# Programmable Sensors of 5-Hydroxymethylcytosine

Grzegorz Kubik, Sabrina Batke, and Daniel Summerer*

Department of Chemistry, Zukunftskolleg, and Konstanz Research School Chemical Biology, University of Konstanz, Universitätsstraße 10, 78457 Konstanz, Germany

**S** *Supporting Information*

**ABSTRACT:** 5-Hydroxymethylcytosine (hmC), the sixth base of the mammalian genome, is increasingly recognized as an epigenetic mark with important biological functions. We report engineered, programmable transcription-activator-like effectors (TALEs) as the first DNA-binding receptor molecules that provide direct, individual selectivities for cytosine (C), 5-methylcytosine (mC), and hmC at user-defined DNA sequences. Given the wide applicability of TALEs for programmable targeting of DNA sequences in vitro and in vivo, this provides broad perspectives for epigenetic research.

The epigenetic nucleobase 5-methylcytosine (mC, Figure 1A) plays important roles in gene expression regulation, genome stability, development, and disease.[1] It was recently
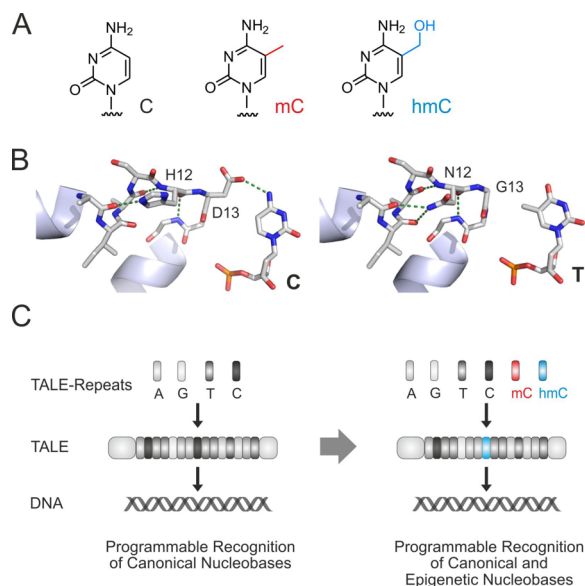


**Figure 1.** Direct, programmable differentiation between C, mC, and hmC in user-defined DNA sequences by engineered TALEs. (A) Chemical structure of C, mC, and hmC. (B) Interaction of RVD HD (amino acids 12 and 13 of TALE repeat) with cytosine (C, left) and of RVD NG with thymine (T, right) in a crystal structure of a TALE-DNA complex (pdb entry 3V6T).[19] Hydrogen bonds are shown as dotted green lines. (C) Concept of direct, programmable differentiation between canonical nucleobases in DNA by TALEs via modular assembly of TALE repeats for A, G, T, and C (left), and additionally between epigenetic nucleobases enabled by TALE repeats with individual selectivities for C, mC, and hmC (right).

discovered that ten-eleven translocation (TET) proteins catalyze the oxidation of mC to 5-hydroxymethylcytosine (hmC, Figure 1A),[2,3] 5-formylcytosine (fC), and 5-carboxylcytosine (caC).[4-6] fC and caC can be removed from DNA and replaced by cytosine (C) via base excision repair.[5] This cycle of methylation, oxidation, and repair now offers a plausible model for the dynamic epigenetic modification of mammalian DNA.[7] Moreover, whereas the intermediates fC and caC exist at only relatively low levels in DNA,[8] hmC exhibits high levels in many cell types, and emerging data link it to important biological functions: hmC exhibits unique genomic distribution,[8] (cancer) cell-specific occurrences,[7] and altered protein recruitment abilities.[9-11]

Key to a deeper understanding of the biological roles of hmC is its locus-specific detection and the hmC-conditional activity-control of loci. This requires effective and flexible strategies to differentiate (i.e., selectively bind or not bind) between C, mC, and hmC at user-defined sequence positions. However, a direct programmable differentiation by Watson–Crick base pairing, which greatly facilitates canonical DNA sequence analysis, is not available for C, mC, and hmC since these exhibit similar pairing properties. Hence, DNA-pretreatments are employed to first differentiate between these nucleobases, and their canonical sequence position is then revealed by analyses relying on base pairing.[12] Two types of DNA-pretreatment are available. Chemical conversion, exploiting the unique reactivity of hmC in β-glucosyl-transfer-,[13-16] in oxidation-,[17] and in bisulfite-reactions,[18] often employed in combination. Alternatively, direct binding by antibodies can be used.[20] As single exceptions, DNA polymerases[21,22] and nanopores[23-25] can differentiate between all four canonical nucleobases as well as between mC and hmC but are limited to single molecule setups in vitro.

Here, we report engineered transcription-activator-like effector (TALE) proteins[26,27] as the first DNA binding receptor molecules that can directly differentiate between C, mC, and hmC in user-defined DNA sequences. TALEs consist of multiple concatenated repeats, each of which selectively recognizes one nucleobase through one of two variable amino acids (repeat variable diresidue, RVD). This recognition follows a simple code with the RVDs NI, NN (NH), NG, and HD (amino acid positions 12 and 13 within the TALE repeat) preferentially binding A-, G-, T-, and C-, respectively.[28-30] We have recently reported the direct sensing of mC in user-defined DNA sequences in vitro[31] based on the ability of RVD HD (which interacts with the 4-amino group of C via a hydrogen bond with the aspartate carboxyl group, Figure 1B)[19,33] to differentiate between C and mC.[34,35] This provides sensitive detection of the

status and level of mC at single positions, in various sequence contexts, and at various positions within the TALE-DNA complex, and enabled single mC detection in the zebrafish genome.[31,32]

To expand this concept to hmC, we aimed to define a toolbox of TALE repeats, each of which provides individual selectivity for C, mC, or hmC (Figure 1C). For engineering and in vitro testing of TALEs, we constructed vector pTRX-ENTRY that is compatible with widely used protocols for the hierarchical assembly of TALE repeat arrays (Figure 2A).[36] pTRX-ENTRY
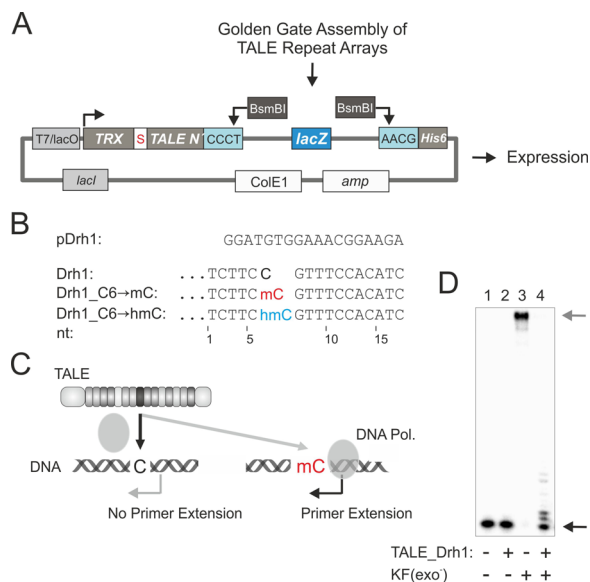


**Figure 2.** Construction and analysis of TALEs. (A) Overview of pTRX-ENTRY for the assembly of vectors for high-level expression and purification of TALEs in *E. coli*. TRX, *E. coli* thioredoxin; S, S-tag. (B) Used oligonucleotide primer and templates with C, mC, or hmC at template position 6 (only 3′-terminal 17 nt of 79 nt templates are shown). (C) Principle of TALE-controlled primer extension shown for C and mC. TALE and DNA polymerase compete for the same binding site in the primer template complex. Binding and nonbinding of TALE to DNA is shown with a black and gray arrow, respectively. (D) PAGE analysis of primer extension reactions as shown in panel C containing 8.325 nM primer−template complex in the presence or absence of 416 nM TALE_Drh1 and 25 mU KF(exo⁻) as indicated. Experiments with TALE_Drh1(HD) contained 832.5 nM TALE. Primer and extension products are marked with black and gray arrows, respectively.

serves as final entry vector for the assembly/*lacZ* screening (>99% success rate, see Supporting Information (SI) Figure 1) of constructs that enable high-level expression and purification of TALEs in *E. coli* via an N-terminal thioredoxin domain and a C-terminal His6 tag (based on a *Xanthomonas axonopodis* TALE scaffold, see SI).[31] Using this approach, we designed and expressed TALE_Drh1 (SI Figure 2), targeting the 17 nt sequence Drh1 (Figure 2B). To study the selectivity of TALE repeats, we employed an assay based on the ability of TALEs to control DNA replication, since this is a process that underlies a large variety of DNA detection methodologies.[31] The assay allows for quantitative analysis of TALE−DNA interactions and previously enabled the detection of single genomic mC (Figure 2C shows the setup with DNA containing a single C or mC). For this, the 5′-³²P-labeled primer pDrh1 is hybridized to a 3-fold excess of DNA oligonucleotide templates containing sequence Drh1 (Figure 2B) by heating to 95 °C for 5 min and cooling to room temperature over 30 min. This complex is incubated with

TALE protein for 30 min, and subsequently, 100 $\mu$M dNTP and Klenow fragment of *E. coli* DNA polymerase I (3′-5′-exo⁻, KF(exo⁻) are added (concentrations of TALE and KF(exo⁻) varied and are indicated in the figures). The mixture is incubated for 15 min at room temperature, denatured by the addition of formamide/EDTA and then resolved by denaturing polyacrylamide gel electrophoresis (PAGE). This enables quantitative analysis of TALE binding by KF(exo⁻) inhibition through quantification of primer extension product (Figure 2D).

A possible approach for the design of custom TALE repeats is to exploit the differential potential of C, mC, and hmC to build hydrogen bonds with polar RVDs, arising from different shielding of the amino group and the unique availability of a hydroxyl group in the vicinity to the 5-position of hmC (Figure 1A,B). We performed primer extensions with TALE_Drh1(HD) bearing RVD HD in TALE repeat 6, and templates Drh1, Drh1_C6 → mC, and Drh1_C6 → hmC bearing a single C, mC, or hmC at position 6 of the TALE_Drh1 binding sequence (Figure 2B). As expected, RVD HD exhibited strong binding to C, but not to mC (Figure 3A; for full data see SI Figure 3).[31] Moreover, it did not bind to hmC, which suggests that D13 of RVD HD is not able to undergo a stabilizing interaction with the hydroxyl group of hmC, establishing HD as a fully selective RVD for C. To design hmC-selective RVDs, we increased the conformational flexibility of the carboxylic acid linker at position 13 to potentially facilitate hydrogen bonding. We analyzed
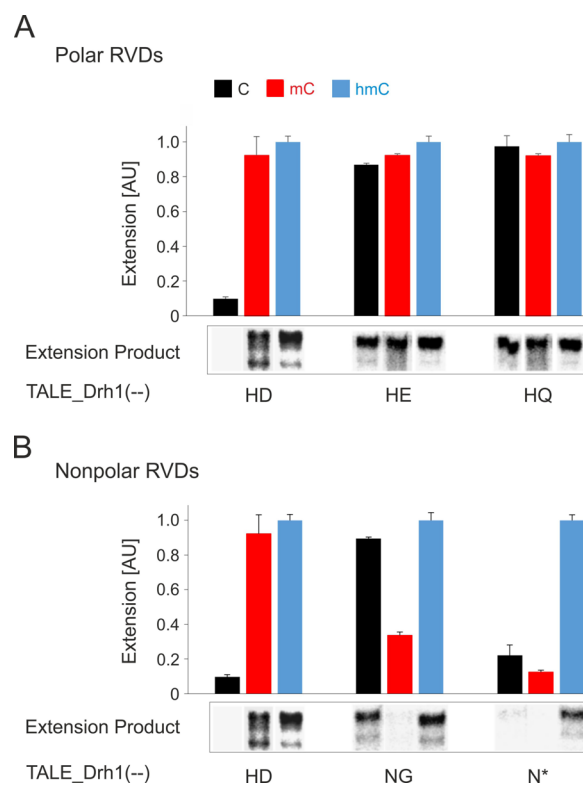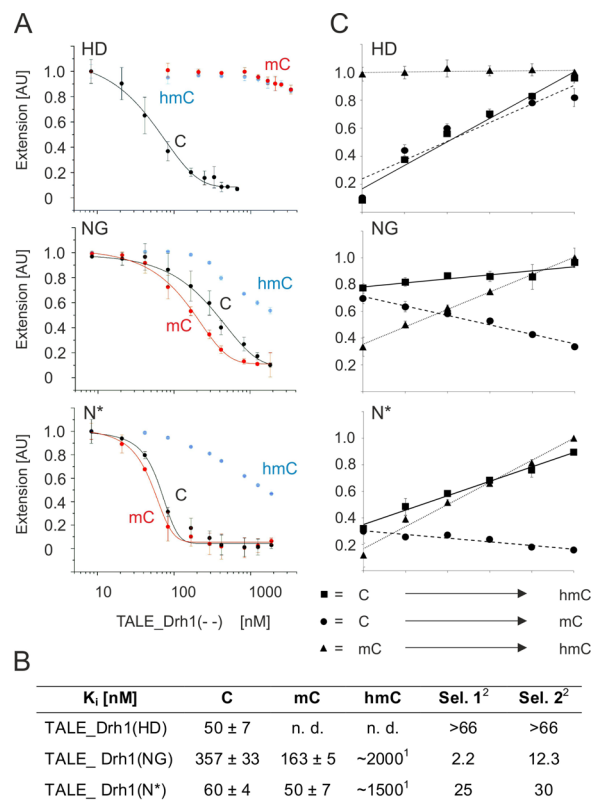


**Figure 3.** Designer TALE repeats for the differentiation between C, mC, and hmC in DNA. (A) TALEs bearing RVD HD and polar RVDs HE and HQ for the differentiation between C (black), mC (red), and hmC (blue) in the template, analyzed in primer extension assays. (B) TALEs bearing RVD HD and nonpolar RVDs NG and N* for the differentiation between C, mC, and hmC in the template, analyzed in primer extension assays. Variable RVD was present in TALE repeat 6. Error bars are from duplicate experiments. Data are normalized for hmC to facilitate comparison of selectivities.

TALE_Drh1(HE) as above, bearing a D13 → E mutation in TALE repeat 6. However, this resulted in a significantly reduced binding to C without promoting the binding to mC or hmC and, thus, in an inability to differentiate between any of the three nucleobases (Figure 3A). A similar result was obtained when an amide group instead of a carboxyl group was installed (TALE_Drh1(HQ), Figure 3A). These data suggest that selective hydrogen bonding of the tested TALE repeats to hmC may be prevented by unfavorable steric interactions between the enlarged side chains of E13 and Q13 and C, mC, and hmC.

An alternative approach for the design of selective TALE repeats is to exploit the different steric demand of the 5-substituents that increases from C over mC to hmC by constructing RVDs with gradually decreasing size. In fact, the 5-methyl group of mC can be accommodated by nonpolar RVDs without side chain at position 13 (i.e., RVD NG) or even with a complete deletion of this amino acid (i.e., RVD N*, * = deletion), indicating that this could be a viable approach.[35,32]

We analyzed TALE_Drh1(NG) bearing RVD NG as above. Compared to TALE_Drh1(HD), this led to significantly reduced binding to C, whereas, in contrast to all tested polar RVDs, binding to mC was now observed (Figure 3B; for full data see SI Figure 3). Strikingly, no binding to hmC was observed, suggesting that the further increased steric demand of the 5-hydroxymethyl group led to unfavorable interactions. This resulted in an overall positive selectivity of RVD NG for mC, in addition to the already established overall positive selectivity of RVD HD for C (for full selectivity profiles of the RVDs, see SI Figure 15). We next asked, if our approach could also be exploited for the design of TALE RVDs selective for hmC. We replaced RVD NG by RVD N* to further reduce the steric demand of the RVD. This indeed resulted in a significantly increased binding to both C and mC, whereas still no significant binding to hmC was observed (TALE_Drh1(N*), Figure 3B). This suggests that the deletion of G13 enables accommodation of 5-hydrogen atoms and 5-methyl groups, but not of 5-hydroxymethyl groups, resulting in an overall negative selectivity of RVD N* for hmC (RVD S* exhibited a related hmC selectivity, but with somewhat reduced binding to C and with lower overall affinity; see SI Figure 16). Taken together, these experiments defined a toolbox of three TALE repeats with individual selectivities for C, mC, and hmC in user-defined DNA sequences (selectivities of RVDs HD, NG, and N* were confirmed in a second sequence context at a different position in the TALE−DNA complex; see SI Figure 17).

To gain quantitative insights into the affinities and selectivities of RVDs HD, NG, and N*, we performed primer extension reactions with varying TALE concentrations (Figure 4A shows the inhibition profiles; for full data see SI Figures 4−9). TALE_Drh1(HD) exhibited a $K_i$ of 50 nM for sequence Drh1 bearing a C at position 6 (Figure 4B). In contrast, only slight inhibition was observed for mC and hmC even at the highest possible TALE concentrations (close to the solubility limit), indicating a very high selectivity (>66-fold, Figure 4A), but preventing $K_i$ determination (for a closely related TALE, a 75-fold difference in inhibition for C and mC opposite RVD HD was observed using a qPCR-based assay with increased sensitivity and dynamic range; see SI Figures 13 and 14). TALE_Drh1(NG) exhibited a 3-fold higher $K_i$ for its cognate nucleobase mC, with a moderate selectivity over C, but a pronounced 12-fold selectivity over hmC (Figure 4B). Strikingly, TALE_Drh1(N*) exhibited a low $K_i$ for C and mC, comparable to the one of TALE_Drh1-



**Figure 4.** Quantitative analysis of affinity and selectivity of TALEs with nonpolar TALE repeats. (A) KF(exo⁻) inhibition profiles by TALE_Drh1(HD), (NG), and (N*) with primer template complexes bearing a single C, mC, or hmC at position 6 (Figure 2B). Data were normalized for hmC and fitted using a dose response function (fits are shown as line). Error bars are from duplicate experiments. (B) Inhibition constants $K_i$ of TALE_Drh1(HD), (NG), and (N*) with primer template complexes of panel A. (C) Inhibition experiments with TALE_Drh1(HD), (NG), and (N*) with primer template complexes of panel A employed in mixtures as indicated by arrows (left and right end of x-axis corresponds to 100% of respective nucleobase). Y-axis is the same as that in panel A. Error bars are from duplicate experiments.

Figure 4B table:

| $K_i$ [nM] | C | mC | hmC | Sel. 1[2] | Sel. 2[2] |
|---|---|---|---|---|---|
| TALE_Drh1(HD) | 50 ± 7 | n. d. | n. d. | >66 | >66 |
| TALE_Drh1(NG) | 357 ± 33 | 163 ± 5 | ~2000[1] | 2.2 | 12.3 |
| TALE_Drh1(N*) | 60 ± 4 | 50 ± 7 | ~1500[1] | 25 | 30 |

[1] Estimated
[2] Selectivities, i.e. $K_i$ ratios for cognate vs. two noncognate nucleobases.

(HD) to C. Moreover, a significant increase of its $K_i$ to ~1.5 μM for hmC resulted in a very high overall selectivity of 25-fold and 30-fold versus C and mC, respectively (Figure 4B).

Since C, mC, and hmC can be present as mixtures with different modification levels in genomic samples, we tested if RVDs HD, NG, and N* would allow for selective sensing of single C, mC, and hmC even in mixed DNA samples. We performed analyses as above with template mixtures (Figure 4C; for full data, see SI Figures 10−12). TALE_Drh1(HD) did not exhibit inhibition in mC/hmC mixtures. However, in mixtures containing its cognate nucleobase C, it exhibited a fully linear dependence on the C modification level over the complete range. Similarly, both TALE_Drh1(NG) and TALE_Drh1(N*) showed a strong linear dependence on the presence of their respective cognate nucleobase mC and hmC in the presence of both noncognate nucleobases. In contrast, only a slight response to one of the noncognate nucleobases (C and mC, respectively) was observed (Figure 4C). This shows that RVDs HD, NG, and N* enable the selective sensing of the modification levels of their cognate nucleobases C, mC, and hmC in the presence of noncognate nucleobases. Hence, for the analysis of an unknown

nucleotide position, targeting with RVD N* reveals hmC, and the additional targeting with RVD HD (or NG) reveals C (or mC). Comparison of the two assays then fully reveals C, mC, and hmC.

In conclusion, we report engineered TALE proteins as the first molecules that provide direct, programmable sensing of C, mC, and hmC in user-defined DNA sequences. We define a toolbox of TALE repeats, each of which provides selectivity for its cognate nucleobase even in the presence of the two noncognate nucleobases. Given the subtlety of the structural differences between large 17mer DNA duplexes containing a single C, mC, and hmC nucleobase, the observed selectivities are surprisingly high. TALEs are fully programmable and genetically encoded and thus applicable both in vitro and in a large number of organisms. Moreover, TALEs can be combined with functional domains such as fluorescent proteins, transcriptional activators and repressors, nucleases, and TET proteins.[26,27,34,37,38] We therefore anticipate that our findings will enable the use of engineered TALEs for diverse epigenetic technologies, ranging from the detection of mC and hmC to the activity control and modification of genes conditional on their epigenetic modification.

## ASSOCIATED CONTENT

### Ⓢ Supporting Information

Oligonucleotide sequences, plasmid construction, expression and purification of TALE proteins, and biochemical assays. This material is available free of charge via the Internet at http://pubs.acs.org.

## AUTHOR INFORMATION

### Corresponding Author
*daniel.summerer@uni-konstanz.de

### Notes
The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

## REFERENCES

(1) Law, J. A.; Jacobsen, S. E. *Nat. Rev. Genet.* 2010, 11, 204.
(2) Tahiliani, M.; Koh, K. P.; Shen, Y.; Pastor, W. A.; Bandukwala, H.; Brudno, Y.; Agarwal, S.; Iyer, L. M.; Liu, D. R.; Aravind, L.; Rao, A. *Science* 2009, 324, 930.
(3) Kriaucionis, S.; Heintz, N. *Science* 2009, 324, 929.
(4) Ito, S.; Shen, L.; Dai, Q.; Wu, S. C.; Collins, L. B.; Swenberg, J. A.; He, C.; Zhang, Y. *Science* 2011, 333, 1300.
(5) He, Y. F.; Li, B. Z.; Li, Z.; Liu, P.; Wang, Y.; Tang, Q.; Ding, J.; Jia, Y.; Chen, Z.; Li, L.; Sun, Y.; Li, X.; Dai, Q.; Song, C. X.; Zhang, K.; He, C.; Xu, G. L. *Science* 2011, 333, 1303.
(6) Pfaffeneder, T.; Hackner, B.; Truss, M.; Munzel, M.; Muller, M.; Deiml, C. A.; Hagemeier, C.; Carell, T. *Angew. Chem., Int. Ed.* 2011, 50, 7008.
(7) Wu, H.; Zhang, Y. *Cell* 2014, 156, 45.
(8) Song, C. X.; Yi, C. Q.; He, C. *Nat. Biotechnol.* 2012, 30, 1107.
(9) Mellen, M.; Ayata, P.; Dewell, S.; Kriaucionis, S.; Heintz, N. *Cell* 2012, 151, 1417.
(10) Spruijt, C. G.; Gnerlich, F.; Smits, A. H.; Pfaffeneder, T.; Jansen, P. W.; Bauer, C.; Munzel, M.; Wagner, M.; Muller, M.; Khan, F.; Eberl, H. C.; Mensinga, A.; Brinkman, A. B.; Lephikov, K.; Muller, U.; Walter, J.; Boelens, R.; van Ingen, H.; Leonhardt, H.; Carell, T.; Vermeulen, M. *Cell* 2013, 152, 1146.
(11) Iurlaro, M.; Ficz, G.; Oxley, D.; Raiber, E. A.; Bachman, M.; Booth, M. J.; Andrews, S.; Balasubramanian, S.; Reik, W. *Genome Biol.* 2013, 14, R119.
(12) Plongthongkum, N.; Diep, D. H.; Zhang, K. *Nat. Rev. Genet.* 2014, 15, 647.
(13) Yu, M.; Hon, G. C.; Szulwach, K. E.; Song, C. X.; Zhang, L.; Kim, A.; Li, X.; Dai, Q.; Shen, Y.; Park, B.; Min, J. H.; Jin, P.; Ren, B.; He, C. *Cell* 2012, 149, 1368.
(14) Pastor, W. A.; Pape, U. J.; Huang, Y.; Henderson, H. R.; Lister, R.; Ko, M.; McLoughlin, E. M.; Brudno, Y.; Mahapatra, S.; Kapranov, P.; Tahiliani, M.; Daley, G. Q.; Liu, X. S.; Ecker, J. R.; Milos, P. M.; Agarwal, S.; Rao, A. *Nature* 2011, 473, 394.
(15) Song, C. X.; Szulwach, K. E.; Fu, Y.; Dai, Q.; Yi, C.; Li, X.; Li, Y.; Chen, C. H.; Zhang, W.; Jian, X.; Wang, J.; Zhang, L.; Looney, T. J.; Zhang, B.; Godley, L. A.; Hicks, L. M.; Lahn, B. T.; Jin, P.; He, C. *Nat. Biotechnol.* 2011, 29, 68.
(16) Robertson, A. B.; Dahl, J. A.; Vagbo, C. B.; Tripathi, P.; Krokan, H. E.; Klungland, A. *Nucleic Acids Res.* 2011, 39, e55.
(17) Booth, M. J.; Branco, M. R.; Ficz, G.; Oxley, D.; Krueger, F.; Reik, W.; Balasubramanian, S. *Science* 2012, 336, 934.
(18) Huang, Y.; Pastor, W. A.; Zepeda-Martinez, J. A.; Rao, A. *Nat. Protoc.* 2012, 7, 1897.
(19) Deng, D.; Yan, C.; Pan, X.; Mahfouz, M.; Wang, J.; Zhu, J. K.; Shi, Y.; Yan, N. *Science* 2012, 335, 720.
(20) Ficz, G.; Branco, M. R.; Seisenberger, S.; Santos, F.; Krueger, F.; Hore, T. A.; Marques, C. J.; Andrews, S.; Reik, W. *Nature* 2011, 473, 398.
(21) Flusberg, B. A.; Webster, D. R.; Lee, J. H.; Travers, K. J.; Olivares, E. C.; Clark, T. A.; Korlach, J.; Turner, S. W. *Nat. Methods* 2010, 7, 461.
(22) Summerer, D. *ChemBioChem* 2010, 11, 2499−501.
(23) Wallace, E. V.; Stoddart, D.; Heron, A. J.; Mikhailova, E.; Maglia, G.; Donohoe, T. J.; Bayley, H. *Chem. Commun.* 2010, 46, 8195.
(24) Schreiber, J.; Wescoe, Z. L.; Abu-Shumays, R.; Vivian, J. T.; Baatar, B.; Karplus, K.; Akeson, M. *Proc. Natl. Acad. Sci. U.S.A.* 2013, 110, 18910.
(25) Laszlo, A. H.; Derrington, I. M.; Brinkerhoff, H.; Langford, K. W.; Nova, I. C.; Samson, J. M.; Bartlett, J. J.; Pavlenok, M.; Gundlach, J. H. *Proc. Natl. Acad. Sci. U.S.A.* 2013, 110, 18904.
(26) Boch, J.; Bonas, U. *Annu. Rev. Phytopath.* 2010, 48, 419.
(27) Bogdanove, A. J.; Voytas, D. F. *Science* 2011, 333, 1843.
(28) Moscou, M. J.; Bogdanove, A. J. *Science* 2009, 326, 1501.
(29) Boch, J.; Scholze, H.; Schornack, S.; Landgraf, A.; Hahn, S.; Kay, S.; Lahaye, T.; Nickstadt, A.; Bonas, U. *Science* 2009, 326, 1509.
(30) Yang, J.; Zhang, Y.; Yuan, P.; Zhou, Y.; Cai, C.; Ren, Q.; Wen, D.; Chu, C.; Qi, H.; Wei, W. *Cell Res.* 2014, 24, 628.
(31) Kubik, G.; Schmidt, M. J.; Penner, J. E.; Summerer, D. *Angew. Chem., Int. Ed.* 2014, 53, 6002.
(32) Kubik, G.; Summerer, D. *ChemBioChem* 2014, DOI: 10.1002/cbic.201402408.
(33) Mak, A. N. S.; Bradley, P.; Cernadas, R. A.; Bogdanove, A. J.; Stoddard, B. L. *Science* 2012, 335, 716.
(34) Bultmann, S.; Morbitzer, R.; Schmidt, C. S.; Thanisch, K.; Spada, F.; Elsaesser, J.; Lahaye, T.; Leonhardt, H. *Nucleic Acids Res.* 2012, 40, 5368.
(35) Valton, J.; Dupuy, A.; Daboussi, F.; Thomas, S.; Marechal, A.; Macmaster, R.; Melliand, K.; Juillerat, A.; Duchateau, P. *J. Biol. Chem.* 2012, 287, 38427.
(36) Cermak, T.; Doyle, E. L.; Christian, M.; Wang, L.; Zhang, Y.; Schmidt, C.; Baller, J. A.; Somia, N. V.; Bogdanove, A. J.; Voytas, D. F. *Nucleic Acids Res.* 2011, 39, 7879.
(37) Maeder, M. L.; Angstman, J. F.; Richardson, M. E.; Linder, S. J.; Cascio, V. M.; Tsai, S. Q.; Ho, Q. H.; Sander, J. D.; Reyon, D.; Bernstein, B. E.; Costello, J. F.; Wilkinson, M. F.; Joung, J. K. *Nat. Biotechnol.* 2013, 31, 1137.
(38) Thanisch, K.; Schneider, K.; Morbitzer, R.; Solovei, I.; Lahaye, T.; Bultmann, S.; Leonhardt, H. *Nucleic Acids Res.* 2014, 42, e38.