



Plant_SNP_TATA_Z-tester: a Web Service that Unequivocally Estimates the Impact of Proximal Promoter Mutations on Plant Gene Expression

Dmitry Rasskazov¹, Irina Chadaeva¹, Ekaterina Sharypova¹, Karina Zolotareva¹, Bato Khandaev¹, Petr Ponomarenko¹, Nikolay Podkolodnyy^{1,2}, Natalya Tverdokhleba¹, Oleg Vishnevsky¹, Anton Bogomolov¹, Olga Podkolodnaya¹, Ludmila Savinkova¹, Elena Zemlyanskaya¹, Vadimir Golubyatnikov³, Nikolay Kolchanov¹, Mikhail Ponomarenko^{1,*}

¹ Institute of Cytology and Genetics, Novosibirsk 630090, Russia

² Institute of Computational Mathematics and Mathematical Geophysics, Novosibirsk 630090, Russia

³ Sobolev Institute of Mathematics, Novosibirsk 630090, Russia

* Correspondence: pon@bionet.nsc.ru. Tel.: +7 (383) 363-4963 ext. 1311 (M.P.)

Human_SNP_TATA_Z-tester as the prototype used in this work to empirically solve an ill-posed inverse problem about how to uniformly estimate effects of mutations in proximal promoters on gene expression in plant grown under various environmental conditions during development (Figure 2b)

For our previously developed Web-service Human_SNP_TATA_Z-tester [28] shown in Figure 2b (hereinafter, see the main text), the input is two 90 bp DNA sequences, which are $S_{wt} = \{s^{wt}_{-90} \dots s^{wt}_i \dots s^{wt}_{+1}\}$ and $S_m = \{s^{min}_{-90} \dots s^{min}_i \dots s^{min}_{+1}\}$ immediately before the transcription start site (TSS, $s^{wt}_0 = s^{min}_0$, where: $s_i \in \{a, c, g, t\}$) of the human gene proximal promoter carrying either ancestral (wt) or minor (min) allele, respectively, of an arbitrary SNP under study. In this figure, readers can see them in two textboxes “1st promoter” and “2nd promoter”, respectively.

After clicking the “Calculate” button (Figure 2b), using our three-step model [29] of the TBP-promoter binding (i.e., TBP slides along DNA [18] \leftrightarrow molecular co-recognition between TBP and TBP-site met [30] \leftrightarrow DNA-bend stabilizes TBP-promoter complex [32]) as proven experimentally [33], our Web-service Human_SNP_TATA_Z-tester [28] estimated two “ $-\ln[K_D(S_\bullet)]$ ”-values expressed on the natural-logarithm scale (ln-units), which evaluate the TBP-promoter affinity upon each DNA sequences (S_\bullet) independently from one another, such as:

$$-\ln[K_D(S_\bullet)] = 10.9 - 0.2 \{\ln[K_{SLIDE}(S_\bullet)K_{STOP}(S_\bullet)K_{BEND}(S_\bullet)]\}, \quad (S1)$$

where K_D is the equilibrium dissociation constant estimation (in moles per liter, M); 10.9 (ln-units) seems to numerically match nonspecific TBP–DNA affinity (10 μ M) as measured independently [62]; 0.2 is a stoichiometric coefficient of the three-step TBP–promoter binding, as determined elsewhere by means of the difference in the length of the TBP consensus site and the region of TBP sliding along DNA [29].

First of all, within Eq. S1, $-\ln[K_{STOP}(S_\bullet)]$ is an estimation of the equilibrium dissociation constant of the mutual recognition between TBP and the most probable TBP-site encountered at the second step among the three steps in question:

$$-\ln[K_{STOP}(S_\bullet)] = \text{MAX}_{-90 \leq i \leq -20; k \in \{-1, +1\}} \{\sum_{i-1 \leq j \leq i+13} w\{i, s_{jk}^\bullet\}\}, \quad (S2)$$

where $w\{i, s_{jk}^\bullet\}$ is Bucher’s weight of nucleotide s_{jk}^\bullet at the j th position of the TBP-site [31]; k is an indicator of either a direct (+1) or complementary (−1) strand of the double-stranded B-helical DNA of the promoter under study.; $\text{MAX}(\zeta)$ is the highest ζ -value observed.

Besides, in Eq. S1, $-\ln[K_{\text{SLIDE}}(S_{\bullet})]$ is an estimate of the equilibrium dissociation constant of an interaction between TBP and the promoter DNA during their sliding one over the other at the first step among the three within this bioinformatics model, as:

$$-\ln[K_{\text{SLIDE}}(S_{\bullet})] = \text{MEAN}_{[\xi-7;\xi+19]; k \in \{-1; +1\}} (35.1\mu + 0.8[\text{TA}]), \tag{S3}$$

where ξ is the position of the most probable TBP-site according to Bucher’s criterion [31] (i.e., Eq. S2); the μ value of the minor-groove width of the B-helical DNA at this site’s center was determined elsewhere [63]; [TA] is the concentration of dinucleotide TA; 0.8 and 35.1 are linear regression coefficients [64].

Finally, in Eq. S1, $-\ln[K_{\text{BEND}}(S_{\bullet})]$ is an estimation of the equilibrium dissociation constant of intermediate short-lived complexes between TBP and each of two DNA strands of the TBP-site separately from one another during DNA melting leading to the bend that fixes the TBP–promoter complex [32] at the last step of their binding, as follows:

$$-\ln[K_{\text{BEND}}(S_{\bullet})] = \text{MEAN}_{[\xi-7;\xi+19]; k \in \{-1; +1\}} (0.9[\text{TA}, \text{AA}, \text{TG}, \text{AG}] + 2.5[\text{TA}, \text{TC}, \text{TG}] + 14.4), \tag{S4}$$

where 0.9, 2.5, and 14.4 are linear regression coefficients [64].

After that, examining all the possible mutations, $s^* \rightarrow \varphi$, at each j th position among 26 positions of the most probable TBP-site according to Eq. S2, our Web-service Human_SNP_TATA_Z-tester [28] estimated standard error of the mean SEM_{\bullet} of the $-\ln[K_D(S_{\bullet})]$ values calculated using Eq. S1, as:

$$\text{SEM}_{\bullet} = \{(\sum_{\xi-7 \leq j \leq \xi+19} \sum_{\varphi \in \{a,c,g,t\}} \ln[K_D(s^*_{\xi-7} \dots s_{j-1} \varphi s_{j+1} \dots s^*_{\xi+19})] / K_D(s^*_{\xi-7} \dots s_{j-1} s_j s_{j+1} \dots s^*_{\xi+19})\}^2 / ((3*26)(3*26 - 1))\}^{1/2}. \tag{S5}$$

Using both sequences S_{wt} and S_{min} and Eqs. S1 – S5, this toolbox found two paired value sets $\{-\ln[K_D(S_{\text{wt}})] \pm \text{SEM}_{\text{wt}}\}$ and $\{-\ln[K_D(S_{\text{min}})] \pm \text{SEM}_{\text{min}}\}$, respectively, which are necessary for Fisher’s Z-score [65], for instance:

$$Z = \text{abs}\{\ln[K_D(S_{\text{wt}}) / K_D(S_{\text{min}})] / [\text{SEM}_{\text{wt}}^2 + \text{SEM}_{\text{min}}^2]^{1/2}\}. \tag{S6}$$

Eventually, with the help of the R software [65], using this Z-value, our Web-service Human_SNP_TATA_Z-tester [28] found a p value of the probability of the tested hypothesis “ $H_0: K_D(S_{\text{wt}}) \neq K_D(S_{\text{min}})$ ” so that if it is statistical significant ($p > 0.95$), it made the decision:

```
IF      {INEQUALITY “ $K_D(S_{\text{min}}) < K_D(S_{\text{wt}})$ ” is statistically significant },
THEN   {PREDICTION is “the minor allele of the gene considered is overexpressed relative to the ancestral one”};
ELSE   IF      {INEQUALITY “ $K_D(S_{\text{min}}) > K_D(S_{\text{wt}})$ ” is statistically significant},
THEN   {DECISION is “the minor allele of this gene is underexpressed relative to the ancestral one”},
OTHERWISE {DECISION is “the expression change of this gene is insignificant”}.
```

Our Web service Human_SNP_TATA_Z-tester [28] presents this decision (Eq. S7) in the “Decision” line of the “Result” textbox, while all the intermediate results are in the other lines of this textbox, as readers can see in Figure 2b. Table S1 presents how we experimentally selectively confirmed this from article to article [28, 35-39, 41, 42, 66, 67]. Figure S1 does this in graphical form.

Table S1. An electrophoretic mobility shift assay (EMSA)-based in vitro verification of the complex of TBP with synthetic 26 bp oligodeoxyribonucleotides (ODNs) identical to natural human gene promoters near each SNP tested

#	Human gene, dbSNP ID [66]	WT min	26 bp oligodeoxyribonucleotides, direct strand 5'-ODN-3'	Prediction, ln units		Experiment, <i>in vitro</i> , ln units		Reference
				$-\ln(K_D)$	$\Delta \ln(K_D)$	$-\ln(K_D)$	$\Delta \ln(K_D)$	
1	LEP	WT	gatcggggccGCTATAAGAgggggcggg	19.43		16.37		
2	rs34104384	A-30T	gatcggggccGCTATAAGTgggggcggg	19.70	0.27	16.43	0.06	[35]
3	rs201381696	A-35G	gatcggggccGCTGTAAGAgggggcggg	18.22	-1.21	15.28	-1.09	
4	rs200487063	G-38A	gatcggggccACTATAAGAgggggcggg	19.85	0.42	16.96	0.59	
5	ABCA9	WT	aattattttgTATATTTctgagcatac	19.66		16.47		[39]
6	rs367781716	T-37C	aattattttgCATATTTctgagcatac	18.78	-0.88	15.59	-0.88	
7	F9	WT	tttggTACAACTAATcgaccttacca	18.86		15.05		[41]
8	rs750827465	C-34A	tttggTACAAATAATcgaccttacca	19.32	0.46	16.22	1.17	

Table S1. Cont.

#	Human gene, dbSNP ID [66]	WT min	26 bp oligodeoxyribonucleotides, direct strand 5'-ODN-3'	Prediction, ln units		Experiment, <i>in vitro</i> , ln units		Reference
				-ln(K _D)	Δln(K _D)	-ln(K _D)	Δln(K _D)	
9	<i>HBD</i>	WT	acaggaccagC A TAAAAggcagggca	19.29		17.14		[36]
10	rs34166473	T-30C	acaggaccagCA T AAAAggcagggca	18.27	-1.02	15.02	-2.12	
11	rs35518301	A-31G	acaggaccagC G TAAAAggcagggca	18.65	-0.64	16.12	-1.02	
12	<i>MBL2</i>	WT	catctatttctTA T ATAGcctgcaccc	20.17		17.20		
13	rs72661131	T-39C	catctatttctTA C ATAGcctgcaccc	19.30	-0.87	16.73	-0.47	
14	<i>IL1B</i>	WT	ttttgaaagc C ATAAAAAacagcgagg	19.22		17.73		
15	rs1143627	C-31T	ttttgaaagc T ATAAAAAacagcgagg	20.16	0.94	19.15	1.42	
16	<i>TPI1</i>	WT	cgcggcgctcTATA T AAgtgggcagt	20.70		19.31		
17	rs1800202	T-26G	cgcggcgctcTATA G AAgtgggcagt	19.29	-1.41	16.12	-3.13	
18	<i>F3</i>	WT	gccggcccTTTATAg c gcgcggggca	19.60		16.45		
19	rs563763767	c-21T	gccggcccTTTATAg T gcgcggggca	20.02	0.42	17.47	1.02	
20	<i>HBB</i>	WT	cagggctggg C A T AAAAgtcagggca	19.20		16.81		[42]
21	rs34598529	A-28G	cagggctgggCATAg G AAgtcagggca	18.34	-0.86	14.40	-2.41	
22	rs34598529	A-28C	cagggctgggCATAg C AAgtcagggca	18.63	-0.57	14.51	-2.30	
23	rs281864525	A-25C	cagggctgggCATAAAA C gtcagggca	18.73	-0.47	15.71	-1.10	
24	rs63750953	Δ-25AA	ccagggctgggCATAg A gtcagggcag	18.61	-0.59	15.71	-1.10	
25	rs33980857	T-29A	cagggctgggCA A AAAAgtcagggca	17.70	-1.50	16.02	-0.79	
26	rs33980857	T-29C	cagggctgggCA C AAAAgtcagggca	18.17	-1.03	15.78	-1.03	
27	rs33980857	T-29G	cagggctgggCA G AAAAgtcagggca	17.67	-1.53	16.02	-0.79	
28	rs33931746	A-27T	cagggctgggCATAg T AAgtcagggca	19.75	0.55	16.63	-0.18	
29	rs34598529	A-28G	cagggctgggCAT G AAgtcagggca	17.85	-1.35	14.51	-2.30	
30	rs34500389	C-32T	cagggctggg T ATAAAAAgtcagggca	20.18	0.98	17.26	0.45	
31	<i>HBZ</i>	WT	agctccctgTA T ATAAggggaccctg	20.76		18.93		
32	rs11318094	T-29A	agctccctgTA A ATAAggggaccctg	19.28	-1.48	17.03	-1.90	
33	<i>EPOR</i>	WT	cacgtcatc T A T TTTGTctgctacg	18.24		15.05		
34	rs1006576690	T-27A	cacgtcatcTAT A TTTGTctgctacg	19.39	1.15	18.33	2.74	[37]
35	rs971717705	A-29G	cacgtcatcT G TTTTTGTctgctacg	18.11	-0.13	15.71	0.12	
36	rs567946217	c-31A	cacgtcat A TATTTTGTctgctacg	19.33	1.09	17.68	2.09	
37	rs567946217	c-31G	cacgtcat G TATTTTGTctgctacg	18.51	0.27	16.34	0.75	
38	<i>GCG</i>	WT	gctggagagT A TATAAAAgcagtgcg	20.89		18.64		[38]
39	rs183433761	A-41G	gctggagagT G TATAAAAgcagtgcg	20.28	-0.61	17.83	-0.81	

Table S1. Cont.

#	Human gene, dbSNP ID [66]	WT min	26 bp oligodeoxyribonucleotides, direct strand 5'-ODN-3'	Prediction, ln units		Experiment, <i>in vitro</i> , ln units		Reference
				-ln(K _D)	Δln(K _D)	-ln(K _D)	Δln(K _D)	
40	ASMT	WT	ggtgaccttttgtGcccagaataggt	18.18		14.33		
41	rs1402972626	G-30A	ggtgaccttttgtAcccagaataggt	18.93	0.75	13.82	-0.51	
42	CDY2A	WT	agaatgttccataTaatcgtcatagc	19.27		15.65		
43	rs20067072	T-24C	agaatgttccataCaatcgtcatagc	18.76	-0.51	14.51	-1.14	
44	GTPBP6	WT	atcacgagcacgtGatgaggagcggc	17.30		13.41		[28]
45	rs1393008234	G-24T	atcacgagcacgtTatgaggagcggc	18.68	1.38	13.48	0.07	
46	SHOX	WT	gaggtcgccgcgtAataatagtgaga	20.31		17.06		
47	rs1452787381	A-45G	gaggtcgccgcgtGataatagtgaga	19.21	-1.10	15.16	-1.90	
48	ZFY	WT	ggcggagggggccCaactaccatccc	17.67		13.82		
49	rs1452787381	C-56T	ggcggagggggccTaactaccatccc	18.18	0.51	13.12	-0.70	
50	GRIN1	WT	tggagggggACAAAGACagggtggtg	17.59		14.95		
51	rs1402667001	g-34a	tggaggaggACAAAGACagggtggtg	17.74	0.15	15.42	0.47	
52	ASCL3	WT	tcgaaaaaTAAAAAATAAataaacat	19.04		18.24		[67]
53	rs1049743008	T-45C	tcgaaaaaTAAAAACAAAataaacat	18.75	-0.29	17.90	-0.34	
54	NOS1	WT	tgtttcctGATAGAAAAaaaaaatgg	18.56		18.77		
55	rs1195040887	G-27A	tgtttcctGATAAAAAaaaaaatgg	18.96	0.34	18.95	0.18	

Note. For each TBP-ODN complex, $-\ln(K_D)$ and $\Delta\ln(K_D) = \ln(K_{D;WT}/K_{D;min})$ are absolute and relative estimates (i.e., compared to those of the wild-type allele, WT), respectively, of the equilibrium dissociation constant expressed in natural-logarithm units (ln units). Human genes: *ABCA9*, ATP-binding cassette subfamily A member 9; *ASCL3*, Achaete-Scute family BHLH transcription factor 3; *ASMT*, acetylserotonin O-methyltransferase; *CDY2A*, chromodomain Y-linked 2A; *EPOR*, erythropoietin receptor; *F3* and *F9*, coagulation factors III and IX, respectively; *GCG*, glucagon; *GRIN1*, glutamate ionotropic receptor NMDA type subunit 1; *GTPBP6*, GTP-binding protein 6; *HBB*, *HBD*, and *HBZ*, hemoglobin subunits β , δ , and ζ , respectively; *IL1B*, interleukin 1 β ; *LEP*, leptin; *MBL2*, mannose-binding lectin 2; *NOS1*, nitric oxide synthase 1; *SHOX*, short stature homeobox; *TPI1*, triosephosphate isomerase 1; *ZFY*, zinc finger protein Y-linked.

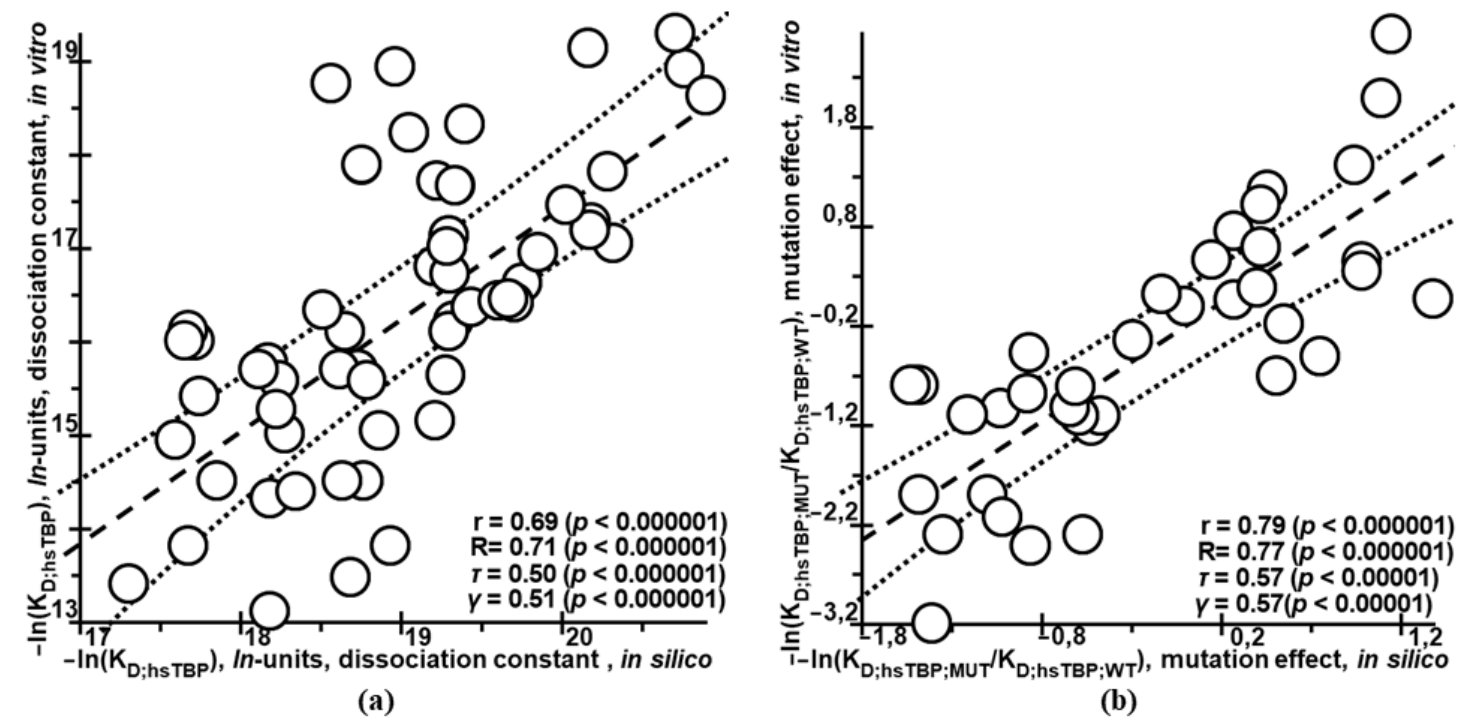


Figure S1. The significant correlations between the *in silico*-predicted (X-axis) and *in vitro*-measured (Y-axis) K_D values of the equilibrium dissociation constant of the TBP-ODN complex, as graphical representation of Table S1. Legend: (a) and (b): absolute and relative estimates (i.e., compared to those of the wild-type allele, WT), respectively, of equilibrium dissociation constant K_D expressed in natural-logarithm units (ln units). Dashed and dotted lines denote linear regression and boundaries of its 95% confidence interval, as calculated in the Statistica software (Statsoft™, Tulsa, OK, USA). Circles denote the ancestral (WT) and minor alleles (dbSNP ID [66]) of the SNP listed in Table S1; r , R , τ , γ , and p are linear correlation, Spearman's rank correlation, Kendall's rank correlation, Goodman-Kruskal generalized correlation, and their statistical significance levels, respectively.

