1     **Competition between binding sites determines gene expression at**

2     **low transcription factor concentrations**

3

4     **Authors:** David van Dijk[1,2,3,†], Eilon Sharon[5,2,3,†], Maya Lotan-Pompan[2,3], Adina

5     Weinberger[2,3], Eran Segal[2,3,#,*], and Lucas B. Carey[4,#,*]

6

7     **Affiliations:**

8     1 Department of Biological Sciences, Department of Systems Biology, Columbia

9     University, New York, NY, USA

10     2 Department of Computer Science and Applied Mathematics, Weizmann Institute of

11     Science, Rehovot, Israel

12     3 Department of Molecular Cell Biology, Weizmann Institute of Science, Rehovot, Israel

13     4 Department of Experimental and Health Sciences, Universitat Pompeu Fabra, Dr.

14     Aiguader 88, 08003 Barcelona, Spain.

15     5 Current address: Department of Genetics, Stanford University, Stanford, California,

16     United States of America

17     † These authors contributed equally to this work.

18     # These authors contributed equally to this work.

19     * Correspondence to: E. Segal (eran.segal@weizmann.ac.il) or L. Carey (lucas.carey@upf.edu)

20

21

22

# **Abstract**

The response of gene expression to intra- and extra-cellular cues is largely mediated

through changes in the activity of transcription factors (TFs), whose sequence specificities

are largely known. However, the rules by which promoters decode the amount of active

TF into gene expression are not well understood. Here, we measure the activity of 6500

designed promoters at six different levels of TF activity in budding yeast. We observe that

maximum promoter activity is determined by TF activity and not by the number of sites.

Surprisingly, the addition of an activator-binding site often reduces expression. A

thermodynamic model that incorporates competition between neighboring binding sites

for a local pool of TF molecules explains this behavior and accurately predicts both

absolute expression and the amount by which addition of a site increases or reduces

expression. Taken together, our findings support a model in which neighboring binding

sites interact competitively when TF is limiting but otherwise act additively

1

## **Significance Statement**

3   In response to intracellular and extracellular signals organisms alter the concentration

4   and activity of transcription factors (TFs), proteins that regulate gene expression.

5   However, the molecular mechanisms that determine the response of a target promoter to

6   changes in the number of active TF molecules are not well understood. By combining

7   mathematical modeling with measurements of TF dose-response curves for thousands

8   of designed promoters, we show that competition for active TF molecules is a major factor

9   in determining gene expression. At low TF concentrations additional activator-binding

10   sites within a promoter can actually reduce expression. Thermodynamic modeling

11   suggests that steric hindrance between neighboring binding sites cannot explain this

12   behavior, but that competition for limiting TF molecules can.

13

## 1 Introduction

2       Cells respond to internal and external changes by controlling their gene expression

3 programs. A major mechanism by which this is achieved is by modulating the activity of

4 transcription factors (TFs) that bind to specific sites in gene promoters where they activate

5 or repress transcription(Struhl, 1995). For example, in the budding yeast *S. cerevisiae,*

6 almost half of the genome changes in expression level in response to amino acid

7 starvation. A single transcription factor, Gcn4, is responsible for the activation of over 500

8 of these genes(Natarajan et al., 2001). While transcriptome and Chromain-IP (ChIP)

9 studies are useful for understanding the wiring of these large regulatory networks, they

10 are not informative about how the quantitative relationship between TF and target gene

11 expression are encoded in the DNA. It is still not well understood how promoter

12 architecture determines how each target of a TF will respond to changes in the amount

13 of active TF ([TF]). Furthermore, many targets of the same transcription factor are

14 expressed at different levels in the absence of that TF, and the fold-induction of the target

15 is largely independent of its expression at low or high TF (Carey et al., 2013; Rajkumar

16 et al., 2013). However, the molecular mechanisms that enable this decoupling are largely

17 unknown.

18

19       In order to understand how promoters encode the function that maps the amount

20 of active TF to a gene transcription level output, we measure the dose response curves

21 for 6500 synthetically designed promoters. We have used a synthetic approach(Sharon

22 et al., 2012) in which pairs of promoters differ by a single regulatory element, in contrast

4

1    to native promoters that have many differences between them, preventing systematic

2    investigation of the effect of individual DNA sequence elements on expression response.

3

4    We observe a wide range of dose response curves that differ in both TF

5    independent expression, and TF dependent (dynamic range) expression. The first, TF

6    independence, can be attributed to differences in the predicted nucleosome occupancy

7    of the promoters. In contrast, TF dependent (dynamic range) differences are mostly

8    determined by the number and affinity of binding sites. Surprisingly, we find that

9    expression level saturates with number of binding sites, at a level determined by the

10   amount of active TF ([TF]) and not by the number of binding sites. Moreover, in many

11   cases adding an activator can actually reduce expression, especially at low [TF]. In order

12   to quantitatively understand the results, we test several hypotheses using a

13   thermodynamic model that incorporates cooperative and competitive interactions

14   between TF binding sites. Our results suggest that expression synergism (i.e. the additive

15   or reductive effect of adding an additional homotypic binding site) is due to 'sharing' a

16   local pool of TF between nearby sites.

17

18   **Results**

19

20   **Promoter DNA sequence can encode a wide range of transcriptional responses to**

21   **changes in the amount of active TF**

22   To systematically measure how transcriptional responses are encoded in promoter

23   DNA sequence we measured the activity of 6500 designed promoters using a

5

1    fluorescence reporter (Sharon et al., 2012) in six growth media that each differ in their

2    concentration of amino acids ([AA]) (see **Methods** for details). The majority of these

3    promoters contain binding sites for Gcn4, Leu3, Met31 or Bas1; TFs involved in amino

4    acid biosynthesis. At high [AA] the TFs Gcn4, Bas1, Leu3 and Met31 are mostly

5    inactive(Struhl, 1992); at low [AA] the concentration of the active form of these TFs

6    increases (their expression and/or ability to activate transcription increases), and their

7    targets increase in expression(Gasch et al., 2000). For these four TFs, the number of

8    active TF molecules ([TF]) increases gradually in response to decreasing [AA] (Ljungdahl

9    and Daignan-Fornier, 2012) **(Supplementary Fig. 1)**. The combinatorial fashion in which

10    TF binding site type, number, affinity, position and accessibility vary in the designed

11    promoter set enables us to systematically investigate the mapping between promoter

12    DNA sequence, [TF] and the induced expression (**Fig. 1A,** see Methods for details).

13      The measurements were carried out using our previously described method that

14    involves FACS sorting and deep sequencing of a barcoded pooled promoter

15    library(Sharon et al., 2012; 2014). Briefly, uniquely barcoded promoters that drive a YFP

16    reporter are FACS sorted into 12 bins of expression that subsequently receive an

17    expression-bin barcode. Deep sequencing thus results in reads that contain both a

18    sequence and an expression barcode. A computational analysis of these reads gives, for

19    each promoter and growth condition, an expression distribution, from which the mean is

20    extracted, resulting in 6500 highly reproducible (**Supplementary Fig. 2**) dose-response

21    curves (**Fig. 1B**). Our promoters encode a wide range of responses with a general trend

22    in which more TF binding sites give a greater dynamic range between low and high [AA]

23    (**Fig. 1B**). We observe that some promoter sequence changes (e.g.: addition of a polyT,

6

1   **Fig. 1C**) affect expression independent of [AA], whereas others (e.g.: addition of Gcn4

2   binding sites, **Fig. 1C**) affect expression in a manner that depends on [AA]. We refer to

3   the former as (active) [TF] independent expression change, and the latter as [TF]

4   dependent expression change.

5

6   **Decoupled control of TF dependent and TF independent expression**

7         In order to distinguish between promoter sequence features that affect expression

8   in a [TF] dependent manner and those that affect expression in a [TF] independent

9   manner, we compare expression at high and low [AA] (see **Methods**) for promoters

10  grouped by DNA sequence features. We find that the number of Gcn4 binding sites affects

11  expression in a TF dependent manner: adding binding sites results, on average, in little

12  increase in expression at high [AA] but a large increase at low [AA] (**Fig. 2A,D**), and thus

13  an increase in the promoter's dynamic range (**Fig. 2E**). The same results are observed

14  for increasing the affinity of the Gcn4 binding site: increasing the affinity results in slightly

15  higher expression at high [AA], much higher expression at low [AA], and an overall

16  increase in the dynamic range of the promoter **(Fig. 2B,F,G,H)**. Thus, both the affinity

17  and number of Gcn4 sites affect a promoter's expression in a manner that depends on

18  the [TF]. In contrast, adding an additional polyT nucleosome disfavoring sequence results

19  in the same fold change in expression at low and high [AA] and no change in the dynamic

20  range of the promoter **(Fig. 2C,I,J,K)**. Adding a binding site for a repressor, changing the

21  position of the binding site, or changing the promoter sequence context to a context with

22  a different predicted nucleosome occupancy also results in no change in the dynamic

7

1    range **(Supplementary Fig. 3)**. Thus, altering the nucleosome occupancy results in a

2    [TF] independent change in expression.

3        Taken together, these results show that sequence mediated expression changes

4    affect the dynamic range of expression when they change binding site affinity or number,

5    but not binding site accessibility, and that both the TF dependent and independent

6    behavior of promoters can be tuned separately.

7

8    **Mutations inside binding sites affect expression in a TF dependent manner**

9        While it is intriguing that addition or removal of entire promoter sequence elements

10    can alter expression either in a TF dependent or independent manner, we wondered if

11    the same independent control can be achieved by single point mutations that are more

12    readily available in an evolutionary context. To determine this, we examined a set of 21

13    3bp scanning mutations made every 3bp across the native HIS3 promoter. We find that

14    19 of these affect expression in a TF independent manner (t-test p=0.83, **Fig. 3A,B**) and

15    that mutations that increase the predicted nucleosome occupancy over the TATA box

16    have lower expression (Pearson R=-0.66 p=9*10$^{-4}$, **Fig. 3C**). Two of the mutations, which

17    fall within the native Gcn4 binding site, appeared to effectively remove response to [AA]

18    change.

19        In addition, we find that systematically mutating the Gcn4 binding site results in a

20    change in dynamic range that is correlated with PSSM score (**Fig. 3D-F, Supplementary**

21    **Fig. 4**). We observe a relatively small increase in expression at high [AA] (Pearson

22    R=0.20, p=0.09), and a much larger increase in expression at low [AA] (Pearson R=0.52

23    p<3*10$^{-6}$), resulting in a net increase in dynamic range with increasing PSSM score

8

1     (Pearson R=0.63 p<$1.3*10^{-8}$ or R=0.77 p<$3*10^{-5}$ when only including values above a

2     previously determined cutoff (Spivak and Stormo, 2012), **Fig. 3E,F, Supplementary Fig.**

3     **4C**). These results are consistent with models in which low-affinity binding sites are

4     always functional, but have a more pronounced effect at high [TF](Carey et al., 2013).

5

6

7     **Maximum expression is set by the amount of active TF and is limited by**

8     **competition for TF molecules**

9     If expression were a simple non-decreasing function of the number of bound TF

10     molecules (Raveh-Sadka et al., 2009; Gertz et al., 2008), we expect expression to

11     increase when either [TF] or number of binding sites in a promoter increases. Thus, a

12     given expression level might be reachable by changing either one or the other, and any

13     promoter, given enough [TF], would be able to reach a level of maximal expression set

14     by the efficiency of transcription initiation. However, this is not what we observe in

15     homotypic promoters. We find that the maximum reachable expression level is

16     determined by [TF] and not by the number of binding sites (**Fig. 4A,B, Supplementary**

17     **Fig. 5**). In all conditions and for all TFs, expression reaches its maximal level at 3-4 sites

18     and then plateaus, decreases, or only slightly increases, depending on the TF, suggesting

19     that this phenomenon is a general consequence of binding site multiplicity and not specific

20     to a particular transcription factor.

21

22     We found that for the set of seven promoters with a single binding site placed at

23     one of seven positions in the promoter, different binding site positions drive different levels

9

1    of expression **(Fig 4C)**. Furthermore, we found that when a binding site that drives high

2    expression (eg: the site at position 51) is added to a promoter with two binding sites

3    (generating a promoter with three sites), expression tends to increase **(Fig. 4D)**. In

4    contrast, when a site that drives low expression (eg: the site at position 93) is added to a

5    promoter with two sites, expression tends to decrease if the expression of the two binding

6    site promoter is already high **(Fig. 4E)**.

7

8         We hypothesized that the observed saturation behavior, which is most pronounced

9    at high [AA] (low [TF]) (**Supplementary Fig. 5)**, is a consequence of competition for

10   limiting TF between binding sites that drive different levels of expression. To compare

11   possible underlying mechanisms we used thermodynamic modeling of gene expression.

12   In short, for each promoter, the model enumerates all possible binding configurations of

13   TF and TBP (TATA binding protein that recruits the transcriptional machinery). The weight

14   of each configuration is based on binding site affinities, Gcn4 concentration and

15   interactions between bound TF molecules, and bound TBP molecules, after which the

16   ratio between weighted TBP bound to TBP unbound configurations determines the

17   expression (see **Methods** for details). We fitted a collection of models of increasing level

18   of complexity to the induction curves of a set of promoters that only contain 0-7 high

19   affinity Gcn4 binding sites. We used a 10-fold cross-validation scheme to assess each

20   model.

21

22        A basic model (see **Methods**) in which binding to each site is independent and

23   each site has either identical contribution to expression or with position-specific driven

10

1    expression, is able to explain an increase in expression with increasing [TF] but does not

2    fit the measured data very well **(Fig. 4F,I,L)**.

3

4        We reasoned that, in order to reproduce the observed saturation, there must be

5    negative interactions between TF binding sites within the same promoter. We examine

6    two alternative mechanisms of binding site interaction: steric hindrance and TF sharing.

7    The steric hindrance model accounts for a previously suggested mechanism in which a

8    bound TF may sterically hinder the binding of a second TF molecule at a neighboring site

9    (Struhl, 1989) by reducing the weight of configurations with multiple bound sites (Raveh-

10   Sadka et al., 2009; Gertz et al., 2008; Giorgetti et al., 2010). The TF sharing model

11   implements competition between neighboring binding sites by dividing the [TF] weight by

12   the total number of binding sites. This mechanism has been observed experimentally,

13   and results from non-specific binding and subsequent 1D sliding; two neighboring binding

14   sites will share their TF capture area and as a consequence have the same effective

15   binding rate as one site (Hammar et al., 2012; Mahmutovic et al., 2015).

16

17       We find that both interaction models can replicate the observed saturation effect,

18   in which, at all [AA], adding a fourth binding site does not result in a large increase in

19   expression (**Fig. 4**). However, quantitatively the TF sharing model better fits the

20   experimental data.

21

22       Taken together, our results suggest that activator binding site multiplicity does not

23   linearly contribute to expression. Our model suggests that this is due to competition

11

1    between binding sites, likely due to neighboring binding sites sharing their capture area

2    as a result of most binding events coming from 1D sliding.

3

4    **Activator binding sites can both increase and decrease expression as predicted**

5    **by a model of TF molecule sharing**

6        The above observations show that multiple binding sites contribute non-linearly to

7    expression. In order to understand the effect of adding or removing individual activator

8    sites in more detail we look at pairs of promoters that differ by only a single binding site.

9    Surprisingly, in 30% of cases adding an additional Gcn4 site reduces expression (**Fig. 5**),

10   and this effect is significantly stronger at high [AA] (55% versus 5% at low [AA], **Fig. 5E,F**),

11   when [TF] is low. However, expression reduction is never below the minimum expression

12   driven by the individual sites (**Supplementary Fig. 6**).

13

14       A comparison of thermodynamic models shows that both steric hindrance and TF

15   sharing can produce expression reduction for activator binding site addition when the

16   added site drives lower expression than the existing site. However, only the TF sharing

17   model shows this effect at low [TF]; steric hindrance shows reduction only when [TF] is

18   high (see **Methods**). Both steric hindrance and TF sharing models predict that the

19   negative interaction between binding sites is stronger at closer distances. Indeed, this is

20   the case in both expression data and in an independent measurement of the same

21   promoter library in which TF binding to promoters was measured *in-vitro* (**Supplementary**

22   **Fig. S7)**. Consistent with the TF sharing model, but not with steric hindrance, this

23   interference is strongest at low TF concentrations both *in-vivo* and *in-vitro*.

12

1

2    The TF sharing model combined with site-specific expression best predicts

3    absolute expression levels as well as synergism, i.e. the change in expression when

4    adding a site (**Fig. 6, Supplementary Fig. 8-11**). In fact, the TF sharing model, given site-

5    specific expression, is the only model tested that can explain expression reduction at high

6    [AA] (low [TF]).

7

8    Taken together, these results show that site addition can either increase or

9    decrease expression. This synergism is concentration dependent. Negative synergism

10   mostly occurs at low [TF], likely ruling out steric hindrance. TF sharing in combination with

11   site-specific expression predicts the observed behavior: more often than not adding an

12   activator-binding site results in a reduction of expression at low [TF].

13

14

15   <span style="color:red">**Discussion**</span>

16   In summary, we presented here a large-scale investigation of the mapping

17   between promoter DNA sequence and dose response curves by measuring the induced

18   gene expression of 6500 designed promoters at six growth conditions in which the

19   regulating TFs are gradually induced.

20   We observe a wide range of dose-response curves in which the dynamic range is

21   altered by changes in the affinity or number of binding sites, and expression is changed

22   induction independently through changes in the accessibility of the promoter.

23

13

1    These results are confirmed by systematic mutations in either the whole promoter

2    or only at the binding site, both affecting overall expression, but only the latter affecting

3    the dynamic range. This suggests that random mutations (that occur more frequently

4    outside of binding sites) are more likely to change overall expression and not the

5    promoter's response.

6

7    Our current and previous (Sharon et al., 2012) observation that expression

8    saturates with increasing number of activator binding sites suggests that either TF binding

9    or pol2 recruitment saturates. However, we observe that while expression cannot be

10   increased by adding binding sites, expression can be increased by increasing [TF]. This

11   argues against  saturation of pol2 recruitment being the cause of the observed saturation

12   of expression level as a function of homotypic binding sites number in each condition. We

13   find that a model that includes competition between binding sites can quantitatively

14   explain our observations.

15

16   We achieved further insight into the non-linear mapping between promoter

17   configuration and dose-response by comparing pairs of promoters that differ by only a

18   single binding site addition. This analysis revealed that at low [TF] adding an activator is

19   more likely to reduce expression than it is to increase expression, suggesting that there

20   is interaction between binding sites.

21

22   Expression of our synthetic Gcn4 targets maxes out at 3-4 binding sites.

23   Interestingly, the vast majority of native Gcn4 targets have 3-4 binding sites (Schuldiner

14

1   et al. 1998). The mechanistic models proposed in this paper may explain the reason for

2   the distribution of binding site numbers in native promoters.

3

4       Our analysis of the observed dose response curves suggest that they are affected

5   mainly by competition for TF (therefore reducing the effective local TF concentration 'seen'

6   by each binding site; referred to as 'TF sharing') rather than steric hindrance between TF

7   molecules. In particular, the two models behave differently with changing [TF]. While

8   steric hindrance will have a stronger effect at high [TF] due to the increased likelihood of

9   bound configurations, 'TF sharing' effects are reduced at high [TF], as the TF is no longer

10  limiting, and this is what we observe.

11

12      To further investigate the possible mechanism that could explain the measured

13  reduction in expression as a function of activator binding site addition, in addition to the

14  thermodynamic model that was fit to data, we developed a toy mathematical model that

15  describes binding site addition from 1 to 2 sites, enabling us to investigate the regimes in

16  which addition will cause a reduction in expression (see supplementary methods). This

17  model shows that expression reduction by steric hindrance will increase with increasing

18  [TF], whereas reduction by TF sharing decreases with increasing [TF]. It is the latter

19  behavior that we observe.

20

21      We note that alternative models are possible; the TF sharing model fits the data

22  but modeling can only show that a given model is wrong, not that a given model is correct.

23  Recently a non-equilibrium promoter-dynamics model was proposed in which TF

15

1    dissociation is fast and actively driven by transcription (Coulon et al. 2013). Our results

2    from the thermodynamic and toy models are independent of assumptions regarding

3    dissociation. Therefore our predictions are independent of whether or not TF unbinding

4    is a an induced non-equilibrium process. One possible alternative model that will

5    reproduce a decrease in expression at high numbers of binding sites, specifically at low

6    [TF], is a combination of additive activation and cooperative repression in which both the

7    activator and repressor compete for the same binding sites. There is evidence suggesting

8    that the transcriptional repressor Mig1 acts cooperatively (Gertz et al. 2009). While no

9    repressors are predicted to bind with high affinity to the Gcn4 binding site (ATGACTCAT),

10   Yap3 and Yap7 are predicted to bind weakly (de Boer et al. 2012). We hypothesized that

11   if the Gcn4 sites have repressive potential, then site addition can cause expression

12   reduction below the level driven by the other sites. We find that, while addition of a binding

13   site often results in expression below the maximum of the expression driven by the

14   individual sites this expression is always greater than the minimum expression driven by

15   the individual sites. The added site can, at most, reduce expression by an amount that

16   the other sites drive, and can never repress beyond that level. In other words, we find that

17   to remove expression first expression has to be added. This is a strong prediction of the

18   TF sharing model and is not predicted by the cooperative repression model.

19

20       A second model that can explain the observation that expression reaches a

21   maximum at around three binding sites is that having more than three bound TF

22   molecules does not increase recruitment of RNA polymerase. It is likely that beyond some

23   number, additional bound transcriptional activators do not contribute to increased

16

1    expression at a single promoter. However, this model cannot qualitatively explain our

2    observation that, for all four TFs, the expression from 3 binding sites is lower at higher

3    [AA] (lower [TF]). The 'activator saturation' model predcts that the same maximal

4    expression could be reached at all amino acid concentrations, but that it might require

5    more binding sites at lower [AA]. This is not what we observe. Moreover, while on average

6    expression saturates at three sites, this is not always the case. Figure 5A shows that

7    going from three to four sites can both increase (green lines) and decrease (red lines)

8    expression; expression rarely remains constant, likely ruling out the 'activator saturation'

9    model.

10

11        Taken together, we have found a strong non-linear mapping between promoter

12    architecture and dose-response, that, by assuming competition between binding sites,

13    we are able to accurately predict from DNA sequence alone. In specific, our model points

14    to a reduction in effective local [TF] (per binding site) due to overlapping capture areas.

15    When [TF] is limiting the effective search time (the time it takes for a TF to find its binding

16    site) is not significantly reduced when another site is added close to an existing one, since

17    search time is dominated by the total capture area. In the regime where [TF] is high more

18    sites bind more TFs and thus have the ability to drive higher expression.

19

20        Our model is also consistent with recent *in-vitro* results performed using the same

21    set of promoters showing that at low [TF], multiple Gcn4 binding sites increase the

22    likelihood of TF binding, but do not increase the number of TF molecules bound to a single

17

1    molecule of promoter, while at high [TF], adding more binding sites does increase the

2    number of bound molecules (Levo et al., 2015).

3

4         Competition for limiting TF, in both one and three dimensional space, may also

5    explain some previously unexplainable results regarding titration by large arrays of

6    extraneous TF binding sites. Lee & Maheshri (Lee and Maheshri 2012) found that

7    contiguous arrays of tetO binding sites bind less TF than do non-contiguous arrays, and

8    that contiguous arrays are less efficient at titrating away TF. Our reanalysis of their data

9    shows that this effect is strongest at low [TF], suggesting that TF sharing may be occurring

10   at these arrays as well, either in 1D space or in 3D space. Splitting the array of decoy

11   binding sites in half results in a larger decrease in expression at low [TF] than at high [TF]

12   (**Supplementary Fig. 12**), as expected from a model in which large number of binding

13   sites spread throughout the genome (at the promoter of interested and at the decoy sites)

14   are sharing a limiting number of TF molecules.

15

16        The yeast genome, which has densely packed genes for a eukaryote, has several

17   promoters (e.g.: Cln3) that are longer than 1kb. Yet 20 TF binding sites could, in theory,

18   be packed into less than 200bp. Intriguingly, the GAL genes, which are highly induced by

19   a large and rapid increase in the active amount of Gal4, tend to have only one or two

20   nucleotides between the sites. In contrast, genes activated at the G1->S transition (e.g.:

21   Cln2) have TF binding sites that are spaced further apart. It was recently shown that Cln3,

22   the protein that activates the TFs bound to the Cln2 promoter is present in limiting

23   concentrations(Wang et al., 2009); the spacing between binding sites may reduce the

18

1     effect of sharing. Binding site spacing is known to be influenced by physical interactions

2     between TFs (Kazemian et al., 2013). Here we suggest that TF sharing between closely

3     spaced binding sites is an additional force acting upon the evolution of promoters. Binding

4     sites for some TFs, especially those with long 1D sliding ranges(Gorman and Greene,

5     2008; Slutsky and Mirny, 2004) may need space for maximal TF occupancy at low [TF].

6     Dense clusters can be used to create a highly responsive behavior (large dynamic range)

7     and less dense clusters might create overall high expression also at low [TF]. Our results

8     suggest that TF sharing can play an important role in determining the response of a

9     promoter to changes in [TF], and therefore influence the evolution of binding site

10     configurations.

11

## Methods

### Promoter library design construction and measurements

14       We used a previously described library of 6500 bar-coded sequences, each with

15     a specific combination of TF binding sites and poly-A tracts(Sharon et al., 2012). These

16     sequences were cloned into a pKT103-derived plasmid upstream of a yellow fluorescent

17     protein (YFP) and transformed into the yeast strain Y8205(Tong and Boone, 2006). We

18     then grew the pooled library in six different growth media each with a different

19     concentration of amino acids: synthetic minimal media with glucose and the amino acids

20     His & Leu with a $2^{11}$, $2^6$, $2^4$, $2^3$, $2^2$ or $2^0$ fold dilution of amino acids. Cells from each growth

21     condition were sorted by FACS according to expression (YFP / mCherry) into 16 bins.

22     DNA from condition and bin was PCR amplified using unique primers so that each

23     resulting sequencing read has three bar codes: promoter, condition, expression bin. The

19

1    distribution of reads across the bins for each promoter and condition enables us to derive

2    an expression level as previously described(Sharon et al., 2014). Thus, we achieve, for

3    each promoter expression across the conditions.

4

5       A Gcn4-GFP ura3::TEFpr-mCherry strain was grown overnight in SCD-HL,

6    resuspended in SCD-His-Leu or SCD, and then the SCD was serial diluated into SCD-

7    HL, resulting in different concentrations of His and Leu. GFP and mCherry were

8    measured using a BD Fortessa flow cytometer using FITC and PE_TexasRed filter sets.

9

10   **Expression normalization**

11      We observed condition specific expression differences that did not appear to stem

12   from biological differences. For example, even the promoters that were not induced (such

13   as Gal4 targets) varied, though slightly, across conditions in a non-monotonic manner.

14   These differences likely stem from day-to-day and experimental variability, as each

15   condition was a separate batch and was sorted on different days. To correct for this effect

16   we subtracted from all promoters the median expression of all Gal4 targets, thus removing

17   this technical variability. All analysis was carried out on the normalized expression values.

18

19   **Growth conditions**

20      Because the two lowest and two highest [AA] conditions induce similar expression

21   we combined them to get a more robust expression measurement. Thus, for the analyses

22   in which we compare low to high [AA] we use the average of the two lowest and the

20

1  average of the two highest [AA] conditions. In the analyses in which we compare four

2  conditions we use the previous two plus the middle two [AA] conditions.

3

## Acknowledgments

9

## Competing Interests

11  The authors declare that they have no competing interests.

## References

13

14  Carey, L.B., van Dijk, D., Sloot, P.M.A., Kaandorp, J.A., and Segal, E. (2013). Promoter

15  sequence determines the relationship between expression level and noise. *11*,

16  e1001528–e1001528.


17  Gasch, A.P., Spellman, P.T., Kao, C.M., Carmel-Harel, O., Eisen, M.B., Storz, G.,

18  Botstein, D., and Brown, P.O. (2000). Genomic expression programs in the response of

19  yeast cells to environmental changes. Mol. Biol. Cell *11*, 4241–4257.


20  Gertz, J., Siggia, E.D., and Cohen, B.A. (2008). Analysis of combinatorial cis-regulation

21  in synthetic and genomic promoters. Nature *457*, 215–218.

1   Gorman, J., and Greene, E.C. (2008). Visualizing one-dimensional diffusion of proteins

2   along DNA. Nat. Struct. Mol. Biol. *15*, 768–774.

3   Hammar, P., Leroy, P., Mahmutovic, A., Marklund, E.G., Berg, O.G., and Elf, J. (2012).

4   The lac repressor displays facilitated diffusion in living cells. Science (New York, NY) *336*,

5   1595–1598.

6   Kaplan, N., Moore, I.K., Fondufe-Mittendorf, Y., Gossett, A.J., Tillo, D., Field, Y., LeProust,

7   E.M., Hughes, T.R., Lieb, J.D., Widom, J., et al. (2009). The DNA-encoded nucleosome

8   organization of a eukaryotic genome. Nature *458*, 362–366.

9   Kazemian, M., Pham, H., Wolfe, S.A., Brodsky, M.H., and Sinha, S. (2013). Widespread

10  evidence of cooperative DNA binding by transcription factors in Drosophila development.

11  Nucleic Acids Res. *41*, 8237–8252.

12  Levo, M., Zalckvar, E., Sharon, E., Dantas Machado, A.C., Kalma, Y., Lotam-Pompan,

13  M., Weinberger, A., Yakhini, Z., Rohs, R., and Segal, E. (2015). Unraveling determinants

14  of transcription factor binding outside the core binding site. Genome Res. gr.185033.114.

15  Ljungdahl, P.O., and Daignan-Fornier, B. (2012). Regulation of amino acid, nucleotide,

16  and phosphate metabolism in Saccharomyces cerevisiae. Genetics *190*, 885–929.

17  Mahmutovic, A., Berg, O.G., and Elf, J. (2015). What matters for lac repressor search in

18  vivo-sliding, hopping, intersegment transfer, crowding on DNA or recognition? Nucleic

19  Acids Res. *43*, 3454–3464.

20  Natarajan, K., Meyer, M.R., Jackson, B.M., Slade, D., Roberts, C., Hinnebusch, A.G., and

22

Marton, M.J. (2001). Transcriptional profiling shows that Gcn4p is a master regulator of gene expression during amino acid starvation in yeast. Mol. Cell. Biol. *21*, 4347–4368.

Rajkumar, A.S., Dénervaud, N., and Maerkl, S.J. (2013). Mapping the fine structure of a eukaryotic promoter input-output function. Nat Genet *45*, 1207–1215.

Raveh-Sadka, T., Levo, M., and Segal, E. (2009). Incorporating nucleosomes into thermodynamic models of transcription regulation. Genome Res. *19*, 1480–1496.

Sharon, E., Kalma, Y., Sharp, A., Raveh-Sadka, T., Levo, M., Zeevi, D., Keren, L., Yakhini, Z., Weinberger, A., and Segal, E. (2012). Inferring gene regulatory logic from high-throughput measurements of thousands of systematically designed promoters. Nat. Biotechnol. *30*, 521–530.

Sharon, E., van Dijk, D., Kalma, Y., Keren, L., Manor, O., Yakhini, Z., and Segal, E. (2014). Probing the effect of promoters on noise in gene expression using thousands of designed sequences. Genome Res.

Slutsky, M., and Mirny, L.A. (2004). Kinetics of protein-DNA interaction: facilitated target location in sequence-dependent potential. Biophys. J. *87*, 4021–4035.

Spivak, A.T., and Stormo, G.D. (2012). ScerTF: a comprehensive database of benchmarked position weight matrices for Saccharomyces species. Nucleic Acids Res. *40*, D162–D168.

Struhl, K. (1989). Molecular mechanisms of transcriptional regulation in yeast. Annu. Rev. Biochem. *58*, 1051–1077.

23

1   Struhl, K. (1995). Yeast transcriptional regulatory mechanisms. Annu. Rev. Genet. *29*,

2   651–674.


3   Struhl, K. (1992). 31 Yeast GCN4 Transcriptional Activator Protein. Cold Spring Harbor

4   Monograph Archive; Volume 22B (1992): Transcriptional Regulation 2.


5   Tong, A.H.Y., and Boone, C. (2006). Synthetic genetic array analysis in Saccharomyces

6   cerevisiae. Methods Mol. Biol. *313*, 171–192.


7   Wang, H., Carey, L.B., Cai, Y., Wijnen, H., and Futcher, B. (2009). Recruitment of Cln3

8   cyclin to promoters controls cell cycle entry via histone deacetylase and other targets.

9   PLoS Biol. *7*, e1000189.


10   De Boer, Carl G., and Timothy R. Hughes. 2012. "YeTFaSCo: A Database of Evaluated

11   Yeast Transcription Factor Sequence Specificities." *Nucleic Acids Research* 40

12   (Database issue): D169–79.

13

14   Gertz, Jason, Eric D. Siggia, and Barak A. Cohen. 2009. "Analysis of Combinatorial Cis-

15   Regulation in Synthetic and Genomic Promoters." *Nature* 457 (7226). Nature Publishing

16   Group: 215–18.

17

18   Giorgetti, L., Siggers, T., Tiana, G., Caprara, G., Notarbartolo, S., Corona, T., Pasparakis,

19   M., Milani, P., Bulyk, M. L. & Natoli, G. Noncooperative interactions between transcription

20   factors and clustered DNA binding sites enable graded transcriptional responses to

21   environmental inputs. *Mol. Cell* **37,** 418–428 (2010).

22

24

1    Schuldiner, O., C. Yanover, and N. Benvenisty. 1998. "Computer Analysis of the Entire

2    Budding Yeast Genome for Putative Targets of the GCN4 Transcription Factor." *Current*

3    *Genetics* 33 (1): 16–20.

4

5    Coulon, Antoine, Carson C. Chow, Robert H. Singer, and Daniel R. Larson. 2013.

6    "Eukaryotic Transcriptional Dynamics: From Single Molecules to Cell Populations."

7    *Nature Reviews. Genetics* 14 (8). Nature Publishing Group: 572–84.

8

9    Lee, T.-H. & Maheshri, N. A regulatory role for repeated decoy transcription factor binding

10   sites in target gene expression. Mol. Syst. Biol. 8, 576 (2012).

11

1   **Figures**
2
3



4
5

6   **Figure 1. Measurements of TF concentration dependent expression for thousands**

7   **of designed promoters. (A)** Schematic depiction of the experimental design. A pooled

8   library of 6500 designed promoters was transformed into yeast and expression levels of

9   all strains in the pooled library were measured in minimal media at each of six different

10  amino acid concentrations (see Methods). **(B)** Promoter expression measurements

11  sorted by dynamic range. For each promoter in the library we obtain an expression

12  measurement at each of the six AA concentrations. For promoters that lack Gcn4, Leu3,

13  Bas1 or Met31 sites, expression does not change with decreasing AA concentration (top

14  of **B**). For promoters with multiple Gcn4 binding sites, expression increases with

15  decreasing [AA]. **(C)** Shown are four representative induction curves showing the effect

16  of changing the number of Gcn4 binding sites (cyan, green, blue), or adding a polyT

17  nucleosome disfavoring sequence (green, red). IDs show library construct identifiers.

26

**Figure 2. The effect of Gcn4 binding site number and polyT nucleosome disfavoring sequences on [TF] dependent and independent expression.** (**A-C**) show expression at high [AA] (X-axis) versus expression at low [AA] (Y-axis) for various promoter sequence features. Dashed lines are the diagonal (slope=1) line that best fit each category of promoters. The black dashed diagonal line (Y=X) represents the regime where expression is constant across conditions. The vertical distance from the Y=X line measures how much any one promoter changes in expression across conditions. Density plots (using ks denstiy estimation) at the X-axis and Y-axis show the distributions of expression values for each promoter at high and low [AA] respectively. (**D-K**) show expression and expression fold change (Y-axis) in box plots as a function of promoter sequence features (X-axis). The dashed black lines in connects the medians of each box. Asterisks denote statistically significant (t-test, p<0.01) changes between subsequent groups. (**A**) Shown are promoters grouped by the number of Gcn4 binding sites. (**B**)

27

1    Shown are promoters with either low or high affinity Gcn4 sites. **(C)** Shown are promoters

2    with either one or two polyT nucleosome disfavouring sequences. **(D,E)** Box plots of the

3    data in (**A**). Promoters are grouped by the number of binding sites. (**D**) Shown is

4    expression at high [AA] (Y-axis). (**E**) Shown is expression fold change (dynamic range,

5    log2(low [AA] / [high AA])). **(F-H)** Box plots of the data in (**B**). Shown are expression at

6    high [AA] (**F**), low [AA] (**G**) and expression fold change (dynamic range) (**H**) for promoters

7    with low or high affinity binding sites. All differences are statistically significant (t-test

8    p<1e-3, p<1e-5, p<1e-4 for **F**,**G**,**H** respectively). **(I,J,K)** Box plots of the data in (**C**). Shown

9    are expression at high [AA] (**I**), low [AA] (**J**) and expression fold change (dynamic range)

10   (**K**) for promoters with either one or two polyT sequences. Expression at high and low

11   [AA] show significant change as a function of polyT number (t-test p<1e-4, p<1e-4

12   respectively), however dynamic range does not change significantly (t-test p=0.77).

13

28

1



2

**Figure 3. The effect of promoter and binding site mutations on [TF] dependent and**

**independent expression. (A-C)** Expression data for a set of 20 sequences that only

differ by a single point mutation made at different locations along the promoter. For each

sequence the nucleosome occupancy over the TATA box was predicted using a

thermodynamic model(Kaplan et al., 2009). **(A)** Shown is the effect of promoter point

mutations on expression at low and high AA concentration; the dashed blue line is the

best fit of a line with slope=1 to the data. The two scanning mutations that affect the Gcn4

binding site are identified with arrows. These mutations result in removal of the effect of

[TF] on expression change. **(B)** The same data as in (A) with dynamic range (log2(low

[AA] / [high AA])) graphed against expression at high [AA]. **(C)** Mutations that increase

the predicted nucleosome occupancy over the TATA box decrease expression in a linear

manner (fit red line). **(D)** A set of promoters, each with a single Gcn4 binding site, with

29

1    single base mutations that affect the predicted affinity. Arrows mark the consensus

2    sequence and the native His3 promoter sequence. Point mutations in a set of promoters

3    with a single Fhl1 binding site are shown as a control. The dashed red line (slope=1)

4    marks the sequence that gives the highest dynamic range in expression. **(E)** The same

5    data as in (D) with dynamic range graphed against expression at high [AA]. **(F)** Mutations

6    predicted to have a high Gcn4 binding site affinity have a high dynamic range (dashed

7    red line), but the PSSM has less predictive power at low scores.

8

**Figure 4. A model that incorporates TF sharing with specific position- expression can best explain expression across all amino acid concentrations. (A)** The library contains of promoters with identical Gcn4 binding sites placed at one of seven locations in the promoter. **(B)** Shown are the measured expression levels (Y-axis) as a function of binding site number (different colors) at four AA concentrations (different groups along the X axis) for Gcn4. Each box contains data for all promoters with that number of binding sites and no other features (eg: no nucleosome disfavoring sequences or binding sites for other TFs). The black line shows the median expression level for all promoters with that number of binding sites. **(C)** Shown is the expression for each promoter with a single Gcn4 binding site, normalized so that all conditions have the same mean expression. **(D,E)** Shown is the effect of adding a third binding site (at position 51 or position 93) to a promoter that already has two binding sites. The expression of the two binding site promoter (x-axis) is graphed against the three binding site promoter (y-axis). **(F-L)** Each point shows a single promoter measured at one of four conditions (blue, green, red, cyan in decreasing [AA] order) (x-axis) and the predicted expression levels (y-axis) of that promoter, for the six different models, fitted in cross-validation to the data shown in (A), which are promoters with 1 to 7 high affinity Gcn4 binding sites (ATGACTCAT). $R^2$ values were computed for absolute predicted expression on the test data. Each model includes either position specific expression (a unique weight is associated with each unique binding site position) or non-specific expression (all binding site positions share the same weight), and either no interaction, steric hindrance (a negative weight for multiple bound configurations) or TF sharing (the [TF] weight is divided by the number of sites).
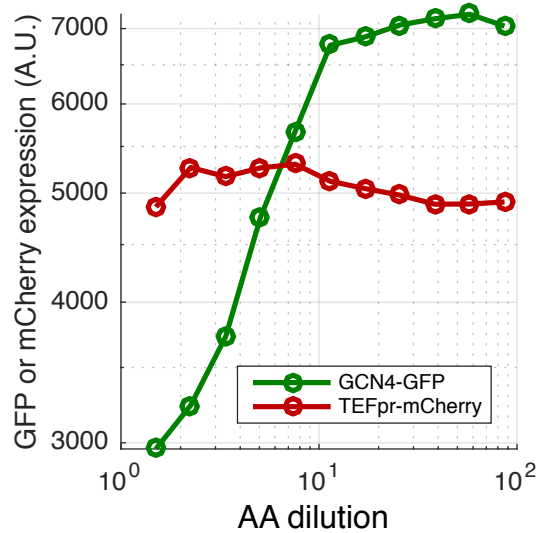
31

1



**Figure 5. At low [TF], addition of a binding site often results in a reduction in expression. (A,B)** Measured expression (and expression change, lines) when adding Gcn4 sites, for high (A) and low (B) [AA] Lines connect promoters that differ by only one binding site. Green lines indicate an expression increase when adding a binding site, red lines indicate that expression decreases. **(C,D)** A scatter plot showing the effect of adding a single Gcn4 site 'B' to the 'A' promoter, where the 'A' promoter has 0-6 binding sites (differently colored points) and the 'AB' promoter has (0-6)+1 sites. If the point is on the diagonal there is no effect of adding an additional binding site. For each AB promoter, synergism is the y-axis distance from the diagonal line. **(E,F)** Box-plots quantifying synergism — the log2 fold-change in expression when adding a binding site to promoters with 0 to 6 sites, for high (E) and low (F) [AA].

32

**Figure 6. TF sharing but not steric hindrance can explain the decrease in expression due to activator binding site addition. (A-F)** Predicted expression as a function of binding site number for six different thermodynamic models, fitted in cross-validation to the Gcn4 measured data. Green lines show a predicted increase in expression upon binding site addition; red lines show a predicted decrease. $R^2$ values were computed for absolute predicted expression on the test data. **(G-I)** Measured versus predicted expression and synergism for the best model at low and high [AA].

33

1

2 **Supplementary Figure 1. Amino acid starvation induces Gcn4 transcription factor**

3 **levels.** Shown is a serial dilution of amino acid concentration and the levels of Gcn4-GFP

4 and a control TEF promoter driving mCherry. The Gcn4-GFP-HIS3 *ura3*::TEFpr-mCherry-

5 URA3 strain is from the GFP collection. Cells were grown in SCD, washed, diluted into

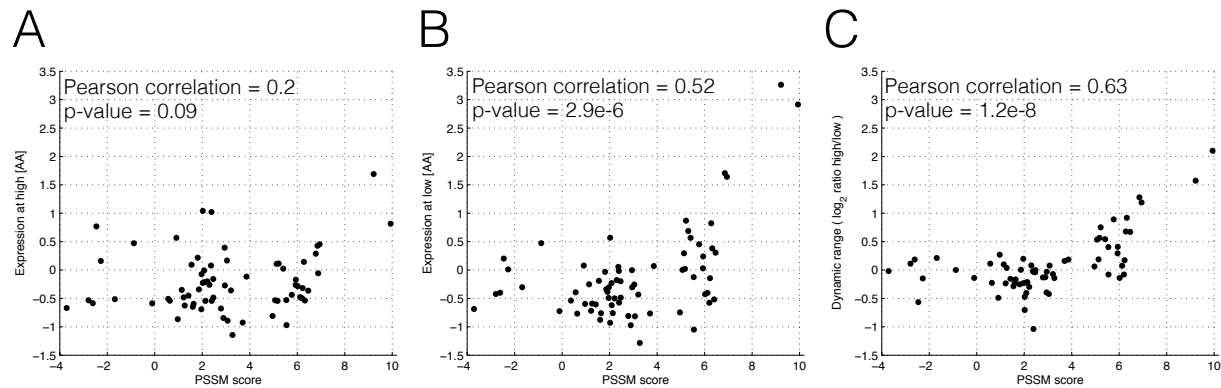6 varying concentrations of SCD-HL media, and grown for eight hours.

1



2

**Supplementary Figure 2. Reproducibility between biological replicates and between isolated and high-throughput measurements.** Shown are biological replicates of the high-throughout measurements (A, B) and isolated strains measured in plate reader (C). (A) shows condition 1 versus 2 and (B) shows condition 5 versus 6. These conditions were measured on different days. (C) shows 96 promoters that were isolated form the pooled library and measured using a plate-reader in high and low [AA] (equivalent to conditions 1 and 6). The X-axis shows expression measured in the high-throughput pooled experiment and the Y-axis shows the expression measured using the plate-reader for individual strains.
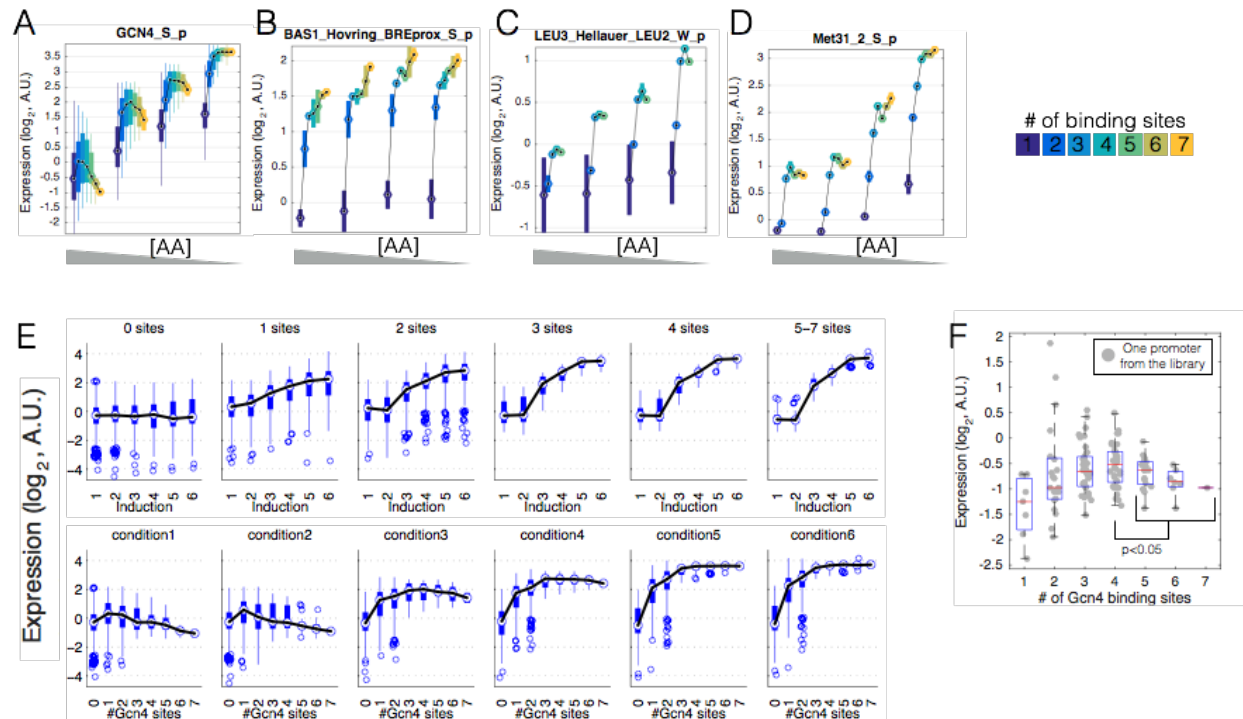
1



2

3

4 **Supplementary Figure 3. Changing the promoter sequence context or addition of a**

5 **binding site for a repressor changes expression without changing dynamic range.**
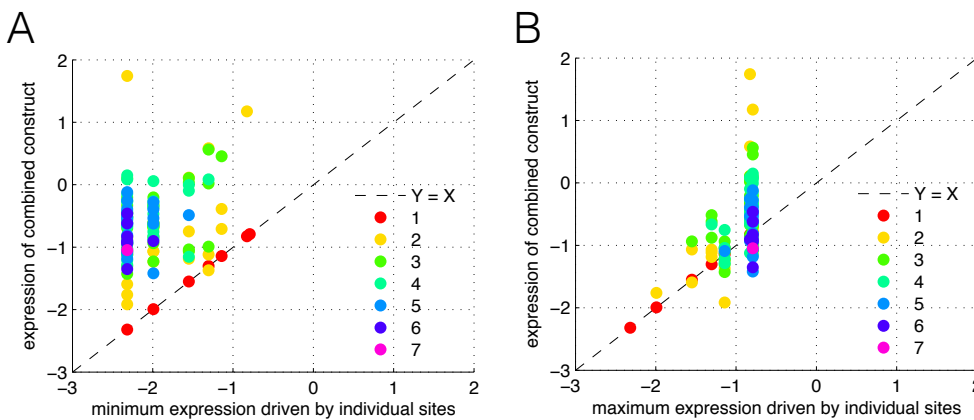
6 **(A)** Shown are expression at high and low [AA] for a set of promoters with a single Gcn4

7 binding site placed at different locations with the HIS3 (red, low nucleosome occupancy)

8 or GAL1,10 (blue, high nucleosome occupancy) promoter context. (B) Shown are

9 expression at high and low [AA] for a set of promoters with a single Gcn4 binding site

10 placed at different locations (blue) and the same promoters with a Mig1/2 repressor-

11 binding site added to the -36 position in the promoter. Addition of a repressor binding site

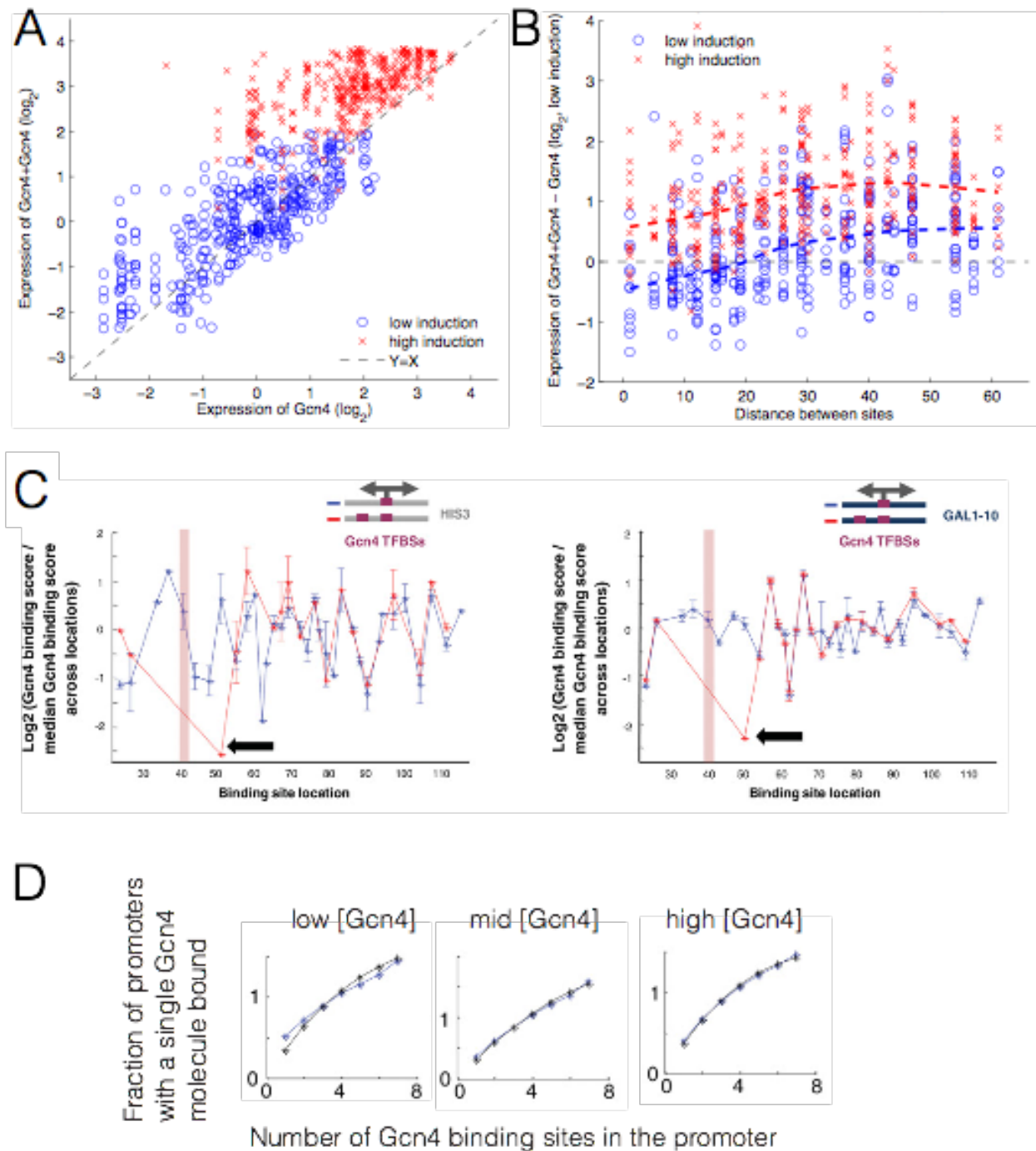12 moves expression along the line X=Y, and therefore affects [TF] independent expression.

13

14

15

36

1

2

3   **Supplementary Figure 4. Mutations in the Gcn4 binding site that decrease the**

4   **PSSM score decrease [TF] dependent expression. (A)** Expression at high [AA] for a

5   set of promoters that differ only by the sequence at the single Gcn4 binding site in the

6   promoter. **(B)** Expression at low [AA] for a set of promoters that differ only by the

7   sequence at the single Gcn4 binding site in the promoter. **(C)** Dynamic range for each of

8   the promoter calculated from the data shown in (A) and (B).

9

1

2



3

4

5 **Supplementary Figure 5. Expression saturates with binding site number for multiple**

6 **transcription factors. (A-D)** Shown are the measured expression levels (Y-axis) as a

7 function of binding site number at four AA concentrations (see Methods) for Gcn4, Bas1,

8 Leu3 and Met31 TFs. Each box is the set of promoters with 0 to 7 binding sites (different

9 colors) measured at one of four conditions (x-axis). Black lines show the median expression

10 per condition along the number of binding sites. **(E)** Various ways of graphing the Gcn4

11 binding site data. **(F)** Expression decreases when changing from 4 to >4 binding sites.
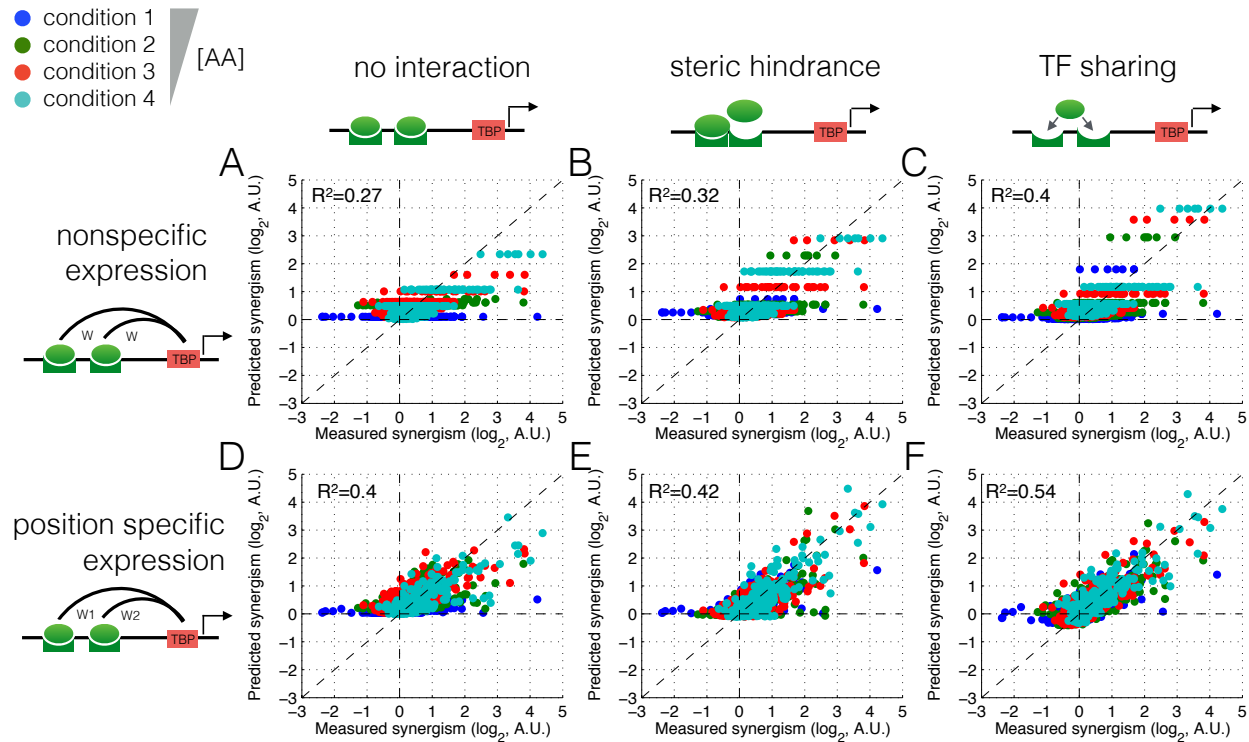
12

13

38

1
2
3



4

**Supplementary Figure 6. Repression is never below the minimum expression driven by the individual sites.** Shown is the minimum (A) and maximum (B) expression of the individual sites (X-axis) versus the expression of the promoter in which the sites are combined (Y-axis). Thus, each point is a promoter with N sites (1-7, colors) for which the expression is shown on the Y-axis. On the X-axis the minimum or maximum expression is shown computed on the expression values of the set of promoters that have only one binding site, namely the same as is present in the promoter that is depicted, e.g. for a promoter with 3 sites S1, S2 and S3 the minimum or maximum is computed on the expression values of the promoters that have only site S1, S2 or S3.
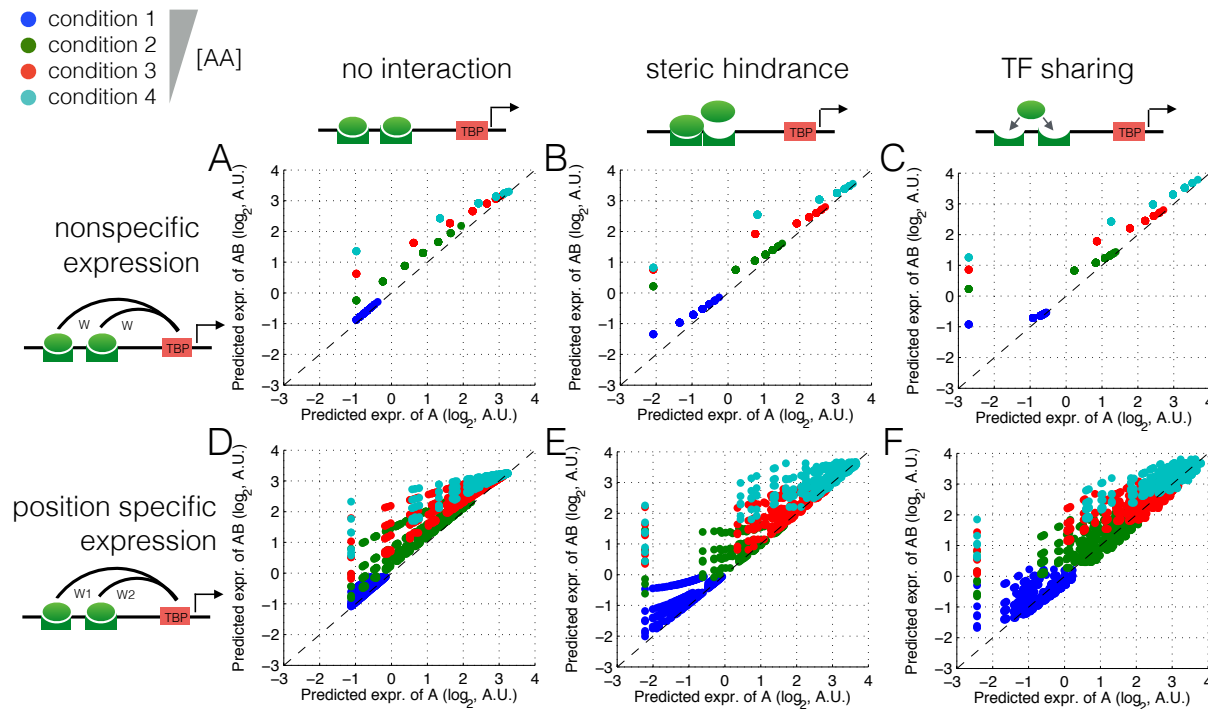
39

1

**Supplementary Figure 7. In-vivo expression and in-vitro binding increase less than expected when adding a second near-by binding site. (A)** Expression change as a function of adding a second Gcn4 site. Each point shows expression for a promoter with one Gcn4 site (X-axis) and the same promoter with a second site added (Y-axis). **(B)** Expression change for each promoter pair is shown (Y-axis) as a function of the distance
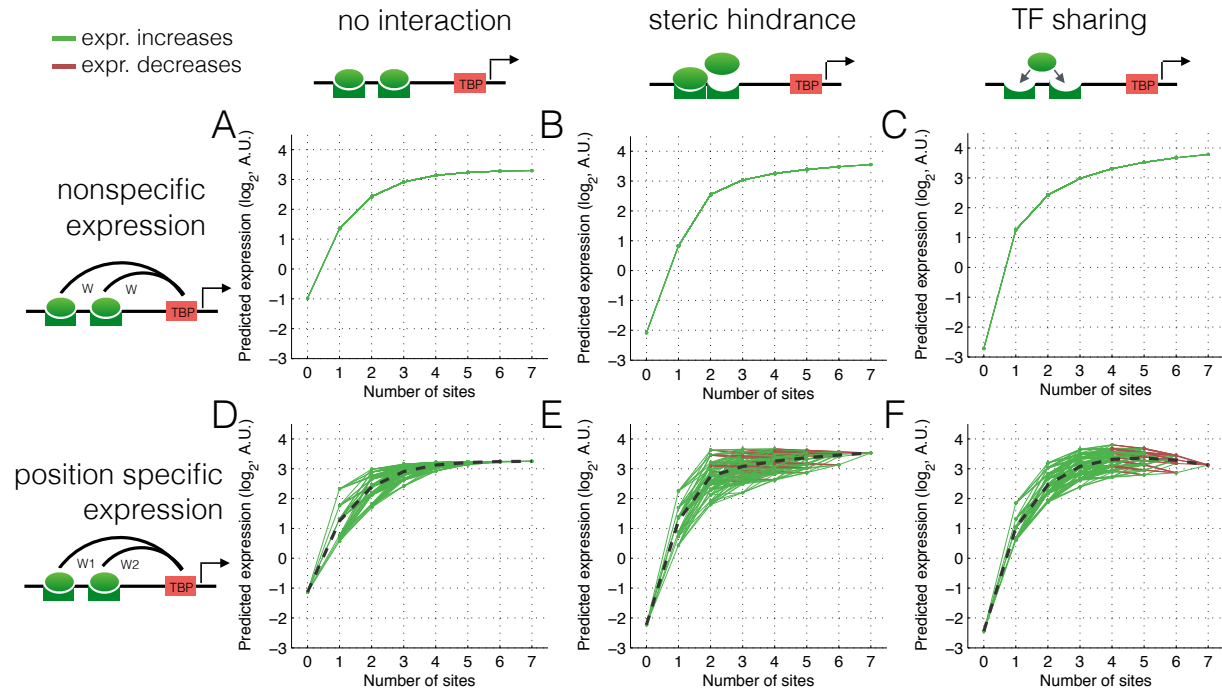
40

1    between the two sites (X-axis). Dashed lines show the median change in expression

2    across all binding site additions. **(C) Levo et al. 2015. Figure 3.** A set of sequences with

3    a strong Gcn4 site placed at different locations along a specific sequence context, either

4    in the presence of an additional strong Gcn4 site located in a fixed location (with the pink

5    rectangle marking the location of the center of this site) (in red) or without this additional

6    site (in blue). Shown is the $\log_2$ of the ratio of the binding score attained by each sequence

7    (with the x-coordinate marking the location of the center of the site) divided by the median

8    binding score across all sequences in this set. The black arrow points to a sequence

9    where the 9-bp sites are separated by a single bp. Sequences with Gcn4 TFBSs of 9 bp

10    placed along the HIS3-derived context (left panel) and along the GAL1-10–derived

11    context (right panel). **(D) Levo et al. 2015. Figure 3.** For a set of sequences with all

12    possible combinations of one to seven binding sites for Gcn4 in seven predefined

13    locations, the average frequency of sequences with a single molecule of bound Gcn4 is

14    shown as a function of the number of sites within the sequence (in blue). The predictions

15    of these dependencies are based on a simple thermodynamic model assuming multiple

16    TF binding events are independent and are also plotted (in black). At low [TF] (left panel),

17    the amount of bound Gcn4 increases less slowly than is predicted from a thermodynamic

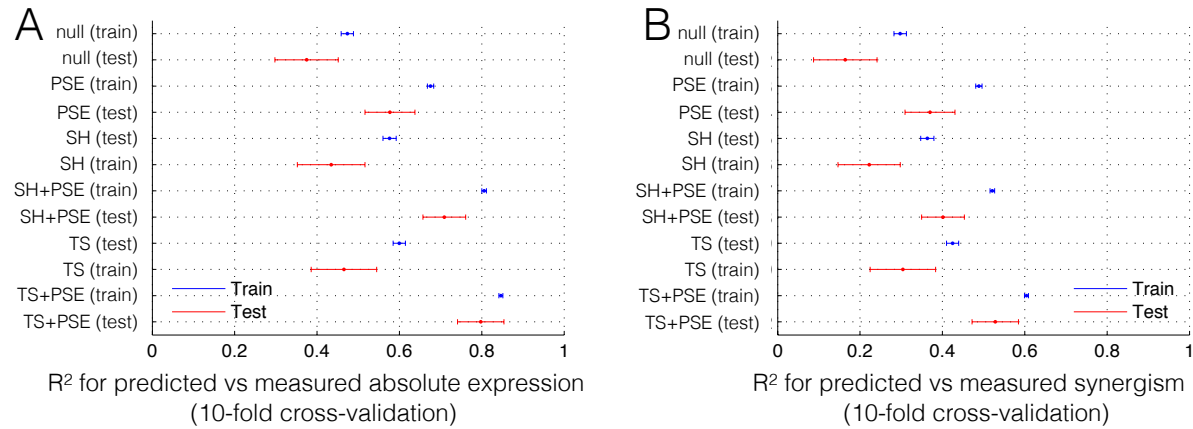18    model lacking negative interactions between TF binding sites.

19

20
21
22

41

1

2

3  **Supplementary Figure 8. TF sharing but not steric hindrance can best explain**

4  **synergism at all amino acid concentrations.** **(A-F)** Measured versus predicted

5  synergism for each [AA] (different colors) for six different thermodynamic models, fitted in

6  10-fold cross-validation to the measured data. Note that TF sharing is the only model that

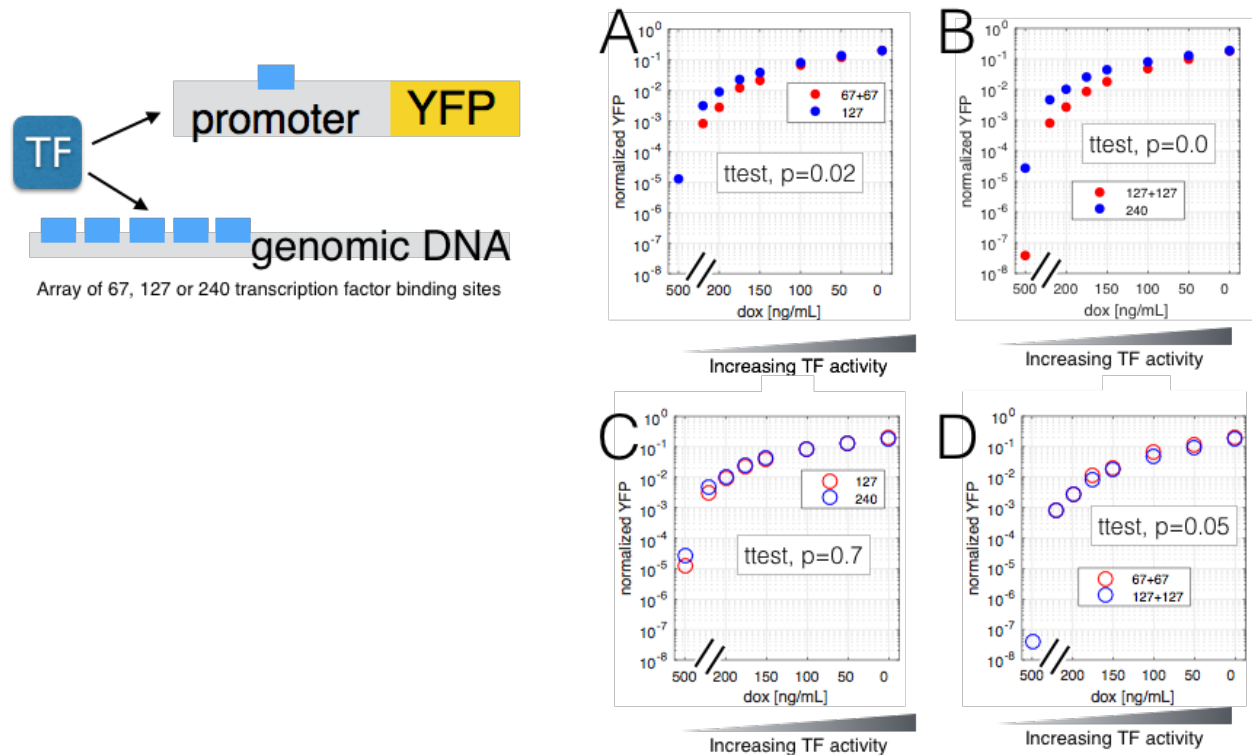7  can generate negative synergism.

8

1

2

3 **Supplementary Figure 9. Only the TF sharing model is capable of producing strong**

4 **negative synergism.** Shown is predicted expression of a single binding site (A) versus

5 expression following the addition of another binding site (AB), for all six thermodynamic

6 models. The top row of models have fewer points visible because they don't have

7 position-specific expression, so all promoters that differ only by the position of binding

8 sites will have the same expression.

9

1

2

3  **Supplementary Figure 10. TF sharing and steric hindrance can explain expression**

4  **saturation and small expression reductions at low [AA]. (A-F)** Predicted expression

5  at low [AA] is shown as a function of binding site number for six different thermodynamic

6  models, fitted in cross-validation to the Gcn4 measured data. Green lines show a

7  predicted increase in expression upon binding site addition; red lines show a predicted

8  decrease.

9

1

2

3 **Supplementary Figure 11. The cross-validated thermodynamic models do not**

4 **exhibit excessive over-fitting.** All six thermodynamic models were fit to absolute

5 expression level at all induction points. Shown is the $R^2$ for predicted versus measured

6 expression level **(A)** and synergism **(B)** for each training (blue) and test (red) data set.

7

8

1

2   **Supplementary Figure 12. Addition of large arrays of excess binding sites titrate**

3   **away limiting TF molecules independent of the number of binding sites in the array.**

4   Lee & Mahreshi used a strain in which YFP is activated by the tTA transcription factor,

5   which is inhibited by doxycycline. Addition of doxycycline reduces the effective TF

6   concentration. Into this strain they integrated into the genome either one or two arrays of

7   extraneous tTA binding sites, with the arrays having 67, 127 or 240 binding sites. Data

8   from Figure 3 were replotted to compare strains with approximately equal numbers of

9   binding sites **(A,B)** or equal numbers of loci containing extraneous sites **(C,D)**. Splitting

10  the binding sites across two loci increases the distance between the binding sites and

11  results in a large decrease in expression (increase in titration of TF). However, adding

12  extra binding sites to an existing locus results in little or no decrease in expression (no

13  increase in TF titration).

46

# 1    Supplement

11


## 12    Thermodynamic model fitted to data

13    Our aim is to predict gene expression from promoter DNA sequence features across several transcriptional-

14    activator concentrations. In specific, our input sequences are 128 promoters that contain 0 to 7 activator

15    Gcn4 binding sites in 7 predefined positions for all ($2^7$) possible configurations. We use a thermodynamic

16    model (Shea & Ackers 1985, Bucheler Hwa 2003, Gertz Cohen 2009) that predicts gene expression by

17    enumerating all promoter configurations and computing the probability of TBP binding (*P(TBP)*), which then

18    is converted to gene expression via a sigmoid function (Eq.1).

19

20                   $$Expression = sigmoid\big(E_{min}, E_{max}, TBP_{mid}, H, P(TBP)\big) \quad (1)$$

21

22    Where $E_{min}$ and $E_{max}$ are the minimum and maximum expression respectively, $TBP_{mid}$ is the $P(TBP)$ at

23    which expression is half maximal and $H$ is the hill-coefficient and describes how switch-like expression is

24    as a function of TBP binding. This sigmoid function describes both transcription and translation.

47

1

2 Both expression (as a function of P(TBP)) and [TF] (as a function of growth condition) are modeled with the

3 same sigmoid function (Eq.2).

4

5 $$sigmoid(ymin, ymax, xmid, h, x) = ymin + \frac{(ymax - ymin)}{1 + (\frac{x}{xmid})^{-h}} \quad (2)$$

6

7 The probability of TBP binding is computed through the relative weight of TBP-bound to TBP-unbound

8 configurations (Eq.3).

9

10 $$P(TBP) = \frac{\sum_c W_c \delta(TBP)}{\sum_c W_c} \quad (3)$$

11

12 Where $W_c$ is the statistical weight of configuration $c$, which is computed as (Eq.4):

13

14 $$W_c(sample = s) = q_{TBP}\delta_{TBP} + \sum_{i=1}^{\#TFBS}\left(q_i([TF_s]) + W_{i,TBP}\delta_{TBP}\right)\delta_i + \sum_{i=1}^{\#TFBS}\sum_{j=1}^{\#TFBS} W_{i,j}\delta_i\delta_j \quad (4)$$

15

16 Where $\delta_i$ is an indicator function of whether the TF is bound in site $i$, $W_{i,j}$ reflects an interaction between

17 bound TF molecules and would be positive for cooperative interaction. In our case we use this as a negative

18 weight to model steric hindrance. $W_{i,TBP}$ is the interaction weight of bound TF to bound TBP and is either

19 specific to each site position or shared across all positions. This weight can be interpreted as the

20 contribution of a bound TF to initiating transcription.

21

22 $$q_i([TF_s]) = \Delta\Delta G + \log\left(\frac{[TF_s]}{\#TFBS}\right) \quad (5)$$

23

24 $\Delta\Delta G$ is the binding affinity per position. The change of Gibbs free energy is similar between sites (same

25 sequence motif) thus we assume it to be constant across positions.

26

48

1 When we assume that sites act independently or for steric hindrance #TFBS = 1. When we assume TF

2 sharing we set #TFBS to the number of binding sites per promoter.

3

4 Finally we assume that the TF concentration follows a sigmoid as a function of the conditions (Eq.6).

5

6
$$[TF_s] = sigmoid([TF]_{min}, [TF]_{max}, [AA]_{mid}, H, [AA]) \qquad (6)$$

7

8 Where $[AA]$ is the amino acid concentration that is varied across conditions (see Methods). $[TF]_{max}$ and

9 $[TF]_{min}$ are the maximal and half-maximal TF concentrations. $H$ is the hill-coefficient.

10

11 For the steric hindrance model we use a negative weight for TF-TF interactions. Thus, $W_{i,j}$ (Eq.4) becomes

12 $W_{coop}$ (see model parameters below). For the TF sharing model we compute an effective TF concentration

13 ($[TF]_{eff}$) by dividing the $[TF]$ by number of binding sites. The amount of sharing is set by $W_{coop}$ as follows:

14 $[TF]_{eff} = [TF] / (N_{tf} * W_{coop})$, where $N_{tf}$ is the number of binding sites.

15

16 **Model parameters**

17 Next we describe the free parameters of the model that were fitted in cross-validation.

18 Accompanying data files (*_params.tab) contain the fitted values for each model.

19 *Binding affinities:*

20 Q_GCN4:        binding affinity of Gcn4  (for sequence motif TGACTCA)

21 Q_TBP:        binding affinity of TBP

22 *TF concentration:*

23 C_max:        Maximum [TF]

24 C_mid:        [AA] for which [TF] is half-maximal

25 C_H:        Hill-coefficient for [TF] as a sigmoid function of [AA]

49

1 C_min:        the minimum [TF] which is fixed to 1, since this parameter is redundant with the binding

2 affinity of Gcn4.

3 *Interaction weights:*

4 W_GCN4_TBP: Interaction weight of bound Gcn4 with bound TBP

5 For position specific expression this is modeled with 7 separate parameters:

6 W_GCN4_TBP_PosX, where X = {22,36,50,64,78,92,106}, which is the position relative to the ATG.

7 *Expression from TBP occupancy:*

8 TBP2Exp_min: minimum expression

9 TBP2Exp_max: maximum expression

10 TB2Exp_mid:    (*P*)TBP for which expression is half-maximal

11 TBP2Exp_h:     Hill-coefficient

12 *Steric-hindrance:*

13 W_coop:        (Wtf-tf) TF-TF interaction weight

14 *TF sharing:*

15 W_coop:        (Wtf-tf) weight that determines to what extent [TF] is shared between neighboring sites.

16

## 17 **Toy model of activator site addition**

18       Both steric hindrance and TF sharing may explain, in certain regimes, reduction of

19 expression from adding an additional binding site. To better understand how and when

20 these models that include interactions between binding sites predict expression reduction,

21 we developed a toy thermodynamic mathematical model that describes expression

22 change as a result from binding site addition. We model the expression change from one

23 to two activator binding sites by enumerating all possible binding configurations of the

50

1 TFs and compute expression from the configurations that have at least one TF bound.

2 Each of the two possible binding sites has its unique contribution to expression (i.e.

3 interaction with the transcriptional machinery). The weight of each configuration is

4 computed from the weight of the TFs (affinity * concentration) (parameter $Wtf$). Per

5 configuration expression is the sum of the product of TF weights and site-specific

6 expression. Total promoter expression is then computed from the sum of expression over

7 all configurations. The unbound configuration has weight 1 (arbitrary constant). We then

8 assume (1) expression driven by each site is independent, or (2) steric hindrance: the

9 double bound configuration has weight 0, or (3) TF sharing, in which, for the promoter

10 with two sites, [TF] = [TF]/2. Next we compute expression change from site addition by

11 subtracting the expression from two sites from the single site promoter. To determine

12 which parameter regimes enable expression to decrease, we solve this equation and get

13 the boundary condition in which site addition does not change expression. We can then

14 easily find the regimes in which expression increases or decreases, namely when the

15 expression driven by the second site (E2) is higher or lower, respectively, than the

16 boundary condition.

17

18      Thus, the expression change as a result of addition of a second site, for

19 independent expression is:

20

21      $Expression\ change\ due\ to\ addition\ of\ a\ second\ site = \frac{E2*Wtf}{Wtf+1}$     (7)

22

23 This value is always positive, thus ruling out independence. For steric hindrance,

1

$$Expression\ change\ = \frac{Wtf*(E2 - E1*Wtf + E2*Wtf)}{(2*Wtf + 1)*(Wtf + 1)} \quad (8)$$

3

4 Solving gives us the boundary condition $E2 = \frac{E1*Wtf}{Wtf+1}$ for which expression does not

5 change. At low [TF] (low Wtf) this limit is low, thus at low [TF] we can only reduce

6 expression when E2 << E1. This limit goes to E2 = E1 when [TF] goes to infinity, thus at

7 high [TF] it becomes easier to reduce expression, i.e. only a small difference between E1

8 and E2 is necessary. So, even at high [TF] reduction will occur when the added site drives

9 lower expression. This is not what we observe. Experimentally, the frequency of reduction

10 goes decreases at high [TF], steric hindrance predicts the opposite, ruling out steric

11 hindrance.

12

13 Finally for TF sharing we get

14

$$Expression\ change\ = \frac{Wtf*(E2 - E1 + E2*Wtf)}{(Wtf + 1)*(Wtf + 2)} \quad (9)$$

16

17 Solving gives us $E2 = \frac{E1}{Wtf+1}$ for which expression does not change. This is the

18 upper limit of E2 for expression reduction, thus when E2 is below this limit we get

19 expression reduction. At low [TF] (low Wtf) this limit is high, thus at low [TF] it's easy to

20 reduce expression. This limit goes to E2 = 0 when [TF] goes to infinity, thus at high [TF]

21 there is no expression reduction. This is the behavior that we observe, suggesting that

52

1    TF sharing is the mechanism behind the expression reduction as a result of activator site

2    addition.

3

4          Taken together, we find that while both steric hindrance and TF sharing have

5    regimes in which site addition decreases expression, only TF sharing shows a decrease

6    when [TF] is low and less or no decrease when [TF] is high. In contrast, the steric

7    hindrance model gives the opposite behavior and shows an increase in expression

8    reduction as the [TF] goes up. Thus, a theoretic investigation of the effect of activator site

9    addition on expression shows that TF sharing and not steric hindrance is, at least

10    qualitatively, able to explain the expression reduction that we observe mostly at high [AA].

11          In order to investigate why adding a binding site can reduce expression we

12    performed the following analysis. We measure the expression driven by the promoter with

13    no Gcn4 binding sites (O), the two promoters each with a single Gcn4 binding site (

14    promoter A, promoter B ), and the promoter with both binding sites ( promoter AB ). We

15    observe that adding each binding sites separately results in an increase in expression

16    compared to the no binding site promoter (A>O and B>O), suggesting that these binding

17    sites, at least in isolation, do not recruit transcriptional repressors. Importantly, these two

18    binding sites in isolation drive different expression levels, with A driving higher expression

19    than B. However, we observe that the two Gcn4 binding sites together often results in

20    lower expression than the highest of the single promoters (AB < max(A,B) , and A>B)).

21    The 'sharing' model suggests that the B (lower expression) binding site steals TF from A

22    (the higher expression binding site), and therefore expression with AB is less than that

53

1    from A (but still greater than expression from B). These results cannot be explained by

2    additive repression.

3

4

5

6    ## Supplementary data

7    Supplementary data are online at

8    https://www.upf.edu/scb/SupplementaryData/vanDijk_Sharon_TFSharing_2015.zip

9

10   **expression_all_constructs.tab**

11          Measured expression values of all constructs measured in the 6 conditions

12   **annotation_all_constructs.tab**

13          Sequence annotations for all promoters in the library

14   **\*_constructs_expression.tab**

15          Predicted expression values for model \*, where \* is: null (basic model), PSE

16   (position specific expression), SH (steric hindrance), S (TF sharing), or a combination of

17   these.

18   **\*_params.tab**

19          Model parameters for model \*, where \* is: null (basic model), PSE (position

20   specific expression), SH (steric hindrance), S (TF sharing), or a combination of these.