



Synthetic STARR-seq reveals how DNA shape and sequence modulate transcriptional output and noise

Stefanie Schöne, Melissa Bothe, Edda Einfeldt, Marina Borschiwer, Philipp Benner, Martin Vingron, Morgane Thomas-Chollier, Sebastiaan H Meijsing

► To cite this version:

Stefanie Schöne, Melissa Bothe, Edda Einfeldt, Marina Borschiwer, Philipp Benner, et al.. Synthetic STARR-seq reveals how DNA shape and sequence modulate transcriptional output and noise. PLoS Genetics, 2018, 14 (11), pp.e1007793. 10.1371/journal.pgen.1007793 . inserm-02155974

HAL Id: inserm-02155974

<https://inserm.hal.science/inserm-02155974>

Submitted on 14 Jun 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

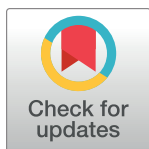
RESEARCH ARTICLE

Synthetic STARR-seq reveals how DNA shape and sequence modulate transcriptional output and noise

Stefanie Schöne¹, Melissa Bothe¹, Edda Einfeldt¹, Marina Borschiwer¹, Philipp Benner¹, Martin Vingron¹, Morgane Thomas-Chollier², Sebastiaan H. Meijsing^{1*}

1 Max Planck Institute for Molecular Genetics, Berlin, Germany, **2** Institut de biologie de l'Ecole normale supérieure (IBENS), Ecole normale supérieure, CNRS, INSERM, PSL Université Paris, Paris, France

* meijsing@molgen.mpg.de



OPEN ACCESS

Citation: Schöne S, Bothe M, Einfeldt E, Borschiwer M, Benner P, Vingron M, et al. (2018) Synthetic STARR-seq reveals how DNA shape and sequence modulate transcriptional output and noise. *PLoS Genet* 14(11): e1007793. <https://doi.org/10.1371/journal.pgen.1007793>

Editor: Timothy E. Reddy, Duke University, UNITED STATES

Received: June 14, 2018

Accepted: October 26, 2018

Published: November 14, 2018

Copyright: © 2018 Schöne et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Data were deposited in ArrayExpress under the accession numbers: E-MTAB-6738 (RNA-seq U2OS-GR18) and E-MTAB-6737 (synSTARR-seq U2OS-GR18). Reviewers access to datasets: STARR-seq data: E-MTAB-6737 Username: Reviewer_E-MTAB-6737 Password: hgeofcho RNA-seq data: E-MTAB-6738 Username: Reviewer_E-MTAB-6738 Password: cieef7tt

Funding: This work was supported by the Deutsche Forschungsgemeinschaft [ME4154/1-1]

Abstract

The binding of transcription factors to short recognition sequences plays a pivotal role in controlling the expression of genes. The sequence and shape characteristics of binding sites influence DNA binding specificity and have also been implicated in modulating the activity of transcription factors downstream of binding. To quantitatively assess the transcriptional activity of tens of thousands of designed synthetic sites in parallel, we developed a synthetic version of STARR-seq (synSTARR-seq). We used the approach to systematically analyze how variations in the recognition sequence of the glucocorticoid receptor (GR) affect transcriptional regulation. Our approach resulted in the identification of a novel highly active functional GR binding sequence and revealed that sequence variation both within and flanking GR's core binding site can modulate GR activity without apparent changes in DNA binding affinity. Notably, we found that the sequence composition of variants with similar activity profiles was highly diverse. In contrast, groups of variants with similar activity profiles showed specific DNA shape characteristics indicating that DNA shape may be a better predictor of activity than DNA sequence. Finally, using single cell experiments with individual enhancer variants, we obtained clues indicating that the architecture of the response element can independently tune expression mean and cell-to-cell variability in gene expression (noise). Together, our studies establish synSTARR as a powerful method to systematically study how DNA sequence and shape modulate transcriptional output and noise.

Author summary

The expression level of genes is controlled by transcription factors, which are proteins that bind to genomic response elements that contain their recognition DNA sequence. Importantly, genes are not simply turned on but need to be expressed at the right level. This is, at least in part, assured by the sequence composition of genomic response elements. Here, we studied how the recognition DNA sequence influences gene regulation by a transcription factor called the glucocorticoid receptor. Specifically, we developed a

to SS]. And by the Max-Planck-Gesellschaft (to SS, MB, EE, MB, PB, MV and SHM. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

method to test the activity of variants in a highly parallelized setting where everything is kept identical except for the sequence of the binding site. The systematic analysis of tens of thousands of sequence variants facilitated the identification of a previously unknown sequence variant with high activity. Moreover, we report how sequence variation of the response element influences cell-to-cell variability in expression levels. Finally, we observe similar activity profiles for distinct sequence variants that share similar three-dimensional DNA shape characteristics arguing that the three-dimensional perception of DNA by the glucocorticoid receptor, modulates its activity towards individual target genes.

Introduction

The interplay between transcription factors (TFs) and genomically encoded *cis*-regulatory elements plays a key role in specifying where and when genes are expressed. In addition, the architecture of *cis*-regulatory elements influences the expression level of individual genes. For example, transcriptional output can be tuned by varying the number of TF binding sites, either for a given TF or for distinct TFs, present at an enhancer [1, 2]. Moreover, differences in its DNA-binding sites can modulate the magnitude of transcriptional activation, as exemplified by the glucocorticoid receptor (GR), a hormone-activated TF [3–5]. The sequence differences can reside within the 15 base pair (bp) core GR binding sequence (GBS) consisting of two imperfect 6 bp palindromic half-sites separated by a 3 bp spacer. Although the effects on activity are more modest than those observed for changes within the core, sequences directly flanking the core also modulate GR activity [3]. However, these sequence-induced changes in activity cannot be explained by affinity [3]. Instead, the flanking nucleotides induce structural changes in both DNA and the DNA binding domain of GR, arguing for their role in tuning GR activity [3].

Notably, the expression level of a gene is typically measured for populations of cells and thus masks that expression levels can vary considerably between individual cells of an isogenic population [6–9]. This variability in the expression level of a gene, called expression noise, results in phenotypic diversity, which can play a role in organismal responses to environmental changes (so called bet-hedging) and in cell fate decisions during development. Expression noise can be explained by the stochastic nature of the individual steps that decode the information encoded in the genome. For example, transcription occurs in bursts [7, 10–12], which can induce variability in gene expression due to differences in burst frequency and in the number of transcripts generated per burst (burst size) [13]. Noise levels are gene-specific, which can be explained in part by differences in the sequence composition of *cis*-regulatory elements [11, 14–16]. For instance, the sequence composition of promoters influences expression variability with high burst size and noise for promoters containing a TATA box [15, 17]. In addition, chromatin and the presence or absence of nucleosome-disfavoring sequences have been linked to transcriptional noise [16–19]. Finally, noise levels can also be tuned by the number and by the affinity of TF binding sites [11, 16].

Many fundamental insights regarding the role of sequence in tuning transcriptional output and noise have come from reporter studies [20, 21]. A key advantage of reporters is that they can provide quantitative information in a controlled setting where everything is kept identical except for the sequence of the region of interest. Until recently, a limitation of reporter studies was that sequence variants had to be tested one at a time. However, the recent development of several parallelized reporter assays allows the simultaneous assessment of many sequence variants [21]. One of these parallelized methods is STARR-seq (Self-Transcribing Active

Regulatory Region sequencing) [22]. In this assay, candidate sequences are placed downstream of a minimal promoter, such that active enhancers drive their own expression and high-throughput sequencing reveals both the sequence identity and quantitative information regarding the activity of each sequence variant. The STARR-seq method has been used to assay enhancer activity genome-wide [22, 23], to study regions of interest isolated either by Chromatin Immunoprecipitation (ChIP) or a capture-based approach [24, 25], and to study the effect of hormones on enhancer activity [25, 26].

Here, we adapted the STARR-seq method to systematically study how sequence variation both within the 15 bp GBS and in the region directly flanking it modulate GR activity. Specifically, we generated STARR-seq libraries using designed synthetic oligos (synSTARR-seq) with randomized nucleotides flanking the core GBS to show that the flanks modulate transcriptional output by almost an order of magnitude. When grouping sequences based on their ability to either enhance or blunt GBS activity, we found that each group contained a broad spectrum of highly diverse sequences, but striking similarities in their DNA shape characteristics. Using the same approach, we also assayed the effect of sequence variation within the core GBS. Finally, using single cell experiments with individual enhancer variants, we study how the sequence composition of the response element influences expression mean and noise. Together, our studies establish synSTARR-seq as a powerful method to study how DNA sequence and shape modulate transcriptional output and noise.

Results

Measuring the activity of thousands of GR binding sequence variants in parallel using the synSTARR-seq approach

To test if we could use the STARR-seq reporter [22] to study how sequence variation of the GR binding site influences GR activity, we first tested if a single GBS is sufficient to facilitate GR-dependent transcriptional activation of the reporter. Therefore, we constructed STARR reporters containing either a single GBS as candidate enhancer (Fig 1A), a randomized sequence or as positive control a larger GBS-containing sequence derived from a GR-bound region close to the GR target gene *FKBP5*. The resulting reporters were transfected into U2OS cells stably expressing GR (U2OS-GR) [27] and their response to treatment with dexamethasone (dex), a synthetic glucocorticoid hormone, was measured. As expected, no marked hormone-dependent induction was observed for the reporter with the randomized sequence. This was true both at the level of RNA (Fig 1B) and at the level of the GFP reporter protein (S1 Fig). In contrast, we observed a robust hormone-dependent activation both at the level of RNA and GFP protein for reporters with either a single GBS or with the larger genomic *FKBP5* fragment (Fig 1B and S1A Fig), showing that a single GBS is sufficient for GR-dependent activation of the STARR-seq reporter.

Our previous work has shown that the sequence directly flanking GBSs can modulate DNA shape and GR activity [3]. For a parallelized and thorough analysis of sequence variants flanking a GBS, we generated STARR-seq libraries for two GBS variants, we previously named Cgt and Sgk, that showed a strong influence of flanking nucleotides on activity [3]. Specifically, we generated libraries using designed synthetic sequences (synSTARR-seq) containing a GBS with five consecutive randomized nucleotides directly flanking the imperfect half site (Fig 1A and S2A Fig). Next, we transfected the GBS flank libraries into U2OS-GR cells to determine the activity of each of the 1024 flank variants present in the library. We performed three biological replicates for each condition and found that the results were highly reproducible ($r \geq 0.91$ for vehicle treated cells, $r \geq 0.98$ for dex treated cells; Fig 1C and S1B–S1E Fig). Notably, we retain duplicate reads in our analysis, which is essential to get quantitative information

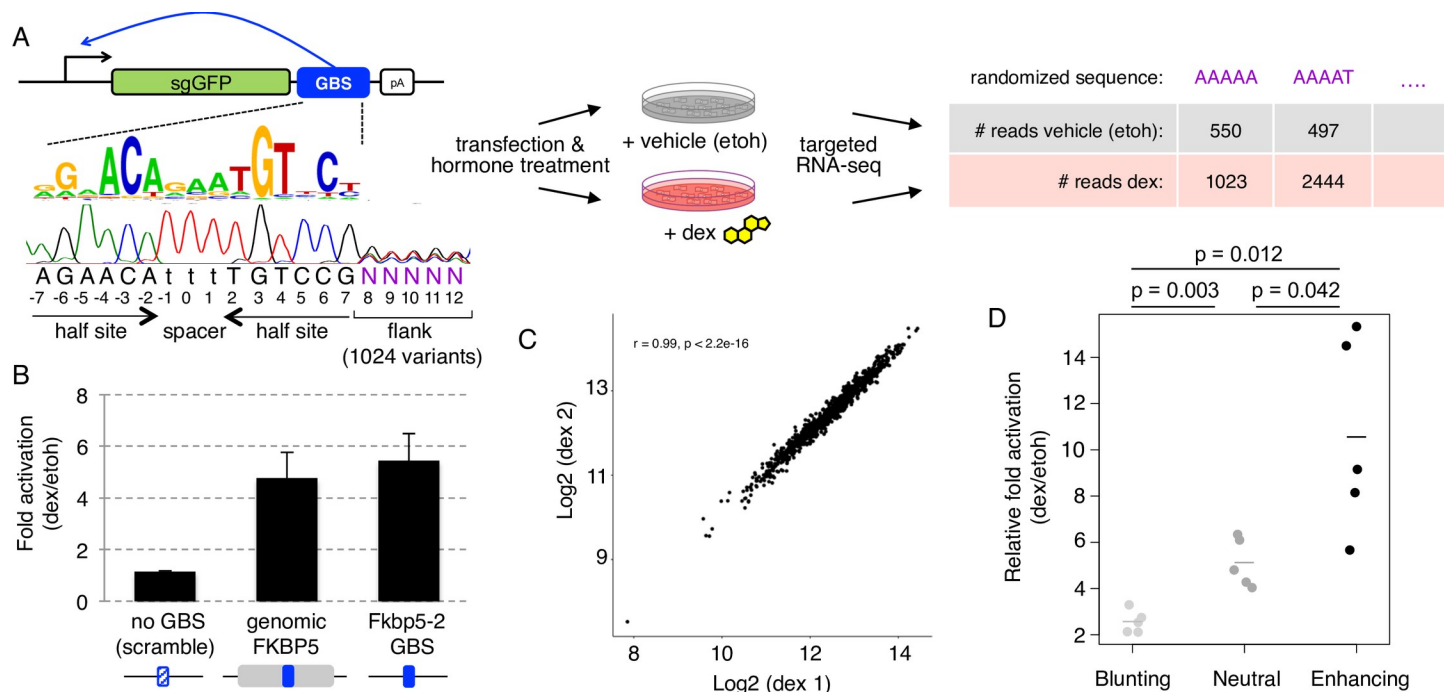


Fig 1. Design and validation of the synSTARR-seq approach. (a) SynSTARR-seq reporter setup using a synthetic library containing a GR Binding Sequence (GBS) flanked by 1024 different flanking sequences (flank library) to screen for flanks that modulate GBS activity. Samples are treated with dexamethasone (dex) or ethanol vehicle (etoh) before targeted RNA-sequencing and counting of the reads. (b) Transcriptional activation of STARR-seq reporter containing candidate enhancer inserts as indicated. Mean fold change upon dexamethasone treatment \pm S.D. ($n = 3$) in U2OS-GR cells is shown. Genomic FKBP5 (211bp region hg19Chr6: 35699789–35699999); FKBP5-2 GBS (single GBS: AGAACAtccGTGCGC); no GBS (AGAAACtccGTGCGC). (c) Representative RNA-seq correlation plot for biological replicates of dexamethasone-treated cells (4h, 1 μ M) transfected with the GBS-flank library. (d) The enhancer activity of blunting ($n = 5$), neutral ($n = 5$) and enhancing ($n = 5$) flank variants was assessed for individually transfected STARR-seq constructs by qPCR. Fold change upon dexamethasone treatment normalized to the activity for the scrambled control plasmid is shown as horizontal line for the mean of each activity group and as dot for each individual construct.

<https://doi.org/10.1371/journal.pgen.1007793.g001>

for individual sequence variants of the library. To calculate the activity for each flank variant, we used DESeq2 [28] to compare the RNA-seq read number between dex- and vehicle (ethanol) treated cells (Fig 1A). This resulted in the identification of 189 flank variants with significantly higher activity (enhancing flanks), 125 flank variants with significantly lower activity (blunting flanks) and 710 flank variants that did not induce significant changes in activity (neutral flanks). To test the accuracy of the synSTARR-seq data, we cloned 5 flank variants from each activity group (enhancing, blunting and neutral) and assayed the activity of each variant individually by qPCR. Consistent with what we observed for the synSTARR library, the activity of blunting flanks was significantly lower than for the neutral flanks whereas the activity of the enhancing flanks was significantly higher (Fig 1D). Notably, all flank variants tested were activated upon dex treatment ranging from 2.1 to 15.3 fold (627% higher) depending on the sequence of the flank. Together, our results show that the synSTARR-seq assay produces reproducible and quantitative information and can be used for a high-throughput analysis of the effect of the flanking sequence on GBS activity.

SynSTARR-seq to assay the effects of flanking nucleotides

To assess how the sequence composition of the flanking region influences GBS activity, we ranked the flank variants by their activity and used a color chart representation to plot the sequence at each position for the Cgt (Fig 2A) and Sgk GBS (S2A Fig), respectively. In addition, we generated consensus sequence motifs for the significantly enhancing and blunting

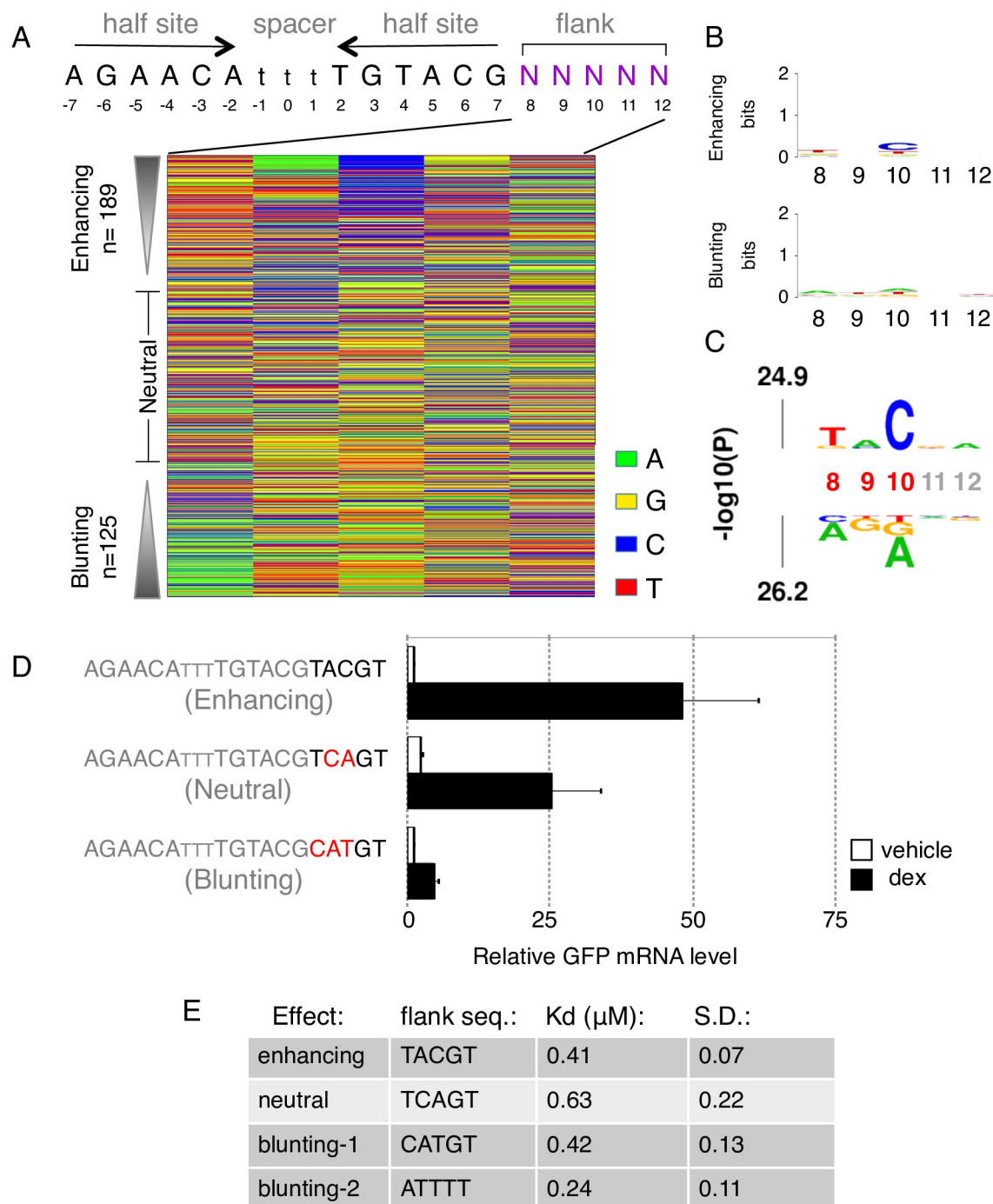


Fig 2. Analysis of the GBS flank library. (a) Color chart summarizing the sequence at each variable position for flank variants ranked by their fold change in response to hormone treatment. (b) Consensus motif for (top) significantly (adjusted p-value < 0.01) enhancing and (bottom) blunting flank variants. (c) kpLogo probability logo (activity logo) for flank variants depicting the p-values from Mann-Whitney U tests of whether GBS variants with a specific nucleotide at a given position are more (displayed above number indicating nucleotide position) or less (displayed below number indicating nucleotide position) active than other GBS variants. Positions with significant nucleotides (p < 0.001) are highlighted (red coordinates). (d) Transcriptional activity of STARR-seq reporters containing candidate flank variants as indicated. Relative RNA levels ± S.E.M. are shown for cells treated with ethanol vehicle and for cells treated overnight with 1 μM dexamethasone (n ≥ 3). (e) Table of EMSA-derived DNA-binding constants (Kd) for flank variants as indicated ± S.D. (n ≥ 3).

<https://doi.org/10.1371/journal.pgen.1007793.g002>

variants (Fig 2B and S2B Fig). Notably, these consensus sequence motifs treat each sequence equally and do not take the quantitative information regarding the activity of each sequence into account. To take advantage of the quantitative information provided by the synSTARR-seq assay, we used *kpLogo* [29], which uses the fold change as weight for each sequence variant, and statistically evaluates the enrichment/depletion of specific nucleotides at each position. The resulting probability logo can be interpreted as an activity logo that visualizes for each position which nucleotides are associated with either higher (letters above the coordinates) or lower (below the coordinates) GBS activity (Fig 2C and S2C Fig). The activity logo, consensus motifs and color chart highlight several sequence features for enhancing and blunting flank variants. For example, high activity is associated with a T at position 8 for both the Cgt and Sgk GBS, which matches what we found previously when we studied the activity of endogenous GR-bound regions [3]. In addition, the most active flank variants preferentially have an A at position 9 followed by a C at position 10 (Fig 2A and S2A Fig). To validate that this “TAC” signature results in high activity, we shuffled the sequence to either TCA or CAT and found that this indeed resulted in markedly lower activity (Fig 2D). For blunting flank variants, we observed a preference for an A at position 8 and a bias against having a C at position 10 (Fig 2A and 2C and S2A and S2C Fig). However, altogether we find that the consensus motifs for enhancing and blunting flanks only have low information content and that a broad spectrum of distinct sequences can enhance or blunt the activity of the adjacent GBS (Fig 2B and S2B Fig).

Our previous work [3] indicates that DNA shape can influence GR activity downstream of binding. Consistent with this notion, we measured similar *K_d* values for flanks variants from the different activity classes (Fig 2E). These findings are also in agreement with published work showing that the nucleotides directly flanking GBSs have little effect on GR affinity [30]. To examine if the flank effects might be explained by differences in DNA shape, we calculated the predicted minor groove width, roll, propeller twist and helix twist [31] for enhancing and blunting flank variants (Fig 3A and S2D Fig and S3 Fig). Consistent with a role for DNA shape in modulating GR activity, we found shape characteristics that differ between enhancing and blunting flanks. For example, we observed a wider minor groove at position 6, and to a lesser degree at position 7 for blunting flanks of the Cgt GBS, when compared to enhancing flanks (Fig 3A and S4A Fig). In addition, blunting flanks for the Cgt GBS have a narrower minor groove than enhancing flanks for positions 8–12 (Fig 3A and S4A Fig), a region with several non-specific minor groove contacts with the C-terminal end of the DNA binding domain of GR [5]. For the Sgk GBS library, we find similar shape characteristics associated with blunting flanks with a wider minor groove at position 6 and a narrower minor groove for positions 8–12 (S2D Fig and S4B Fig). DNA-shape-based hierarchical clustering recapitulates these characteristics in cluster 4, containing many more blunting flanks than any of the other clusters, for both the Cgt and Sgk GBS flank libraries (Fig 3B and 3C and S2E and S2G Fig). Of note, the consensus motifs for cluster 4 and for the other shape clusters have only low information content (Fig 3D and S2F Fig) indicating that distinct sequences can give rise to similar shape characteristics with shared effects on the activity of the adjacent GBS.

Together, these synSTARR-seq experiments uncover how sequence variation in the flanking region of the GBS influences activity and point at a role for DNA shape in modulating GBS activity.

SynSTARR-seq to assay the effects of variation within the core GBS

We next generated an additional synSTARR-seq library to study the effect of variation within the 15bp core sequence. This library contains a fixed GBS half site followed by eight

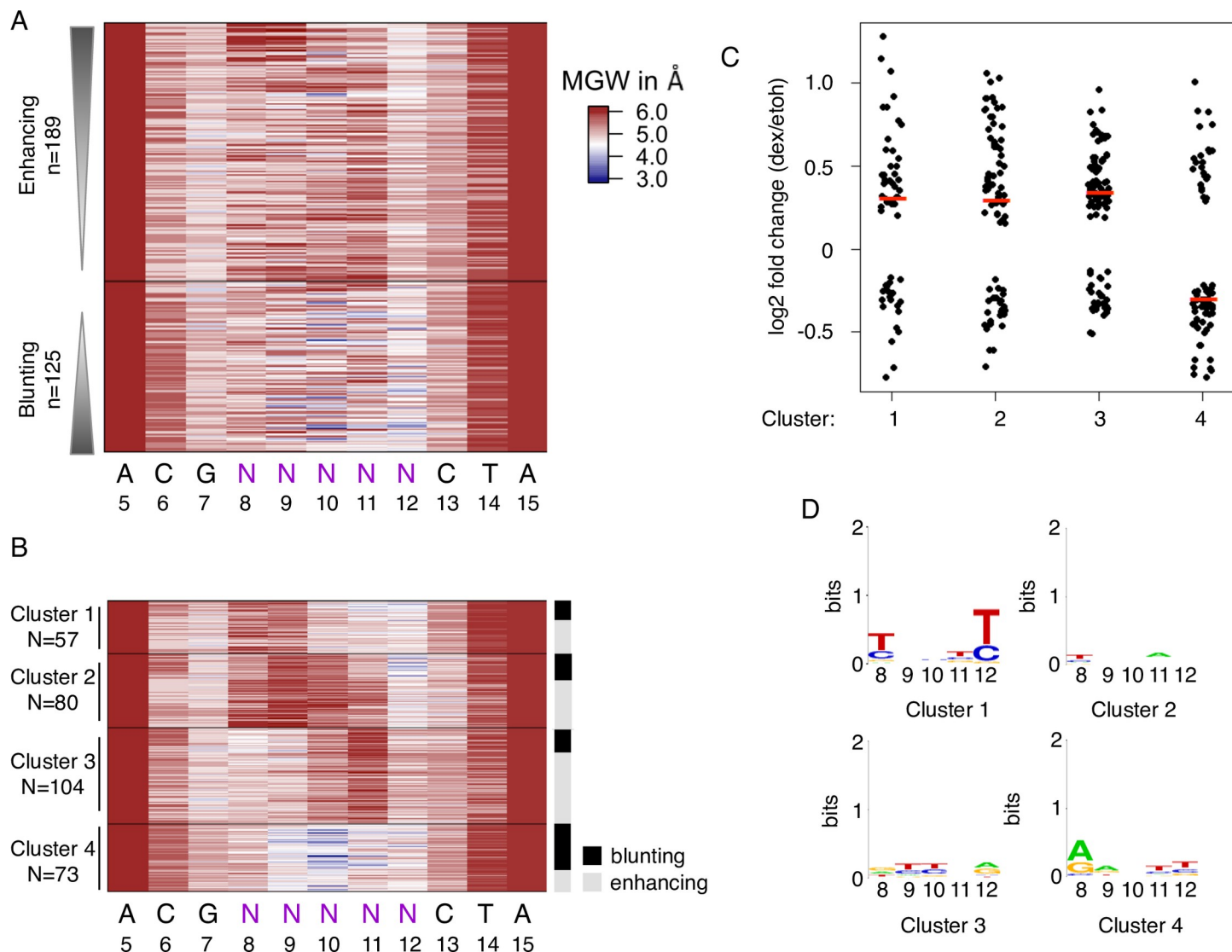


Fig 3. Predicted DNA shape for enhancing and blunting flank variants. (a) Predicted minor groove width (MGW) for significant enhancing and blunting flank variants of the Cgt GBS library ranked by their fold change in response to hormone treatment. (b) K-means clustering based on MGW for significantly enhancing and blunting flank variants. Right side: activating and blunting variants are highlighted in grey and black respectively. (c) Log2 fold change upon dexamethasone treatment for each cluster as indicated. The synSTARR-seq activity for individual sequences is shown as black dots, the median for each cluster as a horizontal red line. (d) Consensus sequence motif for clusters as indicated.

<https://doi.org/10.1371/journal.pgen.1007793.g003>

consecutive randomized nucleotides (Fig 4A). The library, containing over 65,000 variants, was transfected into U2OS-GR cells and the read count for each variant was determined both in the presence and absence of hormone treatment. Compared to the flank library, we observed a lower correlation between experiments, especially for variants with a low read count (S5 Fig). Specifically, when we compared the read count between biological replicates, we found that sequences with a read count below 100 were typically detected in only one of the replicates. Therefore, we decided to remove sequences with a mean read count below 100 across all experiments. Next, we analyzed data from three biological replicates to determine the activity of variants in the library (Fig 4B). To validate the measured activities, we cloned 4 sequences that repress, 4 that show a weak activation (\log_2 fold change < 2) and 8 strongly

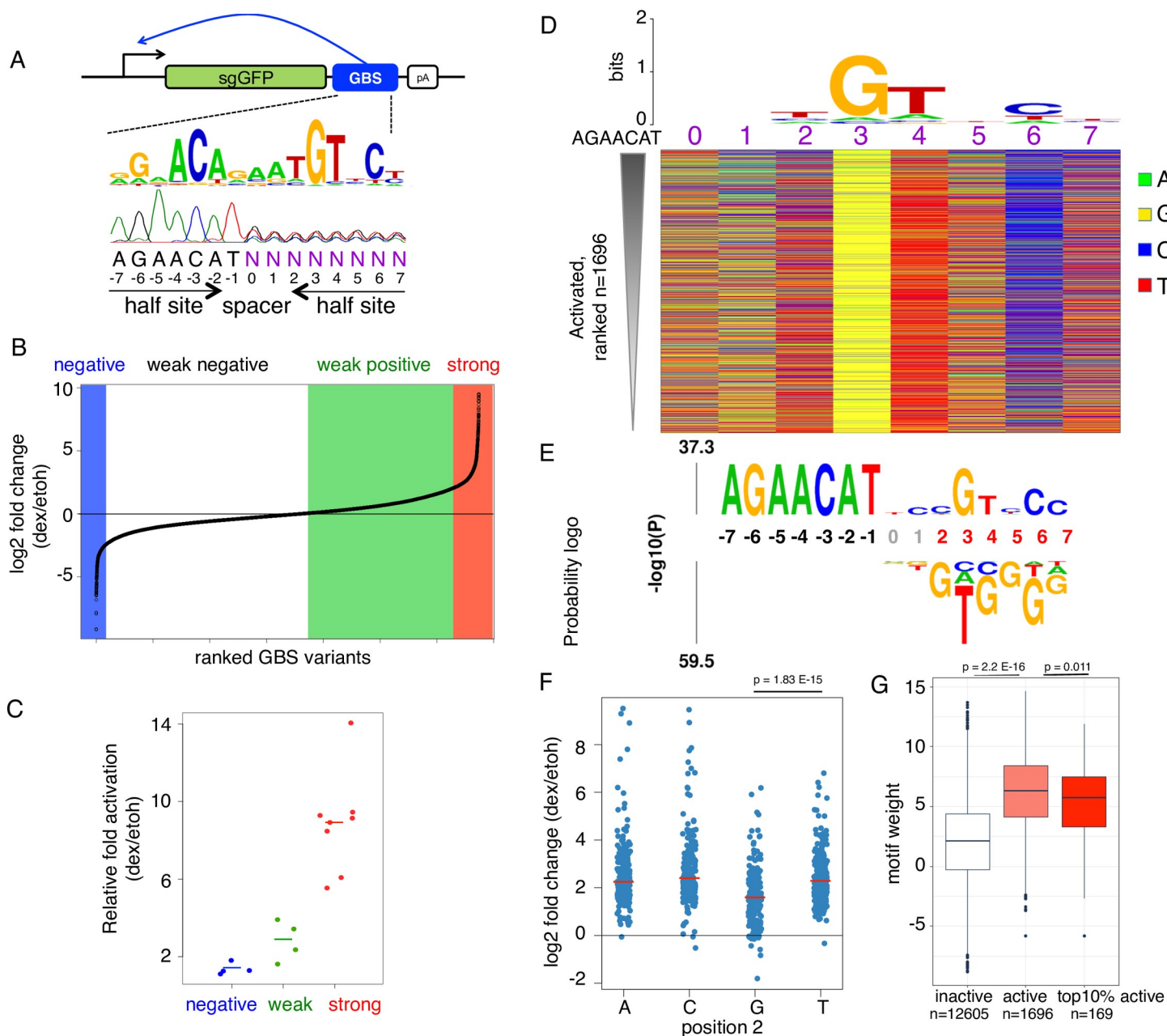


Fig 4. Analysis of the GBS half site library. (a) SynSTARR-seq reporter setup using a synthetic library containing 65,536 candidate GR Binding Sequence (GBS) variants (half site library with 8 variable positions N). (b) Candidate GBS variants were ranked by their fold change in expression in response to hormone treatment (4 h, 1 μ M dex). Only sequences with a mean read count > 100 across all replicates ($n = 3$) for both dex and ethanol vehicle treated cells are shown. Repressed ($\log_2 FC < -2$), weakly active ($0 < \log_2 FC < 2$) and activated GBS variants ($\log_2 FC \geq 2$) are highlighted by a blue, green and red background respectively. (c) The enhancer activity of negative ($n = 4$), weak ($n = 4$) and strong ($n = 8$) GBS variants was assessed by qPCR for individually transfected STARR-seq constructs. Fold change upon dexamethasone treatment normalized to the activity for the scrambled control plasmid is shown. Horizontal line shows the mean for each activity group; dots the values for individual constructs. (d) Top: Consensus motif and below a color chart summarizing the sequence at each variable position for each significantly activated GBS variant (adjusted p-value < 0.01) ranked by their fold change in response to dex treatment. (e) kLogo probability logo (activity logo) for half site variants depicting the p-values from Mann-Whitney U tests of whether GBS variants with a specific nucleotide at a given position are more (displayed above number indicating nucleotide position) or less (displayed below number indicating nucleotide position) active than other GBS variants. Positions with significant nucleotides ($p < 0.001$) are highlighted in red, fixed positions in black. (f) Log₂ fold change upon dexamethasone treatment for GBS-like variants with either an A, C, G or T at position 2 (exact match to AGAACATnnXGTnCN, with X either A, C, G or T). Data for individual sequences are shown as blue dots. Horizontal red lines show the median for each group. p-values were calculated using a Student's t-test. (g) Boxplot of the motif weight (using the truncated 15nt long M00205 motif from Transfac) for inactive ($-0.5 \leq \log_2 \text{fold change} \leq 0.5$; white), active (light red) and the top 10% active (dark red) GBS variants. p-values were calculated using a Student's t-test.

<https://doi.org/10.1371/journal.pgen.1007793.g004>

activating GBS variants. Consistent with the results from our screen, the three groups showed distinct levels of activity (Fig 4B and 4C). However, for the group of repressed GBS variants we did not recapitulate the observed repression in our screen (Fig 4C), indicating that these variants might behave differently in isolation. Alternatively, what looks like repression might be a consequence of issues with data normalization, which assumes that the distribution of the log fold changes is centered on 0, which is not given when GBS variants can activate but not repress gene expression. Notably, a lack of GR-dependent transcriptional repression was also reported in another study using the STARR-seq approach to study the regulatory activity of GR-bound genomic regions [25] indicating that GR might not be able to repress transcription in the STARR-seq context.

Given that the observed repression was not reproducible, we concentrated our analysis on 1696 sequences that facilitated significant GR-dependent transcriptional activation. Consistent with activation, we found that the consensus motif for activating sequence variants recapitulates the known GR consensus sequence with the second half site 3-bp downstream of the fixed first half site of our library (Fig 4D). Accordingly, the GBS motif weight, which serves as a proxy for DNA binding affinity, is higher for activating sequences when compared to sequences that did not respond to hormone treatment (Fig 4G). However, the score for the top 10% most active sequences was not higher than for all active variants (Fig 4G), arguing that higher affinity does not drive the high levels of activation. As expected and consistent with the GR consensus motif, the color chart (Fig 4D) and activity logo (Fig 4E) highlight a strong preference for a G at position 3 and accordingly GBS activity is significantly lower for variants with a nucleotide other than G at this position (S6A Fig). The activity logo also highlights that a G at position 2 is associated with lower activity (Fig 4E and 4F).

Previous studies have shown that the sequence of the spacer can modulate GBS activity [4, 5]. Therefore, we compared the activity of all 16 spacer variants in our library that match the GBS consensus for the second half site at the key positions 3, 4 and 6 (S7A Fig). In line with a role for the spacer in modulating transcriptional output, we find significant differences between the spacer variants (S7B Fig). For example, the activity for variants with an AC spacer is significantly higher than for several other spacer variants (S7B Fig) whereas the activity for GT variants is significantly lower ($p_{\text{adj}} < 0.01$) than either AA, AC or TC variants (S7B Fig).

Unexpectedly, the activity logo and top of the color chart indicated a high activity for variants with a C at position 2 (Fig 4D and 4E), instead of a consensus T observed in the GR consensus motif and from *in vitro* experiments studying the effect of DNA sequence on GR DNA binding affinity [30]. A careful examination of the sequence composition of the most active variants also revealed a preference for TC at the preceding positions within the spacer (Figs 4E and 5A). To test if the high activity for sequences with a C at position 2 depends on the nucleotide composition of the preceding nucleotides, we changed them to GG and found that this resulted in a marked reduction in GR-dependent activation (Fig 5B and S8A Fig). In addition, we compared the activity between variants with a T or a C at position 2. The activity was higher for the C variant when preceded by TC. However, when we changed the preceding nucleotides to GG the activation was stronger for the T than the C variant (Fig 5B and S8A Fig). These experiments indicated that the high activity for the C variant depends on the preceding nucleotides.

Interestingly, the most active variants resemble the sequence composition of the “combi” motif we identified previously [32]. The combi motif contains only a single GR half site followed by TTCC and we found evidence that GR binds this sequence as a monomer in conjunction with a partnering protein [32]. Similar to the combi motif, several of the most active variants (Fig 5A) contain a GR half site followed by TTCC. However, whereas the combi motif lacks a second GR half site, the motif for the 25 most active variants from our screen (named

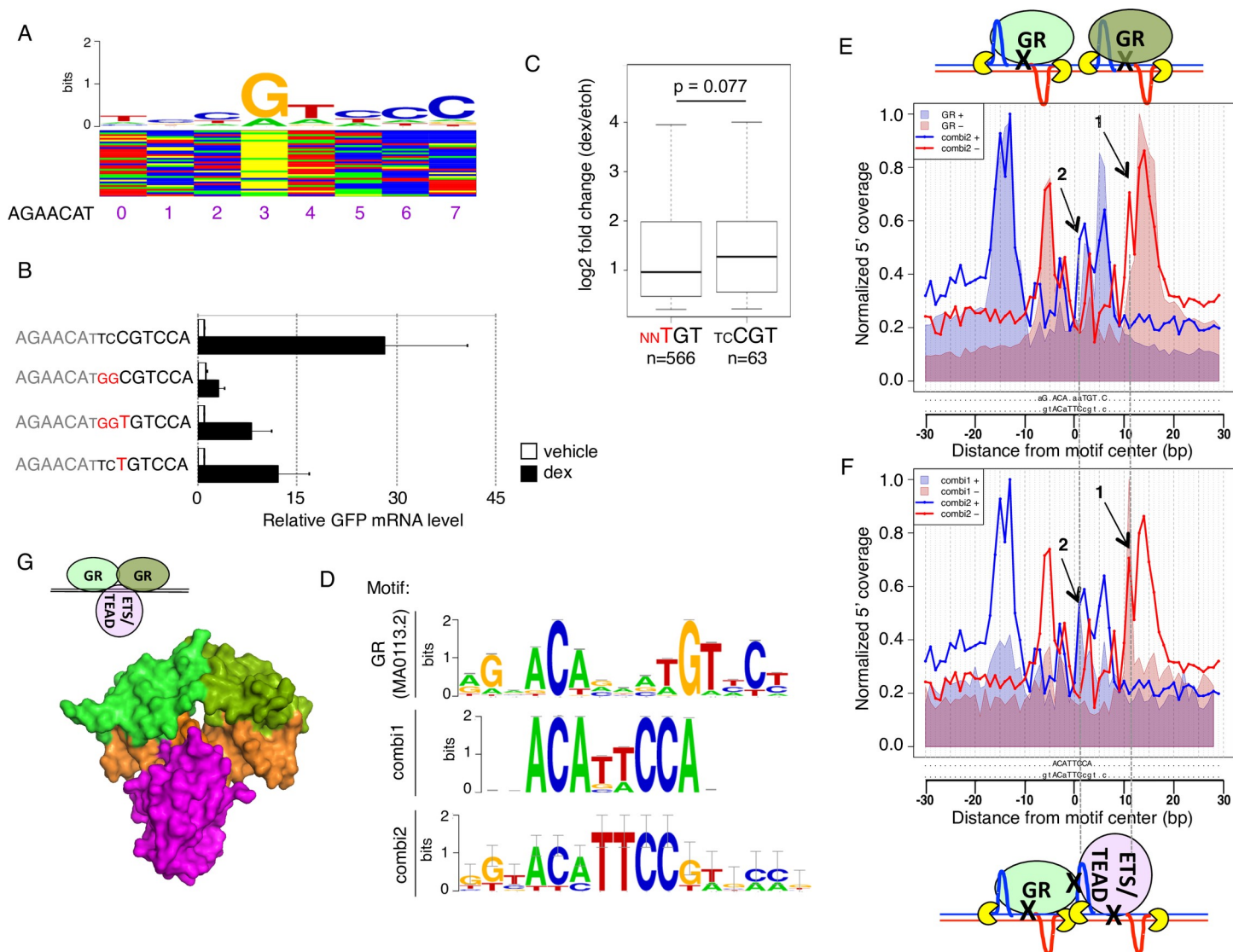


Fig 5. Identification and characterization of the combi2 motif. (a) Color chart for the top activated GBS variants and above the consensus motif for the 25 most active sequences (b) Transcriptional activity of STAR-seq reporters containing candidate GBS variants as indicated. Relative RNA levels \pm S.E.M. are shown for cells treated with ethanol vehicle and for cells treated overnight with 1 μ M dexamethasone ($n = 3$). (c) Boxplot of the log2 fold change upon treatment for 4 h with 1 μ M dexamethasone for genes with a ChIP-seq peak in the region ± 20 kb around the TSS containing either a conventional GBS match (M00205; p value < 0.0001) or a combi2-like sequence (combi2 motif; p value < 0.0001). Center lines show the median. p -value was calculated using the Wilcoxon rank sum test. (d) Motif logo representing the positional weight matrices for the canonical GBS (JASPAR MA0113.2), combi1 and combi2 motif. (e) Alignment of the ChIP-exo footprint profiles for the combi2 (solid lines, blue for the positive strand, red for the negative strand) and the canonical GBS motif (shaded area, blue for positive strand, red for negative strand). Arrows 1 and 2: Additional 5' coverage for the combi2 motif that does not match the conventional GBS footprint. (f) Alignment of the ChIP-exo footprint profiles for the combi1 (shaded area, blue for positive strand, red for negative strand) and combi2 motif (solid lines, blue for the positive strand, red for the negative strand). Arrows 1 and 2: Additional 5' coverage for the combi2 motif when compared to the canonical GBS footprint aligns with signal for the combi1 motif. (g) Structural alignment of combined binding of a GR dimer (green) and ETS1 (purple, middle: PDB 1K79) at the combi2 sequence (orange).

<https://doi.org/10.1371/journal.pgen.1007793.g005>

“combi2”) also contains a recognizable second GR half site (Fig 5A). To gain insight into the mode of GR binding at the combi2 motif, we examined published ChIP-exo data [32]. ChIP-exo is an assay that combines ChIP with a subsequent exonuclease step [33] which results in a base-pair resolution picture of GR binding. The ChIP-exo signal takes the form of sequence-specific peak patterns (footprint profiles), detectable on both strands with the program ExoProfiler [32]. We applied ExoProfiler to scan GR-bound regions with the combi2 motif (Fig

5D and 5E, solid lines). As control, we analyzed the footprint profile for the canonical GR consensus motif (Fig 5D; JASPAR MA0113.2) and recovered peak pairs on the forward and reverse flanks that demarcate the protection provided by each of the monomers of the GR dimer (Fig 5E, shaded area). The signal for the first half site is essentially the same and a similar pattern is also observed for the second half site, indicating that GR binds as a dimer on regions bearing the combi2 motif, however with additional signal (highlighted with black arrows in Fig 5E). In addition, we compared the footprint profile between the original combi (Fig 5D; [32]) and the combi2 motif (Fig 5F). Again, the position and shape of the peaks are compatible for the first half site but the ChIPexo signal for the second half site looks markedly different. The aforementioned additional signal for the combi2 motif aligns with the position of the second peak pair of the combi motif (Fig 5F), indicating that the footprint profile for the combi2 motif appears to be a composite of the signal for homodimeric GR binding at canonical GBSs and the signal for monomeric GR binding together with another protein. Our previous work suggests that this partnering protein on combi motif might be Tead or ETS2. The ChIP-exo profile thus points to three alternative binding configurations on combi2: homodimeric GR, monomeric GR binding with Tead/ETS2 or the simultaneous binding of homodimeric GR complex together with Tead/ETS2. Structural modeling suggests that this third mode is possible given the absence of obvious sterical clashes that would prevent this mode of binding (Fig 5G). However, additional functional studies are needed to determine if GR indeed partners with Tead/ETS2, or possibly with other proteins, at the combi2 motif.

To assess if DNA shape could play a role in modulating GBS activity, we calculated the predicted minor groove width for all 1696 significantly activated sequences ranked by activity (S6B Fig). Comparison of the top 20% most active and bottom 20% least active sequence variants highlighted two regions with significant differences. First, consistent with our findings for the flank library, we find that a wider minor groove at positions 6 and 7 correlates with weaker activity (S6B and S6C Fig). Second, we find that a narrower minor groove in the spacer (position -1 and 0) correlates with weaker activity (S6B and S6C Fig). As we observed for the flank variants, the different activity classes do not show a distinct sequence signature (S6B Fig) again arguing that DNA shape might modulate GBS activity.

Together, the findings for our half site library suggest a role for both DNA shape and sequence in tuning the activity of GBS variants. Moreover, our screen uncovered a novel high-activity functional GR binding sequence variant.

SynSTARR to assay the effects of enhancer sequence composition on noise

Thus far, we analyzed the effect of sequence composition on transcriptional output by analyzing mean expression levels for populations of cells. To test if sequence variation in the enhancer influences cell-to-cell variability in gene expression (noise), we measured GFP levels for individual STARR constructs in single cells (Fig 6A and 6B). Cells were transfected with individual constructs along with an mCherry expression construct to remove extrinsic noise, for example caused by differences in transfection efficiency. We first analyzed sequence variants containing a single GBS (single GBS group) including known GBSs; two variants matching the combi2 sequence motif and the Cgt GBS with an enhancing flank variant. Consistent with previous findings [5], we found that GBS variants from the single GBS group induced different mean levels of GFP expression. For example, the mean GFP level upon dex treatment was lower for the Pal GBS than for the Cgt variant (Fig 6C, orange and red squares). In line with findings by others [16], we observed that transcriptional noise scales with mean expression with lower noise for variants with higher mean expression (Fig 6C). Next, we assayed two additional groups of sequences with distinct binding sites architectures that both result in

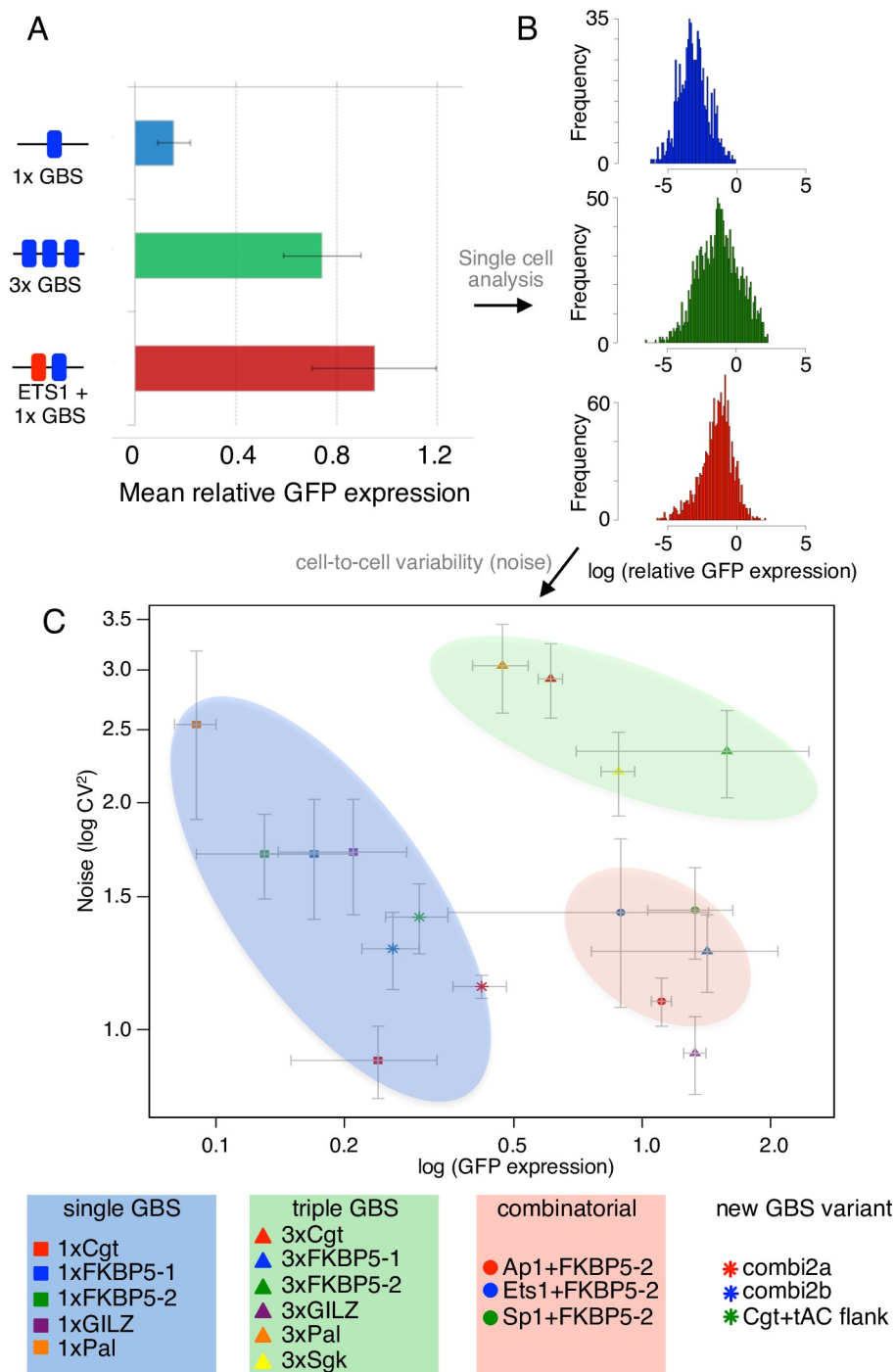


Fig 6. The effects of GBS sequence, number and presence of other TFBS on transcriptional output and noise. (a) Mean GFP expression relative to mCherry of the STARR-seq reporter for cell populations treated overnight with 1 μ M dexamethasone with binding site variant as indicated was determined by flow cytometry. (b) The single-cell distribution of GFP expression relative to co-transfected mCherry was determined for each binding site variant as indicated by flow cytometry. The mean and noise for each binding site variant are extracted from these distributions (see [Methods](#)). (c) Average and S.D. for mean GFP expression and for noise from three biological replicates. Area with mostly single GBS variants is highlighted with a blue background; Area with three GBSs with a green background and area with composite binding sites consisting of a single GBS and a binding site for another TF with a red background.

<https://doi.org/10.1371/journal.pgen.1007793.g006>

more robust GR-dependent activation when compared to single GBS variants (Fig 6A). The first group contained three instead of one GBS copy (triple GBS group) whereas the second group (composite group) contains a GBS flanked by a sequence motif for either AP1, ETS1 or SP1, three sequence motifs that can act synergistically with GR [34, 35]. As expected, the mean GFP expression was higher for each member of both the triple GBS and the composite group when compared to the single GBS group (Fig 6A and 6C). Interestingly, the increase in mean expression we found for the groups of triple GBS and composite enhancers was not accompanied by a decrease in expression noise (Fig 6C). The high noise to mean expression ratio was especially striking for several triple GBS variants (3xPal, 3xCgt, 3xSgk and 3x Fkbp5-2) but observed in general for each member of the groups of triple and composite enhancers when compared to the single GBS group. Furthermore, enhancer variants with similar mean expression levels (*e.g.* 3xSgk and Ets1+FKBP5-2) can have vastly different noise levels indicating that binding sites architecture can independently tune both mean expression and cell-to-cell variability in gene expression with noisier expression for enhancers with multiple GBSs.

Discussion

In this study, we developed a modified version of the STARR-seq method where we used designed synthetic oligonucleotides to assay how sequence variation within and around the GBS influence GBS activity. This facilitated the thorough and parallelized assessment of 1024 flank variants on GBS activity in a highly reproducible and quantitative fashion (Fig 1 and S1 Fig). Similarly, we assessed over 65,000 variants to study how variations in one of the half sites and the spacer influence GBS activity. Taken together, we find that variation in both the half site and in the region flanking the GBS influences GR activity. Quantitatively however, changes in the half site have more profound effects on activity than those in the flanking region (Fig 2D, Fig 4B). A key advantage of using designed sequences over the analysis of genomic regions is that variants can be compared in a context where everything is identical except for the sequence of the GR binding site. Notably, the sequence of the binding site is just one of several signals that are integrated at genomic response elements to modulate GR-dependent transcriptional responses. Therefore, our synSTARR-seq data for a single GBS is unlikely to yield accurate predictions for the activity of GR response elements in the genomic context. Other inputs that could improve predictions include chromatin environment and information regarding how the presence or absence of binding sites for other TFs influences GR-dependent activation. The synSTARR-seq approach can readily be adapted to study how combinations of signals are integrated. For example, principles of combinatorial regulation can be studied using designed sequences for which the GBS is flanked by binding sites for other TFs. Similarly, the assay can be used to investigate the cross-talk between GBS sequence, ligand chemistry, type of core promoter and GR splice isoforms.

Importantly, our findings for the synthetic STARR-seq assay are consistent with GR-dependent regulation of endogenous target genes. Specifically, the nucleotide directly flanking the GBS is preferentially a T for both enhancing flanks in our synSTARR-seq experiments and for the motif we previously found for genomic GR binding sites associated with genes that show the most robust response to GR activation [3]. Moreover, we uncovered a novel functional GR binding sequence variant with high activity, which we called combi2. Consistent with the high activity of the combi2 motif observed in the synSTARR assay, genes with nearby GR-bound peaks matching the combi2 motif were, on average, slightly more activated by GR than genes with peaks matching the consensus motif (based on RNA-seq data we generated for U2OS-GR cells treated for 4h with either 1μM dexamethasone or 0.1% ethanol as vehicle control; Fig 5C). Other sequence preferences we uncovered for flanks that enhance GBS activity include an A

followed by a C at positions 9 and 10 respectively (Fig 2A and 2C and S2A and S2C Fig). One possible explanation for the increased activity is that this sequence generates an additional GR half site or a binding site for another TF. However, the ChIP-exo profile for GBSs flanked by nAC looked essentially the same as the profile for the canonical GBS (S9 Fig), arguing against the binding of an additional factor. Alternatively, the flanking nAC could influence GR's DNA binding affinity. However, a comprehensive analysis of the effect of sequence variation within and in the regions flanking GR binding sites showed that the flanks essentially do not influence the binding affinity of GR [30]. Accordingly, we found similar K_d values for the AC flank when compared to variants with lower activity (Fig 2E) indicating that the change in activity is not driven by affinity. Together, the synSTARR-seq approach uncovered how sequence variation modulates GR activity, which confirmed previous findings based on a small number of sequences but also provided new insights into mechanisms that modulate GR-dependent regulation of endogenous target genes.

We were surprised to find that the consensus motifs for enhancing and blunting flanks displayed low information content indicating that a broad spectrum of distinct sequences can enhance or blunt the activity of the adjacent GBS (Fig 2 and S2 Fig). However, when looking at DNA shape we found specific shape characteristics for each group (Fig 3A, S2 Fig and S3 Fig). This indicates that distinct sequences can induce similar DNA shape characteristics with analogous effects on GBS activity. This finding was corroborated by our analysis of the spacer, which is not directly contacted by GR, yet influences GR activity. Also here we found distinct spacer shape characteristics for the most and least active GBS variants, without a clear sequence signature for each group (S7B Fig). Furthermore, we trained a model to distinguish between high and low activity GBSs based on either DNA sequence or on predicted minor groove width information. Assessment of the accuracy of the models using ROC curves showed that a single shape parameter, minor groove width, can be used to distinguish quite accurately between blunting and enhancing flanks (S10A Fig) and also between the top and bottom 20% active variants from our half site library (S10B Fig). Together, our findings which are based on a systematic analysis of many sequence variants are consistent with previous studies based on a small number of binding sites, showing that GR activity can be modulated by DNA shape [3, 4]. Notably, although the role of DNA shape in modulating the affinity of TFs for DNA has been well documented [36–38], we find that DNA shape modulates GR activity without apparent changes in DNA binding affinity (Fig 2E, [30]). This is consistent with a model where DNA shape acts as an allosteric ligand which induces structural changes in associated TFs which in turn changes the composition and regulatory activity of the complexes formed at the response element [5, 39–41]. Another, not mutually exclusive, explanation for flank-dependent modulation of transcriptional output is that flank variants serve as binding sites for other TFs that act additively or synergistically with GR. Further support for the importance of DNA shape comes from the analysis of the conservation of non-coding regions of the genome. This analysis uncovered greater conservation at the level of DNA shape than on the basis of nucleotide sequence indicating that DNA structure may be a better predictor of function than DNA sequence [42]. Accordingly, incorporation of DNA shape characteristics improves *in vivo* prediction of TF binding sites [43] and, based on our findings, could also improve the prediction of TF binding site activity.

We also explored if GFP protein expression levels of individual cells can be used to study how enhancer architecture influences cell-to-cell variability in gene expression. A similar approach was used to study how sequence variation of the promoter influences transcriptional noise in yeast [16]. Notably, the only difference between the reporters we assayed is their enhancer sequence, which is downstream of the ORF for the GFP protein. For sequences with related enhancer architectures (*e.g.* single GBS variants), we observed that transcriptional noise

scales with mean expression, such that higher expression levels are associated with lower noise (Fig 6C). This is consistent with a two-state promoter model where increases in mean expression are driven by an upsurge in transcription burst frequency [44]. Similarly, the estrogen receptor, a hormone receptor closely related to GR, modulates transcription by changing the frequency of transcriptional bursting [12]. When we compare distinct enhancer architectures, we find that expression mean and noise can be uncoupled. Specifically, the noise to mean expression ratio is higher for response elements harboring multiple TF binding sites when compared to the group of single GBS variants, indicating that the increase in expression might be accompanied by an increase in the number of transcripts produced during each burst. This finding is consistent with studies in yeast showing that increasing the number of binding sites for GCN4 results in increased expression with relatively high noise levels [16]. Notably, both multiple binding sites for GR and a combination of a GR binding site and a binding site for another TF result in an increased noise to mean expression ratio (Fig 6). Genomic GR response elements typically contain multiple GR binding sites and motifs for other TFs [25] arguing that endogenous GR-driven gene expression might be quite noisy with high levels of cell-to-cell variability in gene expression levels. However, if our findings for synthetic response elements reflect what happens at genomic response elements is unclear given that they differ in several ways. For example, in contrast to most endogenous response elements, the GBSs in our reporters are separated by 4 bp. Furthermore, each GBS sequence is identical in our reporters whereas most genomic GBS sequences have a unique sequence composition.

Our results are consistent with a model in which the architecture of the enhancer influences transcriptional burst size and frequency. However, more sophisticated single-cell studies of nascent transcripts are needed for a detailed understanding of the role of enhancer architecture given that our studies are based on the measurement of steady state fluctuations in protein levels. For example, in our experimental approach we cannot rule out that other mechanisms, including differences in RNA stability and translation rates, could contribute to the cell-to-cell variability in expression observed. Nonetheless, our findings argue that differences in enhancer architecture might contribute to gene-specific tuning of expression mean to noise ratios of GR target genes.

Conclusions

Taken together, we present synSTARR, an approach to measure how designed binding site variants influence transcriptional output and noise. The systematic analysis of sequence variants presented here resulted in the identification of a novel functional GR binding sequence and provides evidence for an important role of DNA shape in tuning GR activity without apparent changes in DNA binding affinity. Our simple approach using designed sequences can be applied to other TFs and can be used to systematically unravel how the interplay between sequence and other signaling inputs at response elements modulate transcriptional output.

Materials and methods

Experimental

Plasmids. STARR reporter constructs were generated by digesting the human STARR-seq vector [22] with SalI-HF and AgeI-HF and subsequent insertion of fragments of interest by in-Fusion HD cloning (TaKaRa). All inserts had the following sequence composition: 5'- **TAGA GCATGCACCGGACACTCTTTCCCTACACGACGCTCT**—*INSERT*—**AGATCGGAAG AGCACACGTCTGAACTCCAGTCACTCGACGAATTCGGCC**-3'. Sequence homologous

to the STARR reporter construct in bold; Sequence for p5 and p7 adaptors underlined. The exact sequence of the insert for each construct used in this study is listed in [S1 Table](#).

Cell lines, transient transfections and luciferase assays. U2OS cells stably transfected with rat GR α (U2OS-GR) [27] were grown in DMEM supplemented with 5% FBS. Transient transfections were done essentially as described [5] using either lipofectamine and plus reagents (Invitrogen) or using kit V for nucleofections (Lonza).

Synthetic STARR-seq. Library design and generation: To generate GBS variant libraries, oligos containing degenerate nucleotides (N) at defined positions were ordered from IDT as “DNA Ultramer oligonucleotide” (sequence listed below). The oligonucleotides were made double stranded using Phusion polymerase (NEB; 98°C for 35 sec, 72°C for 5 min) using the revPrimer (GGCCGAATTCGTCGAGTGAC). The resulting double stranded inserts (25ng) were recombined with 100ng linearized (Sall-HF and AgeI-HF) STARR-seq vector [22] by in-Fusion cloning in 5 parallel reactions. After pooling the reactions, the DNA was cleaned up using AMPure XP beads (Beckman Coulter), transformed into MegaX DH10B cells (Invitrogen) and plasmid DNA was isolated using a Plasmid Plus Maxi kit (Qiagen). **STARR-seq:** For STARR-seq experiments, 5 million U2OS-GR cells were transfected with 5 μ g library-DNA by nucleofection using kit V (Lonza). The next day, cells were treated for 4 h with 1 μ M dexamethasone or with 0.1% ethanol as vehicle control. Reverse transcription and amplification of cDNA for subsequence Illumina 50bp paired-end sequencing were done as described [22].

Cgt flank library DNA Ultramer oligonucleotide:

TAGAGCATGCACCGGACACTCTTTCCCTACACGACGCTCTTCCGATCTCAGCGC
AAGAACAAttTGTACGNNNNNCTAGATCGGAAGAGCACACGTCTGAACTCCAGTCA
CTCGACGAATTCGGCC

Sgk flank library DNA Ultramer oligonucleotide:

TAGAGCATGCACCGGACACTCTTTCCCTACACGACGCTCTTCCGATCTCAGCG
CAAGAACAAttTGTCCGNNNNNCTAGATCGGAAGAGCACACGTCTGAACTCCAGTC
ACTCGACGAATTCGGCC

GBS half site library DNA Ultramer oligonucleotide:

TAGAGCATGCACCGGACACTCTTTCCCTACACGACGCTCTTCCGATCTCAGCG
AAAGAACAtnNNNNNNNCGTCGCTAGATCGGAAGAGCACACGTCTGAACTCCAG
TCACTCGACGAATTCGGCC

RNA-seq U2OS-GR cells (Fig 5C). U2OS-GR cells were treated for 4h with either 1 μ M dexamethasone or 0.1% ethanol as vehicle control. RNA was isolated from 1.2 million cells using the RNeasy kit from Qiagen. Sequencing libraries were prepared using the TruSeq RNA library Prep Kit (Illumina). Prior to reverse transcription, poly adenylated RNA was isolated using oligo d(T) beads. Paired end 50bp reads from Illumina sequencing were mapped against the human hg19 reference genome using STAR [45] (options:—alignIntronMin 20—alignIntronMax 500000—chimSegmentMin 10—outFilterMismatchNoverLmax 0.05—outFilterMatchNmin 10—outFilterScoreMinOverLread 0—outFilterMatchNminOverLread 0—outFilterMismatchNmax 10—outFilterMultimapNmax 5). Differential gene expression between dex and etoh conditions from three biological replicates was calculated with DESeq2 [28], default parameters except betaPrior = FALSE.

Electrophoretic mobility shift assays. EMSAs were performed as described previously [3] using Cy-5 labeled oligos as listed in [S2 Table](#).

RNA isolation, reverse transcription and qPCR analysis. RNA was isolated from cells treated for either 4 h or overnight with 1 μ M dexamethasone or with 0.1% ethanol vehicle. Total RNA was reverse transcribed using gene-specific primers for *GFP* (CAAACCTCATCAA TGTATCTTATCATG) and *RPL19* (GAGGCCAGTATGTACAGACAAAGTGG) which was used for data normalization. qPCR and data analysis were done as described [5]. Primer pairs

for qPCR: hRPL19-fw: ATGTATCACAGCCTGTACCTG, hRPL19rev: TTCTTGGTCTCTTC CTCCTTG, GFP-fw: GGCCAGCTGTTGGGGTGTC, GFP-rev: TTGGGACAACTCCAGTG AAGA.

Noise-Measurements. For noise measurements, U2OS-GR cells were transfected using lipofectamine and plus (Invitrogen) essentially as described [5]. In short: The day before transfection, 40,000 U2OS-GR cells were seeded per well of a 24 well plate. The following day, cells were transfected with individual STARR reporter constructs (20ng/well) along with a SV-40 mCherry expression construct (20ng/well) and empty p6R plasmid (100 ng/ well). Transfected cells were treated overnight with either 1 μ M dexamethasone or with 0.1% ethanol vehicle control. Fluorescence intensity was measured using an Accuri C6 flow cytometer (BD Biosciences) and the yellow laser (552nm) and filter 610/20 for mCherry and the deepblue laser (473nm) and filter 510/20 to measure GFP. Gates were set for mCherry and GFP and only cells showing both mCherry and GFP fluorescence were included in the analysis. Relative expression of GFP (GFP/Cherry), from 800–1600 individual dexamethasone-treated cells, was used to calculate mean expression and the standard deviation of cell populations. Mean and standard deviation for noise (CV^2) and for relative GFP expression were derived from three biological replicates.

Computational analyses

Analysis of synSTARR-seq data. RNA-seq reads were filtered and only sequences exactly matching the insert sequence in length and nucleotide composition were included in the analysis. The number of occurrences for each sequence variants was counted for each experimental condition and differentially expressed sequences were identified using DESeq2 [28] using a p adjusted value <0.01 as cut-off. To fit the dispersion curve to the mean distribution, we used the local smoothed dispersion (DESeqwithfitType = "local"). Notably, each of the constructs of the flank libraries contains a functional GBS. Therefore, flanks that blunt activity will appear repressed after hormone treated because their fraction in the total pool of sequences decreases relative to flank variants with higher activities. For the flank libraries, we obtained information for each sequence variant (1024) in the library. For the half site library, we identified 61,582 out of the 65,536 possible variants present in this library. We found that including sequences with low read coverage resulted in many false positive differentially expressed GBS variants. To avoid this, we only included sequences with a mean read count above 100 across all experiments, leaving us with information for 33,689 sequence variants. The pearson correlation coefficient for replicates was calculated using the ggscatter function of the ggpubr library in R.

Boxplots comparing groups of sequence variants as specified in the figure legends show center lines for the median; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the interquartile range from the 25th and 75th percentiles.

Sequence logos to depict the consensus motif for groups of sequences were generated using WebLogo [46]. The probability logo (activity motif) was generated with kpLogo [29] using as input the sequence and fold change (dex/etoh) for each variant and the default settings for weighted sequences.

Motif weight. The motif weight for each variant was calculated using the RSAT *matrix-scan* program [47, 48]. Specifically, the motif weight was calculated using Transfac motif M00205 truncated to the core 15bp, and a custom background model created with RSAT *create background* program, trained on human open chromatin available at UCSC genome browser (<http://genome.ucsc.edu/cgi-bin/hgTrackUi?db=hg19&g=wgEncodeRegDnaseClustered>). Boxplots comparing groups of sequence variants show center lines for the median; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the interquartile range from the 25th and 75th percentiles.

Comparison of ChIP-seq peak height between combi2 and canonical GBS motif. GR ChIP-seq data sets for U2OS-GR cells were downloaded as processed peaks from EBI ArrayExpress (E-MTAB-2731). ChIP-seq peaks in a 40 kb window centered on the transcription start site of differentially expressed genes (RNA-seq data: E-MTAB-6738) were scanned using RSAT *matrix-scan* [47, 48] for the occurrence of either a GBS-match (Transfac matrix M00205, p value cut-off: 10^{-4}) or the combi2 matrix we generated (Fig 5D, p-value cut-off 10^{-4}). Next, peaks were grouped by motif match and median peak height was calculated for each group and the p-value comparing both groups was calculated using a Wilcoxon rank-sum test to produce Supplementary S8B Fig.

Comparison of gene regulation. To compare the level of activation between genes with nearby peaks with either a GBS match (Transfac matrix M00205, p value cut-off: 10^{-4}) or a combi2 match (motif Fig 5D, p-value cut-off 10^{-4}), we first scanned ChIP-seq peaks (U2OS-GR cells: E-MTAB-2731) in a 40 kb window centered on the transcription start site (using all annotated TSSs from Ensembl GRCH37) for motif matches using RSAT *matrix-scan* [47, 48]. Only peaks with an exclusive motif match were retained to generate a boxplot comparing the log2 fold change for genes of each group (RNA-seq data: E-MTAB-2731). Center lines show the median, box limits indicating the 25th and 75th percentiles and whiskers extending 1.5 times the interquartile range from the 25th and 75th percentiles. p-value comparing the log2 fold change for both groups was calculated using a Wilcoxon rank-sum test to produce Fig 5C.

DNA shape prediction. We used DNashapeR [31] to predict the minor groove width, roll, propeller twist or helix twist for sequence variants of interest. Boxplots for individual nucleotide position show center lines for the median; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the interquartile range from the 25th and 75th percentiles. The Wilcoxon rank-sum test was used to calculate the p-values comparing nucleotide position variants between groups. Individual sites were clustered using K-means clustering with $k = 4$ clusters $nstart = 20$ and 100 restarts with the function 'kmeans' from the R 'stats' package.

Classification of GBS activity. To assess classifier performance we generate ROC curves using 10-fold cross-validation. Four different models were tested to classify GBS activity into blunting or enhancing. A mononucleotide model consisting of sequence motifs estimated from relative nucleotide frequencies within the two classes. Class affiliation is predicted with a likelihood ratio test. We also tested a similar model based on dinucleotides. In addition, we tested two random forest (RF) classifiers with 100 trees, based on sequence and shape information. We used the R package "randomForest" for constructing the classifiers [49]. Since RF classifiers are not designed for categorical data, we coded nucleotide sequences using 00 for 'A', 01 for 'C', 10 for 'G', and 11 for 'T'.

ChIP-exo footprint profiles. ChIP-exo footprint profiles were generated using the ExoProfiler package [32] and published ChIP-exo (EBI ArrayExpress E-MTAB-2955) and ChIP-seq (E-MTAB-2956) data for IMR90 cells as input. Peaks were scanned using either the JASPAR MA0113.2 motif [50], the PWM for the combi1 motif [32], the combi2 motif (Fig 5D) or for the AC flank variant, the motif depicted in S9A Fig. Hits were included if the p-value was $<10^{-4}$. Overlay plots for distinct motifs were generated by aligning the profiles on the GBS and normalizing the signal for each motif variant to 1.

Structural alignment of GR:ETS1 complex. Structural alignment of the GR:ETS1 complex on a combi2 sequence was done as described previously [32] except that both GR dimer halves are retained in the resulting model. In short: A structural model of the DNA hybrid sequence (AGAACATTCCGGCACT) was generated using 3D-Dart [51] using the ETS1 structure (PDB entry 1K79) and the GR structure (PDB entry 3G6U). GR and the ETS2 binding motifs were aligned using the CE-align algorithm [52] to the 3D-DART DNA model of the hybrid sequence.

Data access

Data were deposited in ArrayExpress under the accession numbers: E-MTAB-6738 (RNA-seq U2OS-GR) and E-MTAB-6737 (synSTARR-seq U2OS-GR). In addition, we used the previously deposited datasets: E-MTAB-2731 (ChIP-seq U2OS cells), E-MTAB-2955 and E-MTAB-2956 (ChIP-seq and ChIP-exo data IMR90).

Supporting information

S1 Fig. Analysis of individual enhancer variants by flow cytometry and synSTARR-seq reproducibility. (a) Analysis of individual enhancer variants as indicated by flow cytometry showing the side scatter (SSC-A) versus GFP signal for individual mCherry-positive cells. Left: no STARR-seq construct. Right-Top: ethanol, vehicle, treated cells; Right-Bottom: Cells treated overnight with 1 μ M dexamethasone. Numbers in red indicate the percentage of GFP+ (top right side) and GFP- (top left side) cells respectively. Red vertical line demarcates the threshold for being called GFP+. (b) RNA-seq correlation plots for biological replicates of vehicle-treated cells transfected with the GBS-flank library (Cgt flank library). (c) Same as (b) except for biological replicates of dexamethasone-treated cells (4h 1 μ M). (d) RNA-seq correlation plots for biological replicates of vehicle-treated cells transfected with the GBS-flank library (Sgk flank library). (e) Same as (d) except for biological replicates of dexamethasone-treated cells (4h 1 μ M). (TIF)

S2 Fig. Analysis of the Sgk flank library. (a) Color chart summarizing the sequence at each variable position for flank variants ranked by their fold change in response to hormone treatment. (b) Consensus motif for (left) significantly enhancing and (right) blunting flank variants (c) *kp*Logo probability logo (activity logo) for flank variants depicting the p-values from Mann-Whitney U tests of whether GBS variants with a specific nucleotide at a given position are more (displayed above number indicating nucleotide position) or less (displayed below number indicating nucleotide position) active than other GBS variants. Positions with significant nucleotides ($p < 0.001$) are highlighted (red coordinates). (d) Predicted minor groove width (MGW) for significant enhancing and blunting flank variants of the Sgk GBS library ranked by their fold change in response to hormone treatment. (e) K-means clustering based on MGW for significantly enhancing and blunting flank variants. Right side: activating and blunting variants are highlighted in grey and black respectively. (f) Consensus sequence motif for clusters as indicated. (g) Log2 fold change upon dexamethasone treatment for each cluster as indicated. The synSTARR-seq activity for individual sequences is shown as black dots, the median for each group as a horizontal red line. (TIFF)

S3 Fig. DNA shape comparison (propeller twist, helix twist and roll) between blunting and enhancing flanks. (a) Predicted Propeller twist for significant enhancing and blunting flank variants of the Cgt GBS library ranked by their fold change in response to hormone treatment. (b) Propeller twist for selected individual bases for significantly blunting ($n = 189$) and significantly enhancing ($n = 125$) flanks for the Cgt library. p-values were calculated using the Wilcoxon rank-sum test ($*p < 0.01$). (c) Same as (a) except that helix twist is shown. (d) same as (b) except that helix twist is shown. (e) Same as (a) except that roll is shown. (f) same as (b) except that roll is shown. (TIFF)

S4 Fig. MGW comparison between blunting and enhancing flanks. (a) Minor groove width (MGW) for selected individual bases for significantly blunting ($n = 189$) and significantly

enhancing ($n = 125$) flanks for the Cgt library. p-values were calculated using the Wilcoxon rank-sum test. (b) Same as for (a) except for significantly blunting ($n = 162$) and significantly enhancing ($n = 101$) flanks of the Sgk flank library.

(TIFF)

S5 Fig. synSTARR-seq reproducibility for the half site library. (a) Correlation plot between input library (library) and the plasmid library isolated from transfected U2OS-GR cells (input). (b) RNA-seq correlation plots for biological replicates of vehicle-treated cells. (c) Same as for (b) except for biological replicates of dexamethasone-treated cells (4h 1 μ M).

(TIFF)

S6 Fig. Analysis of the GBS half site library. (a) Log2 fold change upon dexamethasone treatment for active GBS variants with either an A, C, G or T at position 3. Data for individual sequences that match consensus second half site at key positions 4 and 6 (exact match to AGAACATnnnXTnCn, with X either A,C,G or T) are shown as blue dots. Horizontal red lines show the average for each group. p-value was calculated using a Student's t-test. (b) Left: Minor groove width (MGW) prediction for GBS variants ranked by activity. Right: Consensus motif for top 20% most active and bottom 20% least active GBS variants. (c) MGW for select individual bases comparing the top 20% most active and bottom 20% least active activated GBS variants. p-values were calculated using the Wilcoxon rank-sum test.

(TIFF)

S7 Fig. Effects of spacer sequence on GBS activity. (a) Motif logo representing the sequence that was used to scan for GBS-matches in the half site library. Black box highlights the two positions in the spacer whose effects on GBS activity was assayed. (b) Boxplot of the log2 fold change upon treatment for 4 h with 1 μ M dexamethasone for GBS matches with spacer variant as indicated. Center lines show the median. Spacer variants with significantly different activity levels (Benjamini-Hochberg corrected p-values calculated using a Student's t-test < 0.01 are highlighted).

(TIFF)

S8 Fig. Characterization of the combi2 motif. (a) Transcriptional activity of STARR-seq reporters containing candidate GBS variants as indicated. Relative RNA levels \pm S.E.M. are shown for cells treated with ethanol vehicle and for cells treated overnight with 1 μ M dexamethasone ($n = 3$). (b) Boxplot showing the peak-height for GR target genes with either a canonical GBS motif match (nnTGT) or a combi2 motif match (tcCGT). Center lines show the median, p value was calculated using a Wilcoxon rank-sum test.

(TIFF)

S9 Fig. Analysis of the nACnn flank. (a) Motif logo representing the positional weight matrix of highly active flank variants that was used to scan for motif-matches to generate the ChIP-exo footprint profile. (b) Alignment of the ChIP-exo footprint profiles for highly active flank variant matches (p value < 0.0001 ; solid lines: blue: positive strand, red: negative strand) and for the conventional GBS motif (M00205; p value < 0.0001 ; shaded areas; blue: positive strand, red: negative strand).

(TIFF)

S10 Fig. Prediction of GBS activity based on DNA sequence or DNA shape. (a) ROC curves analyzing the ability of the models to distinguish between blunting and enhancing flank variants for (left) the Cgt flank library; (right) the Sgk flank library. Mononucleotide: Classifier based on mononucleotide frequencies within the two classes. Dinucleotide: Classifier constructed using dinucleotide frequencies. Sequence Random Forest (RF): Random Forest

classifier trained and tested on coded nucleotide sequences. Shape Random Forest (RF): Random forest classifier based on predicted MGW. (b) Same as for (a) except that model and ROC curves were trained and assessed for their ability to discriminate between the top and bottom 20% significantly active GBS variants from the half site library.
(TIFF)

S1 Table. Inserts plasmids. GBS (or scrambled GBS) underlined.
(PDF)

S2 Table. Oligos for EMSAs.
(PDF)

S1 Data. Numerical data underlying display items.
(XLSX)

Acknowledgments

We would like to thank Marcel Jurk for performing the structural alignment presented.

Author Contributions

Conceptualization: Stefanie Schöne, Sebastiaan H. Meijnsing.

Formal analysis: Stefanie Schöne, Melissa Bothe, Sebastiaan H. Meijnsing.

Funding acquisition: Martin Vingron, Sebastiaan H. Meijnsing.

Investigation: Melissa Bothe, Edda Einfeldt, Marina Borschiwer, Philipp Benner, Morgane Thomas-Chollier, Sebastiaan H. Meijnsing.

Methodology: Stefanie Schöne, Sebastiaan H. Meijnsing.

Resources: Martin Vingron.

Supervision: Morgane Thomas-Chollier, Sebastiaan H. Meijnsing.

Writing – original draft: Morgane Thomas-Chollier, Sebastiaan H. Meijnsing.

Writing – review & editing: Stefanie Schöne, Melissa Bothe, Marina Borschiwer, Martin Vingron.

References

1. Grossman SR, Zhang X, Wang L, Engreitz J, Melnikov A, Rogov P, et al. Systematic dissection of genomic features determining transcription factor binding and enhancer function. *Proceedings of the National Academy of Sciences of the United States of America*. 2017 Feb 14; 114(7):E1291–E300. <https://doi.org/10.1073/pnas.1621150114> PMID: 28137873. Pubmed Central PMCID: 5321001.
2. Schmid W, Strahle U, Schutz G, Schmitt J, Stunnenberg H. Glucocorticoid receptor binds cooperatively to adjacent recognition sites. *The EMBO journal*. 1989 Aug; 8(8):2257–63. PMID: 2792086. Pubmed Central PMCID: 401156.
3. Schöne S, Jurk M, Helabad MB, Dror I, Lebars I, Kieffer B, et al. Sequences flanking the core-binding site modulate glucocorticoid receptor structure and activity. *Nature communications*. 2016 Sep 1; 7:12621. <https://doi.org/10.1038/ncomms12621> PMID: 27581526. Pubmed Central PMCID: 5025757.
4. Watson LC, Kuchenbecker KM, Schiller BJ, Gross JD, Pufall MA, Yamamoto KR. The glucocorticoid receptor dimer interface allosterically transmits sequence-specific DNA signals. *Nature structural & molecular biology*. 2013 Jul; 20(7):876–83. <https://doi.org/10.1038/nsmb.2595> PMID: 23728292. Pubmed Central PMCID: 3702670.
5. Meijnsing SH, Pufall MA, So AY, Bates DL, Chen L, Yamamoto KR. DNA binding site sequence directs glucocorticoid receptor structure and activity. *Science*. 2009 Apr 17; 324(5925):407–10. <https://doi.org/10.1126/science.1164265> PMID: 19372434. Pubmed Central PMCID: 2777810.

6. Maheshri N, O'Shea EK. Living with noisy genes: how cells function reliably with inherent variability in gene expression. *Annual review of biophysics and biomolecular structure*. 2007; 36:413–34. <https://doi.org/10.1146/annurev.biophys.36.040306.132705> PMID: 17477840.
7. Blake WJ, M KA, Cantor CR, Collins JJ. Noise in eukaryotic gene expression. *Nature*. 2003 Apr 10; 422 (6932):633–7. <https://doi.org/10.1038/nature01546> PMID: 12687005.
8. Kaern M, Elston TC, Blake WJ, Collins JJ. Stochasticity in gene expression: from theories to phenotypes. *Nature reviews Genetics*. 2005 Jun; 6(6):451–64. <https://doi.org/10.1038/nrg1615> PMID: 15883588.
9. Raj A, van Oudenaarden A. Single-molecule approaches to stochastic gene expression. *Annual review of biophysics*. 2009; 38:255–70. <https://doi.org/10.1146/annurev.biophys.37.032807.125928> PMID: 19416069. Pubmed Central PMCID: 3126657.
10. Ross IL, Browne CM, Hume DA. Transcription of individual genes in eukaryotic cells occurs randomly and infrequently. *Immunology and cell biology*. 1994 Apr; 72(2):177–85. <https://doi.org/10.1038/icb.1994.26> PMID: 8200693.
11. Suter DM, Molina N, Gattfield D, Schneider K, Schibler U, Naef F. Mammalian genes are transcribed with widely different bursting kinetics. *Science*. 2011 Apr 22; 332(6028):472–4. <https://doi.org/10.1126/science.1198817> PMID: 21415320.
12. Fritsch C, Baumgartner S, Kuban M, Steinshorn D, Reid G, Legewie S. Estrogen-dependent control and cell-to-cell variability of transcriptional bursting. *Molecular systems biology*. 2018 Feb 23; 14(2): e7678. <https://doi.org/10.15252/msb.20177678> PMID: 29476006. Pubmed Central PMCID: 5825209.
13. Cai L, Friedman N, Xie XS. Stochastic protein expression in individual cells at the single molecule level. *Nature*. 2006 Mar 16; 440(7082):358–62. <https://doi.org/10.1038/nature04599> PMID: 16541077.
14. Raser JM, O'Shea EK. Control of stochasticity in eukaryotic gene expression. *Science*. 2004 Jun 18; 304(5678):1811–4. <https://doi.org/10.1126/science.1098641> PMID: 15166317. Pubmed Central PMCID: 1410811.
15. Blake WJ, Balazsi G, Kohanski MA, Isaacs FJ, Murphy KF, Kuang Y, et al. Phenotypic consequences of promoter-mediated transcriptional noise. *Molecular cell*. 2006 Dec 28; 24(6):853–65. <https://doi.org/10.1016/j.molcel.2006.11.003> PMID: 17189188.
16. Sharon E, van Dijk D, Kalma Y, Keren L, Manor O, Yakhini Z, et al. Probing the effect of promoters on noise in gene expression using thousands of designed sequences. *Genome research*. 2014 Oct; 24 (10):1698–706. <https://doi.org/10.1101/gr.168773.113> PMID: 25030889. Pubmed Central PMCID: 4199362.
17. Hornung G, Bar-Ziv R, Rosin D, Tokuriki N, Tawfik DS, Oren M, et al. Noise-mean relationship in mutated promoters. *Genome research*. 2012 Dec; 22(12):2409–17. <https://doi.org/10.1101/gr.139378.112> PMID: 22820945. Pubmed Central PMCID: 3514670.
18. Dey SS, Foley JE, Limsirichai P, Schaffer DV, Arkin AP. Orthogonal control of expression mean and variance by epigenetic features at different genomic loci. *Molecular systems biology*. 2015 May 5; 11 (5):806. <https://doi.org/10.15252/msb.20145704> PMID: 25943345. Pubmed Central PMCID: 4461400.
19. Dadiani M, van Dijk D, Segal B, Field Y, Ben-Artzi G, Raveh-Sadka T, et al. Two DNA-encoded strategies for increasing expression with opposing effects on promoter dynamics and transcriptional noise. *Genome research*. 2013 Jun; 23(6):966–76. <https://doi.org/10.1101/gr.149096.112> PMID: 23403035. Pubmed Central PMCID: 3668364.
20. Raser JM, O'Shea EK. Noise in gene expression: origins, consequences, and control. *Science*. 2005 Sep 23; 309(5743):2010–3. <https://doi.org/10.1126/science.1105891> PMID: 16179466. Pubmed Central PMCID: 1360161.
21. Inoue F, Ahituv N. Decoding enhancers using massively parallel reporter assays. *Genomics*. 2015 Sep; 106(3):159–64. <https://doi.org/10.1016/j.ygeno.2015.06.005> PMID: 26072433. Pubmed Central PMCID: 4540663.
22. Arnold CD, Gerlach D, Stelzer C, Boryn LM, Rath M, Stark A. Genome-wide quantitative enhancer activity maps identified by STARR-seq. *Science*. 2013 Mar 1; 339(6123):1074–7. <https://doi.org/10.1126/science.1232542> PMID: 23328393.
23. Liu Y, Yu S, Dhiman VK, Brunetti T, Eckart H, White KP. Functional assessment of human enhancer activities using whole-genome STARR-sequencing. *Genome biology*. 2017 Nov 20; 18(1):219. <https://doi.org/10.1186/s13059-017-1345-5> PMID: 29151363. Pubmed Central PMCID: 5694901.
24. Vanhille L, Griffon A, Maqbool MA, Zacarias-Cabeza J, Dao LT, Fernandez N, et al. High-throughput and quantitative assessment of enhancer activity in mammals by CapStarr-seq. *Nature communications*. 2015 Apr 15; 6:6905. <https://doi.org/10.1038/ncomms7905> PMID: 25872643.
25. Vockley CM, D'Ippolito AM, McDowell IC, Majoros WH, Safi A, Song L, et al. Direct GR Binding Sites Potentiate Clusters of TF Binding across the Human Genome. *Cell*. 2016 Aug 25; 166(5):1269–81 e19. <https://doi.org/10.1016/j.cell.2016.07.049> PMID: 27565349. Pubmed Central PMCID: 5046229.

26. Shlyueva D, Stelzer C, Gerlach D, Yanez-Cuna JO, Rath M, Boryn LM, et al. Hormone-responsive enhancer-activity maps reveal predictive motifs, indirect repression, and targeting of closed chromatin. *Molecular cell*. 2014 Apr 10; 54(1):180–92. <https://doi.org/10.1016/j.molcel.2014.02.026> PMID: 24685159.
27. Rogatsky I, Trowbridge JM, Garabedian MJ. Glucocorticoid receptor-mediated cell cycle arrest is achieved through distinct cell-specific transcriptional regulatory mechanisms. *Molecular and cellular biology*. 1997 Jun; 17(6):3181–93. PMID: 9154817. Pubmed Central PMCID: 232171.
28. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome biology*. 2014; 15(12):550. <https://doi.org/10.1186/s13059-014-0550-8> PMID: 25516281. Pubmed Central PMCID: 4302049.
29. Wu X, Bartel DP. kpLogo: positional k-mer analysis reveals hidden specificity in biological sequences. *Nucleic acids research*. 2017 Jul 3; 45(W1):W534–W8. <https://doi.org/10.1093/nar/gkx323> PMID: 28460012. Pubmed Central PMCID: 5570168.
30. Zhang L, Martini GD, Rube HT, Kribelbauer JF, Rastogi C, FitzPatrick VD, et al. SelexGLM differentiates androgen and glucocorticoid receptor DNA-binding preference over an extended binding site. *Genome research*. 2018 Jan; 28(1):111–21. <https://doi.org/10.1101/gr.222844.117> PMID: 29196557. Pubmed Central PMCID: 5749176.
31. Chiu TP, Comoglio F, Zhou T, Yang L, Paro R, Rohs R. DNASHapeR: an R/Bioconductor package for DNA shape prediction and feature encoding. *Bioinformatics*. 2016 Apr 15; 32(8):1211–3. <https://doi.org/10.1093/bioinformatics/btv735> PMID: 26668005. Pubmed Central PMCID: 4824130.
32. Starick SR, Ibn-Salem J, Jurk M, Hernandez C, Love MI, Chung HR, et al. ChIP-exo signal associated with DNA-binding motifs provides insight into the genomic binding of the glucocorticoid receptor and cooperating transcription factors. *Genome research*. 2015 Jun; 25(6):825–35. <https://doi.org/10.1101/gr.185157.114> PMID: 25720775. Pubmed Central PMCID: 4448679.
33. Rhee HS, Pugh BF. Comprehensive genome-wide protein-DNA interactions detected at single-nucleotide resolution. *Cell*. 2011 Dec 9; 147(6):1408–19. <https://doi.org/10.1016/j.cell.2011.11.013> PMID: 22153082. Pubmed Central PMCID: 3243364.
34. Strahle U, Schmid W, Schutz G. Synergistic action of the glucocorticoid receptor with transcription factors. *The EMBO journal*. 1988 Nov; 7(11):3389–95. PMID: 2463158. Pubmed Central PMCID: 454837.
35. Pearce D, Yamamoto KR. Mineralocorticoid and glucocorticoid receptor activities distinguished by non-receptor factors at a composite response element. *Science*. 1993 Feb 19; 259(5098):1161–5. PMID: 8382376.
36. Yang L, Orenstein Y, Jolma A, Yin Y, Taipale J, Shamir R, et al. Transcription factor family-specific DNA shape readout revealed by quantitative specificity models. *Molecular systems biology*. 2017 Feb 6; 13(2):910. <https://doi.org/10.15252/msb.20167238> PMID: 28167566. Pubmed Central PMCID: 5327724.
37. Rohs R, Jin X, West SM, Joshi R, Honig B, Mann RS. Origins of specificity in protein-DNA recognition. *Annual review of biochemistry*. 2010; 79:233–69. <https://doi.org/10.1146/annurev-biochem-060408-091030> PMID: 20334529. Pubmed Central PMCID: 3285485.
38. Abe N, Dror I, Yang L, Slattery M, Zhou T, Bussemaker HJ, et al. Deconvolving the recognition of DNA shape from sequence. *Cell*. 2015 Apr 9; 161(2):307–18. <https://doi.org/10.1016/j.cell.2015.02.008> PMID: 25843630. Pubmed Central PMCID: 4422406.
39. Zheng J, Chang MR, Stites RE, Wang Y, Bruning JB, Pascal BD, et al. HDX reveals the conformational dynamics of DNA sequence specific VDR co-activator interactions. *Nature communications*. 2017 Oct 13; 8(1):923. <https://doi.org/10.1038/s41467-017-00978-7> PMID: 29030554. Pubmed Central PMCID: 5640644.
40. Zhang J, Chalmers MJ, Stayrook KR, Burris LL, Wang Y, Busby SA, et al. DNA binding alters coactivator interaction surfaces of the intact VDR-RXR complex. *Nature structural & molecular biology*. 2011 May; 18(5):556–63. <https://doi.org/10.1038/nsmb.2046> PMID: 21478866. Pubmed Central PMCID: 3087838.
41. Hall JM, McDonnell DP, Korach KS. Allosteric regulation of estrogen receptor structure, function, and coactivator recruitment by different estrogen response elements. *Molecular endocrinology*. 2002 Mar; 16(3):469–86. <https://doi.org/10.1210/mend.16.3.0814> PMID: 11875105.
42. Parker SC, Hansen L, Abaan HO, Tullius TD, Margulies EH. Local DNA topography correlates with functional noncoding regions of the human genome. *Science*. 2009 Apr 17; 324(5925):389–92. <https://doi.org/10.1126/science.1169050> PMID: 19286520. Pubmed Central PMCID: 2749491.
43. Mathelier A, Xin B, Chiu TP, Yang L, Rohs R, Wasserman WW. DNA Shape Features Improve Transcription Factor Binding Site Predictions In Vivo. *Cell systems*. 2016 Sep 28; 3(3):278–86 e4. <https://doi.org/10.1016/j.cels.2016.07.001> PMID: 27546793. Pubmed Central PMCID: 5042832.
44. Singh A, Razooky B, Cox CD, Simpson ML, Weinberger LS. Transcriptional bursting from the HIV-1 promoter is a significant source of stochastic noise in HIV-1 gene expression. *Biophysical journal*. 2010

- Apr 21; 98(8):L32–4. <https://doi.org/10.1016/j.bpj.2010.03.001> PMID: 20409455. Pubmed Central PMCID: 2856162.
45. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013 Jan 1; 29(1):15–21. <https://doi.org/10.1093/bioinformatics/bts635> PMID: 23104886. Pubmed Central PMCID: 3530905.
46. Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: a sequence logo generator. *Genome research*. 2004 Jun; 14(6):1188–90. <https://doi.org/10.1101/gr.849004> PMID: 15173120. Pubmed Central PMCID: 419797.
47. Thomas-Chollier M, Defrance M, Medina-Rivera A, Sand O, Herrmann C, Thieffry D, et al. RSAT 2011: regulatory sequence analysis tools. *Nucleic acids research*. 2011 Jul; 39(Web Server issue):W86–91. <https://doi.org/10.1093/nar/gkr377> PMID: 21715389. Pubmed Central PMCID: 3125777.
48. Turatsinze JV, Thomas-Chollier M, Defrance M, van Helden J. Using RSAT to scan genome sequences for transcription factor binding sites and cis-regulatory modules. *Nature protocols*. 2008; 3(10):1578–88. <https://doi.org/10.1038/nprot.2008.97> PMID: 18802439.
49. Wiener ALaM. Classification and Regression by randomForest. *R News*. 2002; 2(3):18–22.
50. Khan A, Fornes O, Stigliani A, Gheorghe M, Castro-Mondragon JA, van der Lee R, et al. JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic acids research*. 2018 Jan 4; 46(D1):D260–D6. <https://doi.org/10.1093/nar/gkx1126> PMID: 29140473. Pubmed Central PMCID: 5753243.
51. van Dijk M, Bonvin AM. 3D-DART: a DNA structure modelling server. *Nucleic acids research*. 2009 Jul; 37(Web Server issue):W235–9. <https://doi.org/10.1093/nar/gkp287> PMID: 19417072. Pubmed Central PMCID: 2703913.
52. Jia Y, Dewey TG, Shindyalov IN, Bourne PE. A new scoring function and associated statistical significance for structure alignment by CE. *Journal of computational biology: a journal of computational molecular cell biology*. 2004; 11(5):787–99. <https://doi.org/10.1089/cmb.2004.11.787> PMID: 15700402.