# Science Advances

## Supplementary Materials for

### Measuring quantitative effects of methylation on transcription factor–DNA binding affinity

Zheng Zuo, Basab Roy, Yiming Kenny Chang, David Granas, Gary D. Stormo

**This PDF file includes:**

- fig. S1. General illustration of the use of M and W nomenclature to represent methylated bases in a DNA sequence.
- fig. S2. Methyl-Spec-seq ePWMs.
- fig. S3. Replicates of FAM and TAMRA anisotropy signals that were used to calculate the effect of mC on the relative binding specificity of ZFP57.
- fig. S4. The relative binding energies of all 64 variants (AP1 libraries) with different methylation profiles, ranked from the strongest (lowest energy) to the weakest binder of the unmethylated library.
- fig. S5. Replicate experiments with HOXB13.
- fig. S6. EMSA sample images for mouse ZFP57 (F1 to F3) and CTCF (F1 to F9).
- fig. S7. EMSA sample images for Gli1, JunB/BATF, and HOXB13.
- fig. S8. Schematic maps of plasmids used for cloning and expression of proteins.
- text S1. DNA oligo sequences for primers and libraries.
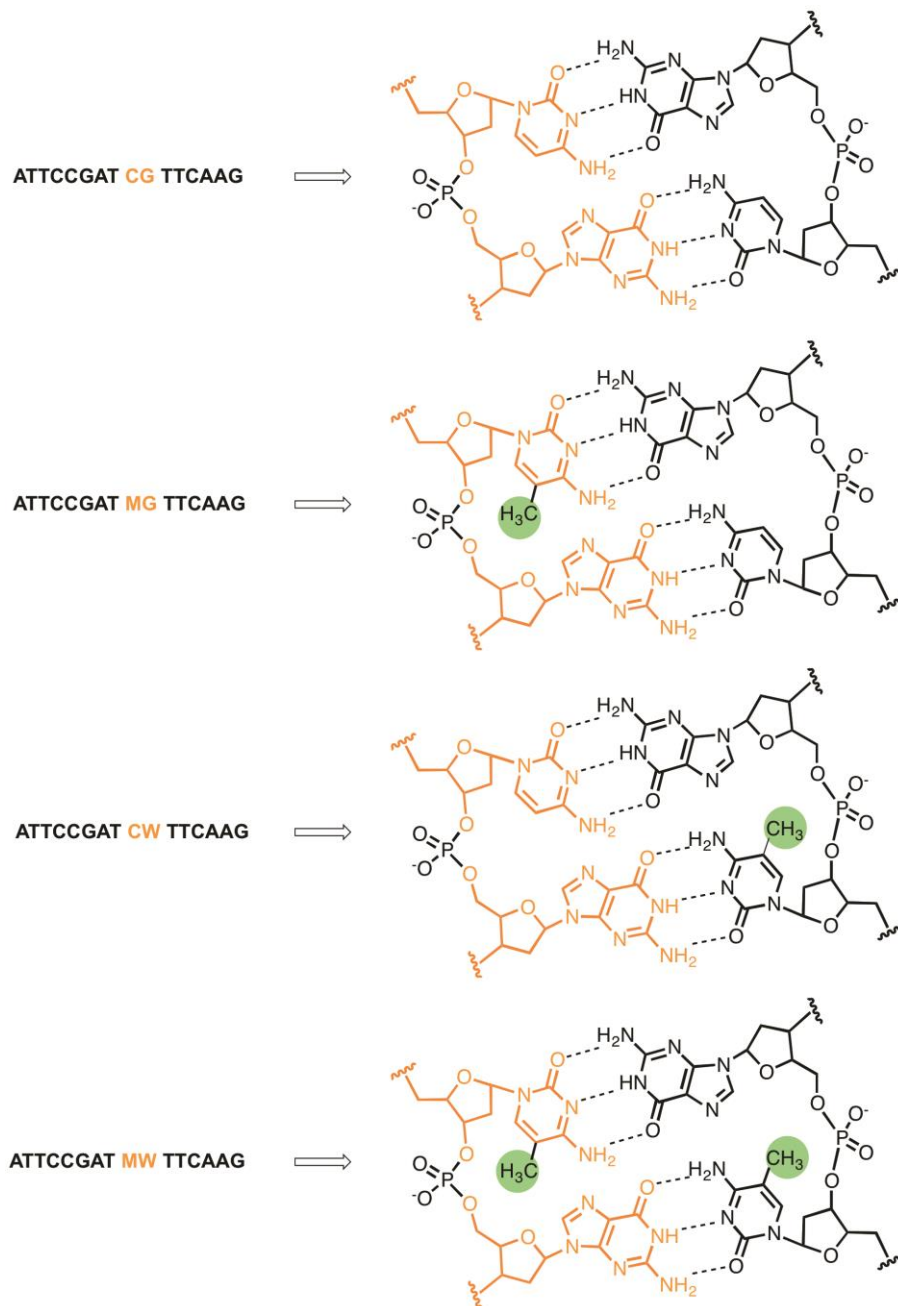- text S2. Instructions for software use.

**fig. S1. General illustration of the use of M and W nomenclature to represent methylated bases in a DNA sequence.** The CG dinucleotide, depending on the methylation status of C, can be represented by one of the four different combinations, such as CG, MG, CW and MW. M represents 5′-methylcytidine, whereas W is 5′-methylcytidine opposing to a G on the reference strand.
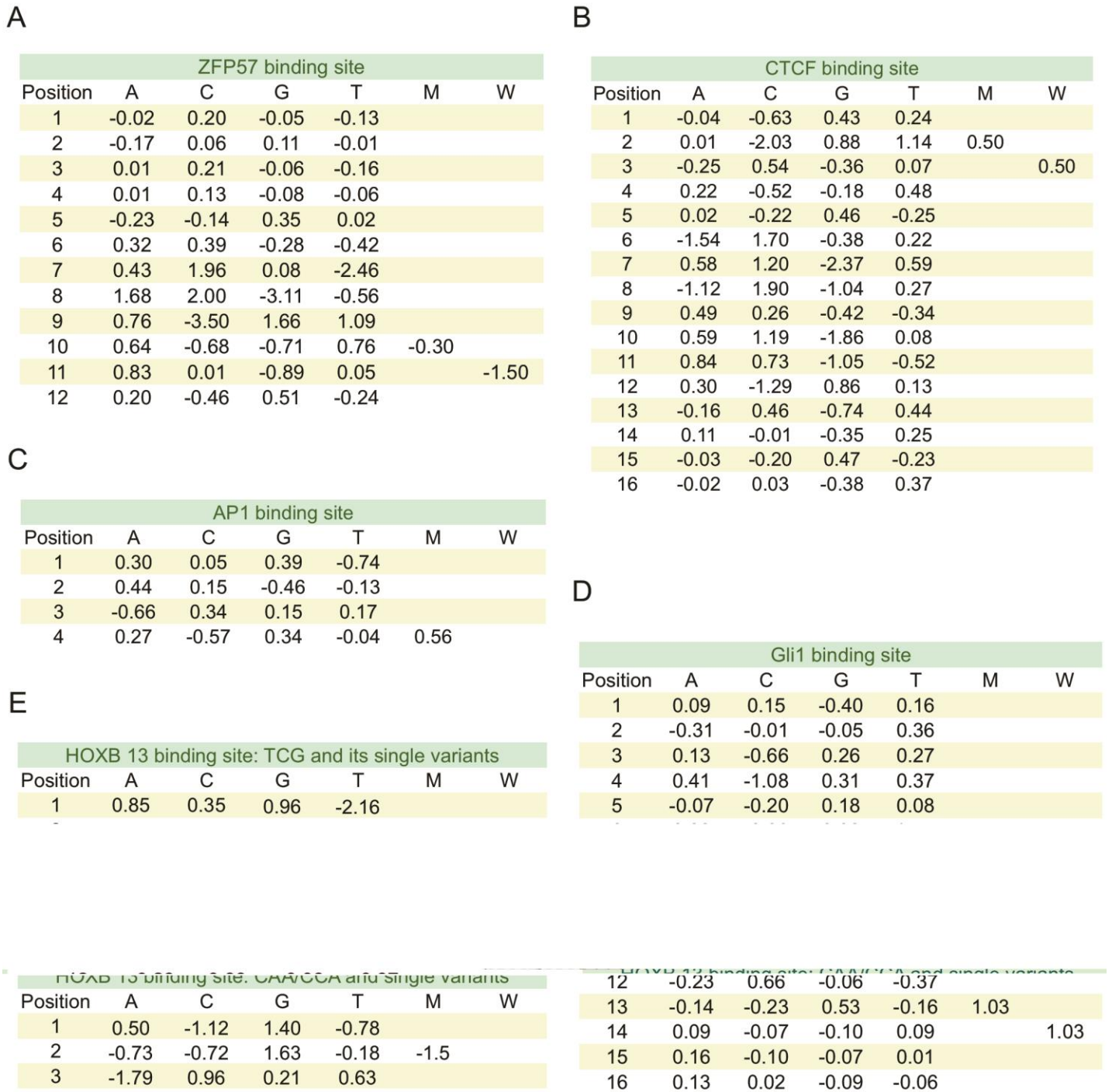
**A**

| ZFP57 binding site | | | | | | |
|---|---|---|---|---|---|---|
| Position | A | C | G | T | M | W |
| 1 | -0.02 | 0.20 | -0.05 | -0.13 | | |
| 2 | -0.17 | 0.06 | 0.11 | -0.01 | | |
| 3 | 0.01 | 0.21 | -0.06 | -0.16 | | |
| 4 | 0.01 | 0.13 | -0.08 | -0.06 | | |
| 5 | -0.23 | -0.14 | 0.35 | 0.02 | | |
| 6 | 0.32 | 0.39 | -0.28 | -0.42 | | |
| 7 | 0.43 | 1.96 | 0.08 | -2.46 | | |
| 8 | 1.68 | 2.00 | -3.11 | -0.56 | | |
| 9 | 0.76 | -3.50 | 1.66 | 1.09 | | |
| 10 | 0.64 | -0.68 | -0.71 | 0.76 | -0.30 | |
| 11 | 0.83 | 0.01 | -0.89 | 0.05 | | -1.50 |
| 12 | 0.20 | -0.46 | 0.51 | -0.24 | | |

**B**

| CTCF binding site | | | | | | |
|---|---|---|---|---|---|---|
| Position | A | C | G | T | M | W |
| 1 | -0.04 | -0.63 | 0.43 | 0.24 | | |
| 2 | 0.01 | -2.03 | 0.88 | 1.14 | 0.50 | |
| 3 | -0.25 | 0.54 | -0.36 | 0.07 | | 0.50 |
| 4 | 0.22 | -0.52 | -0.18 | 0.48 | | |
| 5 | 0.02 | -0.22 | 0.46 | -0.25 | | |
| 6 | -1.54 | 1.70 | -0.38 | 0.22 | | |
| 7 | 0.58 | 1.20 | -2.37 | 0.59 | | |
| 8 | -1.12 | 1.90 | -1.04 | 0.27 | | |
| 9 | 0.49 | 0.26 | -0.42 | -0.34 | | |
| 10 | 0.59 | 1.19 | -1.86 | 0.08 | | |
| 11 | 0.84 | 0.73 | -1.05 | -0.52 | | |
| 12 | 0.30 | -1.29 | 0.86 | 0.13 | | |
| 13 | -0.16 | 0.46 | -0.74 | 0.44 | | |
| 14 | 0.11 | -0.01 | -0.35 | 0.25 | | |
| 15 | -0.03 | -0.20 | 0.47 | -0.23 | | |
| 16 | -0.02 | 0.03 | -0.38 | 0.37 | | |

**C**

| AP1 binding site | | | | | | |
|---|---|---|---|---|---|---|
| Position | A | C | G | T | M | W |
| 1 | 0.30 | 0.05 | 0.39 | -0.74 | | |
| 2 | 0.44 | 0.15 | -0.46 | -0.13 | | |
| 3 | -0.66 | 0.34 | 0.15 | 0.17 | | |
| 4 | 0.27 | -0.57 | 0.34 | -0.04 | 0.56 | |

**D**

| Gli1 binding site | | | | | | |
|---|---|---|---|---|---|---|
| Position | A | C | G | T | M | W |
| 1 | 0.09 | 0.15 | -0.40 | 0.16 | | |
| 2 | -0.31 | -0.01 | -0.05 | 0.36 | | |
| 3 | 0.13 | -0.66 | 0.26 | 0.27 | | |
| 4 | 0.41 | -1.08 | 0.31 | 0.37 | | |
| 5 | -0.07 | -0.20 | 0.18 | 0.08 | | |
| 12 | -0.23 | 0.66 | -0.06 | -0.37 | | |
| 13 | -0.14 | -0.23 | 0.53 | -0.16 | 1.03 | |
| 14 | 0.09 | -0.07 | -0.10 | 0.09 | | 1.03 |
| 15 | 0.16 | -0.10 | -0.07 | 0.01 | | |
| 16 | 0.13 | 0.02 | -0.09 | -0.06 | | |

**E**

| HOXB 13 binding site: TCG and its single variants | | | | | | |
|---|---|---|---|---|---|---|
| Position | A | C | G | T | M | W |
| 1 | 0.85 | 0.35 | 0.96 | -2.16 | | |

| HOXB 13 binding site: CAA/CCA and single variants | | | | | | |
|---|---|---|---|---|---|---|
| Position | A | C | G | T | M | W |
| 1 | 0.50 | -1.12 | 1.40 | -0.78 | | |
| 2 | -0.73 | -0.72 | 1.63 | -0.18 | -1.5 | |
| 3 | -1.79 | 0.96 | 0.21 | 0.63 | | |

**fig. S2. Methyl-Spec-seq ePWMs.** (**A**) ZFP57. (**B**) CTCF. (**C**) JUNB-BATF1 AP1 sites. (**D**) GLI1. (**E**) HOXB13 TCG motif. (**F**) HOXB13 CAA/CCA motif. All values are in kT energy units and the position average (not including M and W) is set to 0. For M and W, energy values are included only if they exceed 0.2 in absolute value, that being the typical variance in measurements.

ZFP57 concentration →

**ME-TAMRA vs. ME-FAM**

| | | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FAM | Ex-1 | 18.68 | 26.54 | 33.89 | 36.31 | 44.49 | 49.81 | 57.21 | 63.84 | 71.73 | 85.11 | 98.41 |
| | Ex-2 | 20.47 | 21.42 | 25.43 | 25.93 | 28.41 | 32.04 | 32.43 | 34.57 | 39.89 | 41.31 | 44.84 |
| | Ex-3 | 19.15 | 23.61 | 24.51 | 26.77 | 30.25 | 34.55 | 36.52 | 39.11 | 46.42 | 53.47 | 57.30 |
| | Ex-4 | 31.91 | 31.62 | 33.42 | 39.58 | 39.56 | 43.39 | | | | | |
| TAMRA | Ex-1 | 102.69 | 108.66 | 114.60 | 115.19 | 125.72 | 130.00 | 136.28 | 143.02 | 148.73 | 158.82 | 168.22 |
| | Ex-2 | 99.90 | 101.65 | 104.97 | 106.67 | 108.52 | 110.87 | 111.93 | 112.08 | 117.35 | 118.75 | 120.99 |
| | Ex-3 | 98.81 | 101.79 | 102.33 | 105.60 | 109.52 | 111.45 | 112.51 | 115.16 | 120.73 | 124.46 | 129.67 |
| | Ex-4 | 112.04 | 110.71 | 112.22 | 118.20 | 116.11 | 120.50 | | | | | |

**ME-TAMRA vs. HM-top-FAM**

| | | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FAM | Ex-1 | 24.69 | 26.09 | 30.28 | 38.04 | 41.50 | 47.81 | 54.25 | 56.59 | 58.44 | 63.18 | 77.26 |
| | Ex-2 | 23.91 | 25.66 | 26.44 | 28.91 | 30.08 | 30.28 | 30.44 | 31.45 | 31.51 | 31.73 | 34.31 |
| | Ex-3 | 30.84 | 30.65 | 29.08 | 29.02 | 30.07 | 32.74 | | | | | |
| TAMRA | Ex-1 | 109.55 | 116.60 | 129.51 | 144.76 | 156.64 | 166.00 | 175.66 | 178.26 | 180.26 | 185.65 | 194.95 |
| | Ex-2 | 108.88 | 113.41 | 115.09 | 121.04 | 123.56 | 127.07 | 129.51 | 129.60 | 129.84 | 133.22 | 137.83 |
| | Ex-3 | 124.61 | 126.49 | 127.28 | 126.80 | 130.24 | 132.88 | | | | | |

**ME-TAMRA vs. HM-bottom-FAM**

| | | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FAM | Ex-1 | 21.90 | 26.08 | 31.91 | 36.97 | 48.75 | 51.47 | 62.83 | 72.60 | 82.71 | 95.24 | 103.84 |
| | Ex-2 | 22.86 | 24.67 | 24.71 | 25.29 | 29.59 | 31.47 | 35.85 | 38.39 | 42.36 | 44.31 | 44.71 |
| | Ex-3 | 30.96 | 34.46 | 38.62 | 41.19 | 43.53 | 43.96 | | | | | |
| TAMRA | Ex-1 | 107.33 | 113.53 | 121.00 | 129.96 | 142.53 | 147.12 | 160.21 | 167.44 | 177.05 | 184.34 | 188.95 |
| | Ex-2 | 108.70 | 109.26 | 111.46 | 112.07 | 117.98 | 118.96 | 123.99 | 125.66 | 129.59 | 131.98 | 132.08 |
| | Ex-3 | 117.68 | 124.17 | 125.96 | 128.44 | 130.82 | 129.78 | | | | | |

**ME-TAMRA vs. UN-FAM**

| | | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FAM | Ex-1 | 25.00 | 26.58 | 29.05 | 36.47 | 39.62 | 44.41 | 46.99 | 51.93 | 56.87 | 61.95 | 71.17 |
| | Ex-2 | 28.00 | 28.87 | 28.93 | 29.87 | 31.44 | 31.91 | 31.96 | 32.85 | | | |
| | Ex-3 | 22.92 | 25.42 | 27.41 | 27.15 | 27.84 | 30.57 | 31.59 | 32.53 | | | |
| TAMRA | Ex-1 | 108.44 | 118.84 | 131.83 | 146.44 | 161.73 | 171.95 | 178.54 | 185.69 | 193.85 | 194.19 | 200.77 |
| | Ex-2 | 114.70 | 119.62 | 122.55 | 125.71 | 135.92 | 140.32 | 144.95 | 145.67 | | | |
| | Ex-3 | 104.25 | 110.21 | 116.10 | 117.01 | 121.82 | 126.90 | 131.06 | 134.44 | | | |

**fig. S3. Replicates of FAM and TAMRA anisotropy signals that were used to calculate the effect of mC on the relative binding specificity of ZFP57.** The DNAs used in this study are listed in Fig. 2B and the FAM (horizontal axis) Vs. TAMRA (vertical axis) plot is shown in Fig. 2C.
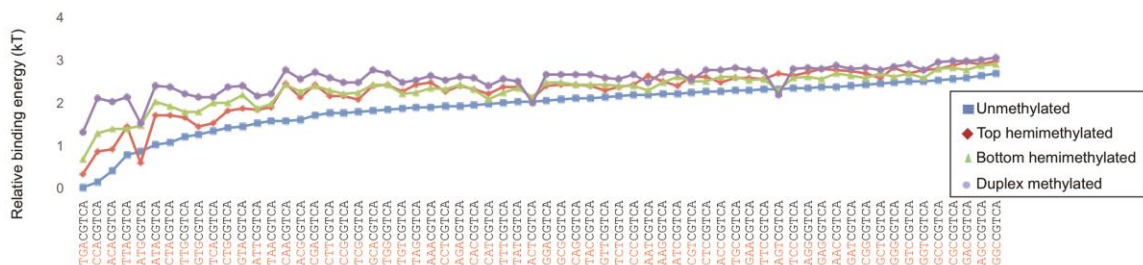


**fig. S4. The relative binding energies of all 64 variants (AP1 libraries) with different methylation profiles, ranked from the strongest (lowest energy) to the weakest binder of the unmethylated library.** The energy of the strongest binder is set to zero.
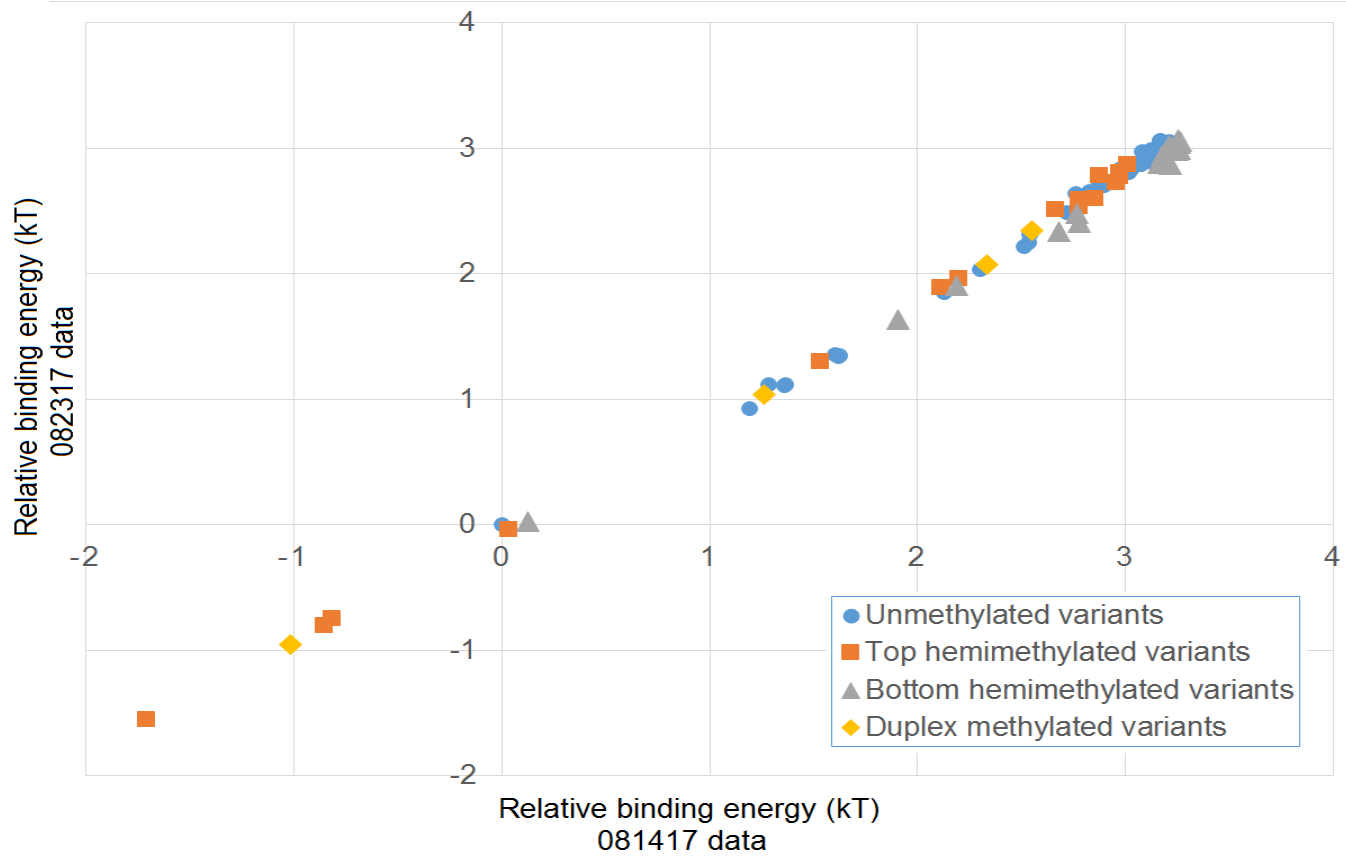
**fig. S5. Replicate experiments with HOXB13.** Shown are the measured binding energies for replicate experiments with the reference sequence, unmethylated TCG, set to 0 energy in each data set. The correlation has $r^2=0.99$ with a slope of 0.95.
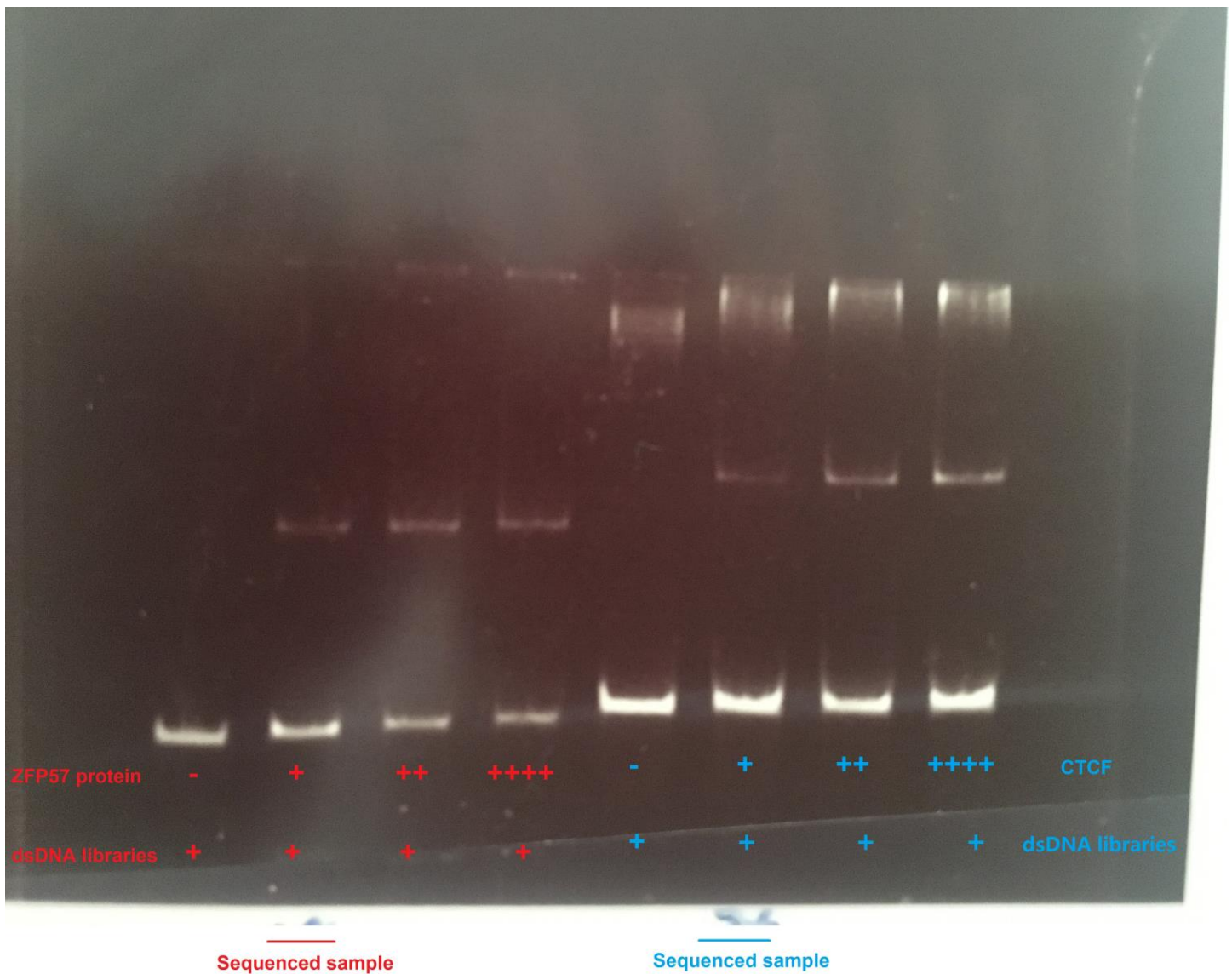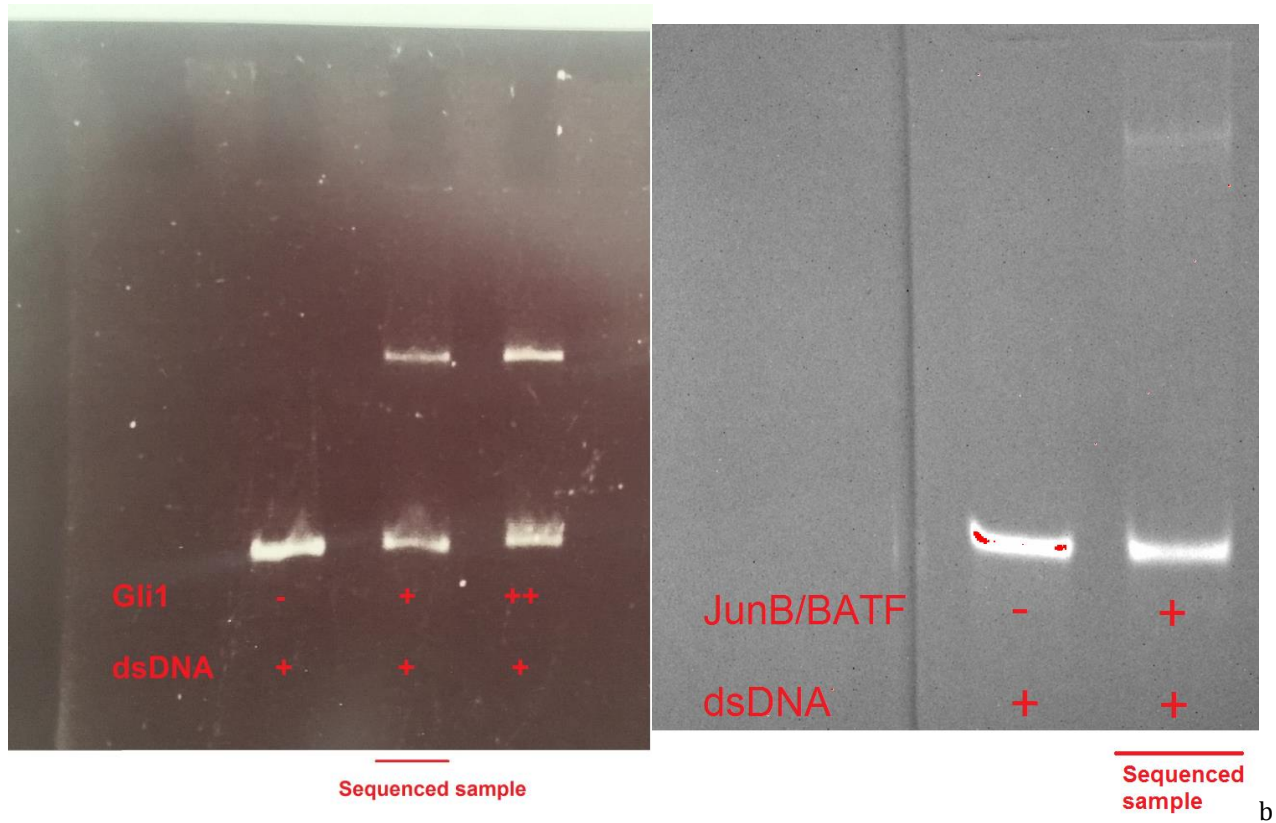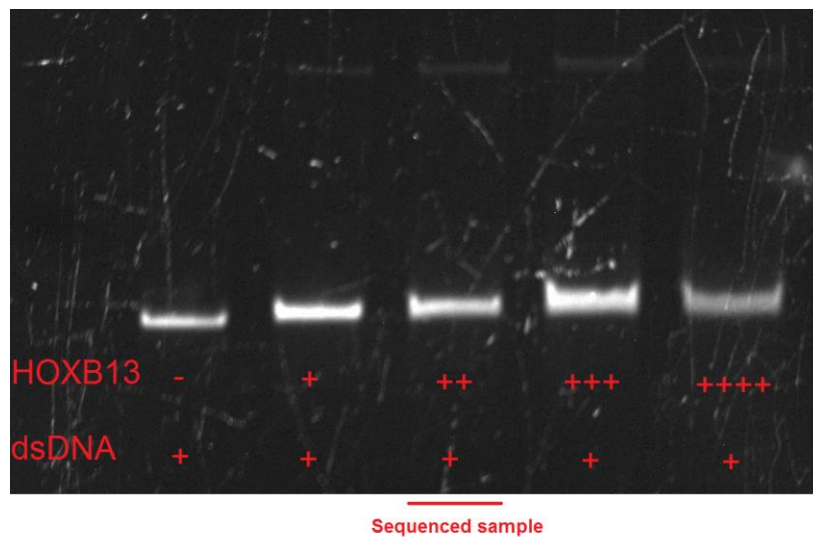
**fig. S6. EMSA sample images for mouse ZFP57 (F1 to F3) and CTCF (F1 to F9).** For ZFP57, 200ng dsDNA libraries were used for each lane; for CTCF, ~400ng dsDNA were titrated into each sample, some abnormal DNA bands at high MW range were excluded for gel cutting and sequencing. In both cases, the highest protein concentration used for binding reactions were estimated to be 1mM. All dsDNA were labeled by FAM on 5' end and imaged for up to 6s exposure. Both samples were run in 9% Tris-glycine gel for 30mins.

(a)



(b)



(c)

**fig. S7. EMSA sample images for Gli1, JunB/BATF, and HOXB13.** The highest protein concentration used for each binding reactions were estimated to be 1mM. All dsDNA were labeled by FAM on 5′-end for gel imaging. (**a**) Gli1 samples, similar condition as CTCF. (**b**) 200 ng dsDNA of AP1 libraries was used for each lane, in the absence and presence of JunB/BATF heterodimer, respectively.(**c**) For HOXB13(mouse), its DNA binding domain (DBD) was expressed by cell-free NEB PURExpress system, instead of E. coli BL21 setting. In vitro synthesized hisSUMO-HOXB13 protein was titrated into individual 20uL binding reaction from low to high (0, 1, 2, 3, 4uL each).

(a)

(b)

**fig. S8. Schematic maps of plasmids used for cloning and expression of proteins.** Schematic vector maps (**a**) NEB DHFR control plasmid was chosen as the original vector backbone, harboring different coding sequences including ZFP57, CTCF, BATF, Gli1, and HOXB13. (**b**) JunB gene was cloned separately into a Kanamycin-resistant, low copy plasmid for co-transformation with BATF construct.

# Supplementary text S1: DNA oligo sequences for primers and libraries.

## General DNA oligo sequences used in Methyl-Spec-seq experiment

**PE1:**
5'-AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT-3'

**PE1-Genetics:**
5'-AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT GATAGTCTCATTTTCACC-3'

**PE1-N-Genetics:**
5'-AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT N GATAGTCTCATTTTCACC-3'

**PE1-TNT-Genetics:**
5'-AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT TNT GATAGTCTCATTTTCACC-3'

**iPE2-42-Physics:**
5'-CAAGCAGAAGACGGCATACGAGATGCTACGCCAGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT AGAACAGATACTGTAATGGAA-3'

**iPE2-43-Physics:**
5'-CAAGCAGAAGACGGCATACGAGATCGTGCAGTCGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT AGAACAGATACTGTAATGGAA-3'

**iPE2-44-Physics:**
5'-CAAGCAGAAGACGGCATACGAGATCGAAATTCTGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT AGAACAGATACTGTAATGGAA-3'

**iPE2-45-Physics:**
5'-CAAGCAGAAGACGGCATACGAGATGCAATCGTCGTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT AGAACAGATACTGTAATGGAA-3'

**Physics-FAM:**
5'-FAM-AGAACAGATACTGTAATGGAA-3'

Note: For each individual sample isolated from EMSA gel, bound or unbound, it was extracted, purified, and further PCR amplified by a particular pair of primers including the PE1 and indexed PE2 ends. In our case, there are three distinct PE1 primers and four indexed PE2 primers, thus totaling 12 different combinations at most for single Illumina pair end sequencing

## Library-specific template and primer pairing for each dsDNA template, including unmethylated, enzymatically methylated, and chemically synthesized methylation libraries.

Note: Barcoding regions to differentiate each methylation state was underlined and labeled green in each template/primer pair, whereas chemically synthesized methylated cytosine was labeled bold red.

**ZFP57-R1-duplex methylated:**
```
5'-GATAGTCTCATTTTCACC CATNNNCAGTGCMGC TTCCATTACAGTATCTGT-3'
                               3'-GMG AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R1-M.sssI methylated:**
```
5'-GATAGTCTCATTTTCACC TAGNNNCAGTGCCGC TTCCATTACAGTATCTGT-3'
                                    3'-AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R2-bottom hemimethylated:**
```
5'-GATAGTCTCATTTTCACC CTAATCNNNTGCCGC TTCCATTACAGTATCTGT-3'
                       3'-GMG AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R2-duplex methylated:**
```
5'-GATAGTCTCATTTTCACC CATATCNNNTGCMGC TTCCATTACAGTATCTGT-3'
                       3'-GMG AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R2-M.sssI methylated:**
```
5'-GATAGTCTCATTTTCACC TAGATCNNNTGCCGC TTCCATTACAGTATCTGT-3'
                                    3'-AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R2-top hemimethylated:**
```
5'-GATAGTCTCATTTTCACC ATCATCNNNTGCMGC TTCCATTACAGTATCTGT-3'
                                   3'AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R2-unmethylated:**
```
5'-GATAGTCTCATTTTCACC AGAATCNNNTGCCGC TTCCATTACAGTATCTGT-3'
                                   3'AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R3-M.sssI methylated:**
```
5'-GATAGTCTCATTTTCACC TAGATCCAGNNNCGC TTCCATTACAGTATCTGT-3'
                                    3'-AAGGTAATGTCATAGACA-5'FAM'
```

**ZFP57-R4-M.sssI methylated:**
```
5'-GATAGTCTCATTTTCACC TAGATCCAGTGCNNN TTCCATTACAGTATCTGT-3'
                                    3'-AAGGTAATGTCATAGACA-5'FAM'
```


**CTCF-R1-M.sssI:**
```
5'-GATAGTCTCATTTTCACC ATNNNNTAGGGGGCACTATGT TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R1-Un:**
```
5'-GATAGTCTCATTTTCACC ATNNNNTAGGGGGCACTATGA TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R2-M.sssI:**
```
5'-GATAGTCTCATTTTCACC ATCCANNNNGGGGCACTATGT TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R2-Un:**
```
5'-GATAGTCTCATTTTCACC ATCCANNNNGGGGCACTATGA TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R3-M.sssI:**
```
5'-GATAGTCTCATTTTCACC ATCCACTANNNNGCACTATGT TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R3-Un:**
```
5'-GATAGTCTCATTTTCACC ATCCACTANNNNGCACTATGA TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R4-M.sssI:**
```
5'-GATAGTCTCATTTTCACC ATCCACTAGGGNNNNCTATGT TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R4-Un:**
```
5'-GATAGTCTCATTTTCACC ATCCACTAGGGNNNNCTATGA TTCCATTACAGTATCTGT-3'
                                          3'-AAGGTAATGTCATAGACA-5'FAM'
```

**CTCF-R5-M.sssI:**

5'-GATAGTCTCATTTTCACC ATCCACTAGGGGGCNNNNTG**T** TTCCATTACAGTATCTGT-3'
                                      3'-AAGGTAATGTCATAGACA-5'FAM'

**CTCF-R5-Un:**

5'-GATAGTCTCATTTTCACC ATCCACTAGGGGGCNNNNTG**A** TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand1-M.sssI:**

5'-GATAGTCTCATTTTCACC **AT**GNNNNACCCAAGATGAA TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'

**Gli1-Rand1-unmethylated:**

5'-GATAGTCTCATTTTCACC **TA**GNNNNACCCAAGATGAA TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'

**Gli1-Rand2-M.sssI:**

5'-GATAGTCTCATTTTCACC **AT**GGACNNNNCAAGATGAA TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'

**Gli1-Rand2-unmethylated:**

5'-GATAGTCTCATTTTCACC **TA**GGACNNNNCAAGATGAA TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand3-M.sssI:**

5'-GATAGTCTCATTTTCACC **AT**GGACCACNNNNGATGAA TTCCATTACAGTATCTGT-3'

                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand3-unmethylated:**

5'-GATAGTCTCATTTTCACC **TA**GGACCACNNNNGATGAA TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand4-M.sssI:**

5'-GATAGTCTCATTTTCACC **AT**GGACCACCCANNNNGAA TTCCATTACAGTATCTGT-3'

                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand4-unmethylated:**

5'-GATAGTCTCATTTTCACC **TA**GGACCACCCANNNNGAA TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand5-M.sssI:**

5'-GATAGTCTCATTTTCACC **AT**GGACCACCCAAGANNNN TTCCATTACAGTATCTGT-3'

                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand5-unmethylated:**

5'-GATAGTCTCATTTTCACC **TA**GGACCACCCAAGANNNN TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**Gli1-Rand6-M.sssI**

5'-GATAGTCTCATTTTCACC **AT**GGACCACCCACNNNGAA TTCCATTACAGTATCTGT-3'
                                        3'AAGGTAATGTCATAGACA-5'FAM'

**Gli1-Rand6-unmethylated:**

5'-GATAGTCTCATTTTCACC **TA**GGACCACCCACNNNGAA TTCCATTACAGTATCTGT-3'
                                        3'AAGGTAATGTCATAGACA-5'FAM'

**Gli1-Rand6-top Hemimethylated:**

5'-GATAGTCTCATTTTCACC **CT**GGACCACCCA**M**NNNGAA TTCCATTACAGTATCTGT-3'

3'AAGGTAATGTCATAGACA-5'FAM'

**Gli1-Rand6-duplex Methylated:**
5'-GATAGTCTCATTTTCACC **GA**GGACCACCCA**M**-3'
3'-CTATCAGAGTAAAAGTGG **CT**CCTGGTGGGTG**M**NNNTT AAGGTAATGTCATAGACA-5'


**AP1-spec:**
5'-GATAGTCTCATTTTCACC CCGTGAAANNNNGTCATTG TTCCATTACAGTATCTGT-3'
                                                3'AAGGTAATGTCATAGACA-5'FAM'

**AP1-Mspec_Both:**
5'-GATAGTCTCATTTTCACC **AT**TTTCAGAGNNN**M**GTCAG TTCCATTACAGTATCTGT-3'
                                  3'-**M**AGTC AAGGTAATGTCATAGACA-5'FAM'

**AP1-Mspec_Bottom:**
5'-GATAGTCTCATTTTCACC **CG**TTTCAGAGNNNCGTCAG TTCCATTACAGTATCTGT-3'
                                  3'-**M**AGTC AAGGTAATGTCATAGACA-5'FAM'

**AP1-Mspec_Top:**
5'-GATAGTCTCATTTTCACC **GC**TTTCAGAGNNN**M**GTCAG TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'

**AP1-Mspec_Non:**
5'-GATAGTCTCATTTTCACC **TA**TTTCAGAGNNNCGTCAG TTCCATTACAGTATCTGT-3'
                                        3'-AAGGTAATGTCATAGACA-5'FAM'


**HOXB13-unmethylated:**
5'-CACGACGCTCTTCCGATCT **AG**CCNNNTAAAC TTCCATGACAGTATCTGT-3'
                                      3'-AAGGTACTGTCATAGACA-5'FAM'

**HOXB13-top hemimethylated:**
5'-CACGACGCTCTTCCGATCT **TC**CCN**M**NTAAAC TTCCATGACAGTATCTGT-3'
                                      3'-AAGGTACTGTCATAGACA-5'FAM'

**HOXB13-bottom hemimethylated:**
5'-CACGACGCTCTTCCGATCT **GA**CCNNGTAAAC TTCCATGACAGTATCTGT-3'
                            3'-**M**ATTTG AAGGTACTGTCATAGACA-5'FAM'


**HOXB13-duplex methylated:**
5'-CACGACGCTCTTCCGATCT **CT**CCN**M**GTAAAC TTCCATGACAGTATCTGT-3'
                            3'-**M**ATTTG AAGGTACTGTCATAGACA-5'FAM'

# Supplementary text S2: Instructions for software use.

### Regression analysis

A website is available for the regression analysis (http://stormo.wustl.edu/cgi-bin/dgranas/motif_mlr.pl). The data are each sequence followed by its binding energy, separated by white space. Sequences can contain any of the letters A, C, G, T, M and W. The first sequence is used as the reference and should not contain M or W. Its energy will be defined as 0 and all of the other energies adjusted to maintain their difference from the reference. The regression returns an energy PWM (ePWM) that provides the best fit values for the energies of every base at every position relative to the reference base defined as energy 0. This includes values for M and W. In the case that M and W always occur together (the typical situation for CpG methylation), the energy for the pair is arbitrarily assigned half to each position.

### Meth-eLogos

Two additional versions of the PWM are produced by the regression website. The first is formatted for use by the Logo website (http://stormo.wustl.edu/EnergyModel). It may differ from the regression ePWM because the values for M and W are the difference between M and C and between W and G in the original ePWM. The numbers are also adjusted so that the mean value of the energies at each position (not including M and W) is 0. That is because the Meth-eLogo plots energy differences from the mean, and M and W as the energy changes when C is methylated (on either strand). That ePWM can be pasted into the Logo page. Designate that it is an energy logo, and if desired include methylation energies.

### Sequence searching

Scoring sites in a sequence is accomplished using the PatSer program (http://stormo.wustl.edu/consensus/cgi-bin/Server/Interface/patser.cgi). From the regression page copy the patser version of the ePWM and paste it into the "matrix" box on the patser page. This ePWM is modified from the regression one in three ways. As with the eLogo PWM it is adjusted to have mean energy of 0 (not including M and W) and it is put in the proper format. In addition, the signs are all changed. That is because lower energy corresponds to higher affinity, but in most bioinformatic analyses higher scores correspond to better binding sites. On the patser page the sequences to be searched are entered in the "sequence file" box, or they can be uploaded using the browser. Check that you have entered a "weight matrix" and set the alphabet to A:T C:G M:W. The ":" defines complementary bases so that both strands can be searched if desired (check the "score complementary sequences" box). You can print just the highest score on each sequence, all scores above some designated score, or the default is to print scores for every position. To search the same sequence with two different methylation states, in one version substitute the methylated Cs and Gs with Ms and Ws. Then a comparison of the two lists of binding energies will identify those positions that are differentially bound under different methylation states.