R, Lubliner S, Sharon E, Kalma Y, Weinberger A, et al.
994 Systematic interrogation of human promoters. Genome Res. 2019;29:171–83.

995 25. Zou Y, Carbonetto P, Wang G, Stephens M. Fine-mapping from summary data with the
996 "Sum of Single Effects" model. PLoS Genet. 2022;18:e1010299.

997 26. Kerimov N, Hayhurst JD, Peikova K, Manning JR, Walter P, Kolberg L, et al. A
998 compendium of uniformly processed human gene expression and splicing quantitative trait
999 loci. Nat Genet. 2021;53:1290–9.

1000 27. Hong CKY, Cohen BA. Genomic environments scale the activities of diverse core
1001 promoters. Genome Res. 2022;32:85–96.

1002 28. Abramov S, Boytsov A, Bykova D, Penzar DD, Yevshin I, Kolmykov SK, et al. Landscape
1003 of allele-specific transcription factor binding in the human genome. Nat Commun.
1004 2021;12:2751.

1005 29. Chen KM, Wong AK, Troyanskaya OG, Zhou J. A sequence-based global map of
1006 regulatory activity for deciphering human genetics. Nat Genet. Nature Publishing Group;
1007 2022;54:940–9.

1008 30. Avsec Ž, Weilert M, Shrikumar A, Krueger S, Alexandari A, Dalal K, et al. Base-
1009 resolution models of transcription-factor binding reveal soft motif syntax. Nat Genet.
1010 2021;53:354–66.

1011 31. Kaplan J, McCandlish S, Henighan T, Brown TB, Chess B, Child R, et al. Scaling Laws
1012 for Neural Language Models. 2020 [cited 2022 Sep 13]; Available from:
1013 http://dx.doi.org/10.48550/arXiv.2001.08361

1014 32. de Almeida BP, Reiter F, Pagani M, Stark A. DeepSTARR predicts enhancer activity
1015 from DNA sequence and enables the de novo design of synthetic enhancers. Nat Genet.
1016 Nature Publishing Group; 2022;54:613–24.

1017 33. Bogard N, Linder J, Rosenberg AB, Seelig G. A Deep Neural Network for Predicting and
1018 Engineering Alternative Polyadenylation. Cell. 2019;178:91–106.e23.

1019 34. Vaishnav ED, de Boer CG, Molinet J, Yassour M, Fan L, Adiconis X, et al. The evolution,
1020 evolvability and engineering of gene regulatory DNA. Nature. 2022;603:455–63.

1021  35. Zhou J. Sequence-based modeling of three-dimensional genome architecture from
1022  kilobase to chromosome scale. Nat Genet. 2022;54:725–34.

1023  36. Fudenberg G, Kelley DR, Pollard KS. Predicting 3D genome folding from DNA sequence
1024  with Akita. Nat Methods. 2020;17:1111–7.

1025  37. Schwessinger R, Gosden M, Downes D, Brown RC, Oudelaar AM, Telenius J, et al.
1026  DeepC: predicting 3D genome folding using megabase-scale transfer learning. Nat Methods.
1027  2020;17:1118–24.

1028  38. Avsec Ž, Kreuzhuber R, Israeli J, Xu N, Cheng J, Shrikumar A, et al. The Kipoi repository
1029  accelerates community exchange and reuse of predictive models for genomics. Nat
1030  Biotechnol. 2019;37:592–600.

1031  39. Athar A, Füllgrabe A, George N, Iqbal H, Huerta L, Ali A, et al. ArrayExpress update -
1032  from bulk to single-cell expression data. Nucleic Acids Res. 2019;47:D711–5.

1033  40. Stovner EB, Sætrom P. PyRanges: efficient comparison of genomic intervals in Python.
1034  Bioinformatics. 2020;36:918–9.

1035  41. Seabold S, Perktold J. Statsmodels: Econometric and statistical modeling with python.
1036  Proceedings of the 9th Python in Science Conference [Internet]. SciPy; 2010. Available from:
1037  https://conference.scipy.org/proceedings/scipy2010/seabold.html