

# Complete decoding of TAL effectors for DNA recognition

*Cell Research* (2014) 24:628-631. doi:10.1038/cr.2014.19; published online 11 February 2014

## Dear Editor,

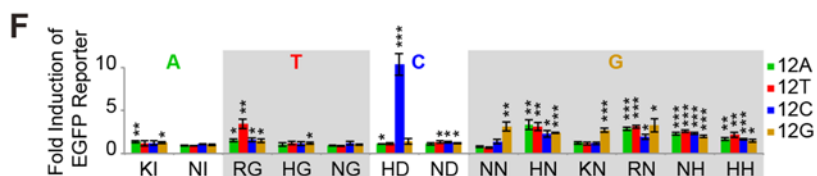
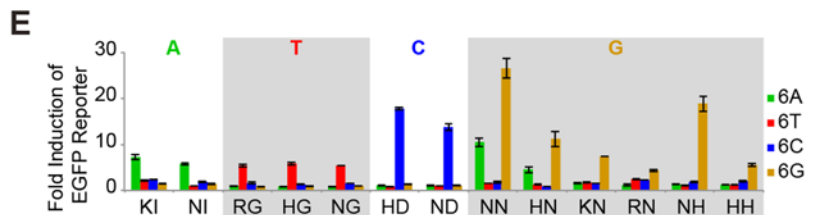
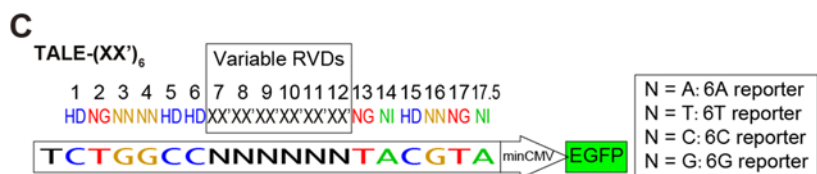
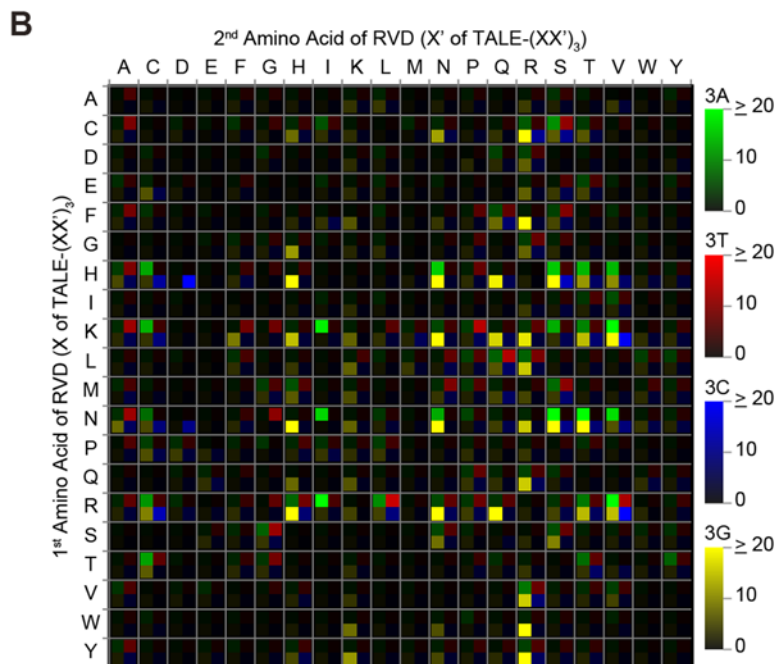
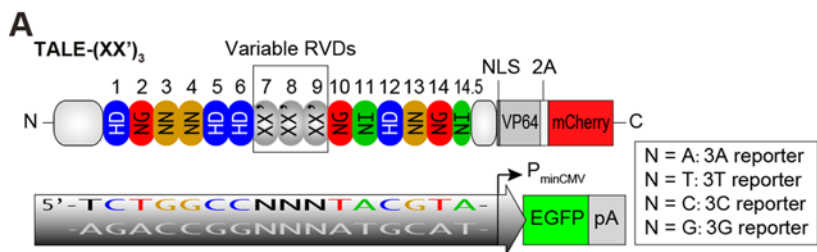
The most striking feature of a transcription activator-like effector (TALE) is the presence of a central DNA-binding region composed of tandem repeats of about 34 amino acids [1]. Two hypervariable residues at positions 12 and 13 (repeat-variable diresidues or RVDs) in each repeat bind to DNA, and this modular DNA-binding feature of TALE repeats has inspired the development of custom-designed TALE repeats for gene editing [2, 3, 4, 5]. The nucleotide recognition preference of the commonly used RVDs has been experimentally or computationally determined [2, 5]. For instance, RVD NN has a high preference for both G and A. The rare RVDs, NK and NH, have better specificity for guanine than NN, but their affinity is relatively lower [3, 6, 7]. We thus decided to conduct a thorough investigation of potential RVDs, which cover all possible combinations of amino acid diresidues, for their DNA recognition capabilities.

We set up a screening platform composed of an artificial TALE-VP64-mCherry construct, which expresses RVD (XX') in 3-tandem repeat format (from 7<sup>th</sup> to 9<sup>th</sup>, TALE-(XX')<sub>3</sub>), and 4 corresponding EGFP reporter constructs, in which potential TALE-(XX')<sub>3</sub>-binding sites composed of 3 consecutive nucleotides (A, T, C or G) are located in front of a minCMV promoter and its downstream *EGFP* gene (Figure 1A, Supplementary information, Figure S1A and Data S1). To test this system, we made a control TALE (TALE-Ctrl) that is identical to TALE-(XX')<sub>3</sub> except for the repeat domain (Supplementary information, Figure S1A), and confirmed that it could not activate any of the 4 EGFP reporters (Supplementary information, Figure S1B), thus serving as the control for basal activity. We then constructed 4 TALE-(XX')<sub>3</sub> expression plasmids by placing the common RVDs (NI, NG, HD and NN) in the middle to target the 3A, 3T, 3C and 3G EGFP reporters, respectively (Supplementary information, Figure S1C). These TALE-(XX')<sub>3</sub> constructs were individually introduced into HEK293T cells together with 1 of the 4 EGFP reporter plasmids to examine their specificities, which were determined

by the fold induction of EGFP expression compared with the basal level (Supplementary information, Data S1). The identity of XX' determined TALE-(XX')<sub>3</sub> specificity on different EGFP reporters: NI, NG, HD and NN predominantly recognized A, T, C and G or A, respectively. This result is consistent with the current knowledge regarding the base preference of these 4 common RVDs, demonstrating that this artificial system is suitable for testing the DNA recognition ability of RVDs (Supplementary information, Figure S1D-S1E).

To quantitatively measure the base preference of all theoretical RVDs, we created a library of TALE-(XX')<sub>3</sub> constructs, which covers a total of 400 types of RVDs, following a special protocol combined with the ULtiMATE assembly method [8] (Supplementary information, Figure S2 and Data S1). X and X' correspond to the 12<sup>th</sup> and 13<sup>th</sup> amino acids in a classical TALE module, respectively. We introduced each of the 400 TALE-(XX')<sub>3</sub> constructs (Supplementary information, Tables S1 and S2) individually into HEK293T cells together with 1 of the 4 EGFP reporter plasmids and measured both the EGFP and mCherry levels by FACS analysis. We then determined the base-recognition efficiencies of the 400 diresidues. A total of 1 600 data points were summarized in 3 formats: heat map (Figure 1B), histograms categorized by the 13<sup>th</sup> residue (X') (Supplementary information, Figure S3) and the 12<sup>th</sup> residue (X) (Supplementary information, Figure S4). The results obtained from this screening provide substantial information regarding the base preference of all theoretical RVDs. In addition to NI, NG, HD and NN, all the natural RVDs and a few artificial RVDs showed base-recognition preferences that were similar to those reported previously [2, 5, 6, 7] (Supplementary information, Table S3). Besides these 25 RVDs, we determined the DNA base-recognition preference of the remaining 375 RVDs that did not evolve naturally and have not been previously examined.

Notably, many of these artificial RVDs showed a distinct preference for DNA bases compared with the 25 reported RVDs, and only a few of them start with 1 of the 2 frequently occurring amino acids, Asn and



**Figure 1** A Complete assessment of TALE RVD efficiencies and specificities. **(A)** Design of the screening system for novel TALE RVDs. **(B)** A heat map generated from library screening of TALE-(XX')<sub>3</sub> with four reporters (3A, 3T, 3C, and 3G) reflecting the base preference of 400 RVDs. EGFP activities from different reporters were coded by different colors representing the reporter identities (3A, green; 3T, red; 3C, blue; 3G, yellow), and the brightness of the colors indicates the fold induction of reporters by TALE-(XX')<sub>3</sub> compared to the basal levels. The single-letter abbreviations for the amino acids are used. **(C)** Design of TALE-(XX')<sub>6</sub> and its corresponding reporters. **(D)** Design of TALE-(XX')<sub>12</sub> and its corresponding reporters. **(E-F)** Base preference of RVDs in TALE-(XX')<sub>6</sub> **(E)** and TALE-(XX')<sub>12</sub> **(F)**. RVDs were clustered by base preference. The x-axis labels indicate the variable RVDs tested in TALE-(XX')<sub>6</sub> or TALE-(XX')<sub>12</sub>. Data are means ± SD, *n* = 3; \**P* < 0.05, \*\**P* < 0.01, and \*\*\**P* < 0.005.

His (Figure 1B, Supplementary information, Figures S3 and S4). From these artificial RVDs, we selected those that showed potential base-recognition preference based on the criteria shown in Supplementary information, Data S1 for further intensive analyses. We found that the adenine recognition ability of KI and RI was similar to that of NI (Supplementary information, Figure S5A). For thymine recognition, we identified 3 additional RVDs aside from NG, which all end with Gly (RG, KG and HG), and seven RVDs that all end with Ala (KA, CA, FA, YA, RA, PA, and AA), but appeared to have higher background, especially for C recognition (Supplementary information, Figure S5B). HD and ND, as reported previously [2, 5, 7], were optimal RVDs for C recognition, with almost no non-specific recognition of other bases (Supplementary information, Figure S5C). Five groups of RVDs were identified to recognize guanine, with each group sharing the same 13<sup>th</sup> residue: Asn (N), His (H), Arg (R), Gln (Q), or Lys (K). Most of these RVDs predominantly recognized guanine except for HN and NN (Supplementary information, Figure S5D). These data support the prediction from previous TALE structural studies suggesting that the 13<sup>th</sup> residues of TALE repeats make the base-specific contact [9, 10]. Nevertheless, our data indicate that the 12<sup>th</sup> residue also affects RVD specificity. For example, with the same N<sub>13</sub>, KN and RN only recognized G, whereas HN and NN recognized both A and G, and LN and MN preferred T and C. Similarly, HQ, KQ and RQ preferentially recognized G, whereas LQ preferred T (Supplementary information, Figures S3-S6).

To further examine the base-recognition preference of RVDs, we created two additional artificial platforms with increased stringency, in which multiple TALE repeats carrying the same RVDs were aligned in tandem: TALE-(XX')<sub>6</sub> and its corresponding EGFP reporter constructs (6A, 6T, 6C, and 6G) (Figure 1C) were used to test RVDs in 6-tandem repeat format, and TALE-(XX')<sub>12</sub> and its corresponding EGFP reporter constructs (12A, 12T, 12C, and 12G) (Figure 1D) were used to test RVDs in 12-tandem repeat format (Supplementary information, Table S1 and Figure S2).

In addition to the 4 most common RVDs, which were used as controls, we mainly chose those that demonstrated outstanding base-recognition specificities from the initial screening. We found that KI and NI functioned similarly with respect to A recognition in the 6-repeat format (Figure 1E). The activities of TALE-(RG)<sub>6</sub> and TALE-(HG)<sub>6</sub> were similar to TALE-(NG)<sub>6</sub> for 6T recognition (Figure 1E), whereas TALE-(KG)<sub>6</sub> showed reduced specificity for 6T (Supplementary information, Figure S7). HD and ND again demonstrated strong C preference (Figure 1E). KN, RN, NH and HH showed specific G recognition with variable efficiencies in 6-tandem repeats, whereas NN and HN recognized both G and A as in the 3-repeat format (Figure 1E), and the 6G preference of TALE-(XX')<sub>6</sub> containing either NR, FR, KH, NK, FK or RQ was significantly reduced (Supplementary information, Figure S7). Interestingly, only TALE-(XX')<sub>12</sub> with RG (for T), HD (for C), NN (for G) and KN (for G) in 12-tandem repeats maintained recognition efficiency and specificity (Figure 1F). This result is somewhat surprising for RG as it is assumed that RVDs ending with Gly cannot form hydrogen bonds with thymine [10]. Consistent with previous reports [6, 7], neither TALE-(NH)<sub>12</sub> nor TALE-(HH)<sub>12</sub> could support 12G reporter activation. Considering the strong preference of NH for G in the 6-repeat format, it is unclear why TALE-(KN)<sub>12</sub> but not TALE-(NH)<sub>12</sub> retained activity for the 12G reporter. By the same token, it is also unclear why TALE-(ND)<sub>12</sub> completely lost its preference for the 12C reporter (Figure 1F). Although the combination of the 12<sup>th</sup> and 13<sup>th</sup> amino acids determines the ultimate binding activity of TALE, the increase of repeat number also leads to the decrease or even complete loss of DNA-recognition activity of TALE, which is likely due to either steric or static repulsion between consecutive TALE repeat units.

To further evaluate these novel RVDs, we applied KN and RG in TALEN assembly in place of NN and NG, respectively, and compared them with conventional RVDs in TALENs-mediated DNA cleavage by measuring indel rates. TALENs<sub>KN</sub> for G-targeting showed similar efficiency in creating indels

as TALENS<sub>NN</sub> in 2 independent tests, and both of them performed better than TALENS<sub>NH</sub>. On the contrary, TALENS<sub>RG</sub>, although functional, were less effective than TALENS<sub>NG</sub> (Supplementary information, Table S4). It is possible that other diresidues newly revealed in this study could function as valid RVDs in recognizing DNA bases with high specificity. However, rigorous tests are needed in order to more accurately determine their DNA recognition capabilities.

In addition, we identified a significant number of RVDs that target multiple DNA bases (Supplementary information, Table S5). The availability of RVDs that target different combinations of bases in a degenerate manner may provide certain flexibility in future application such as engineering of sophisticated genetic circuitry [11].

By further deciphering the DNA base preference of all RVDs, natural or artificial, we can achieve a clear understanding of the mechanism that guides the base preference of TALE RVDs. Comprehensive information regarding the specific DNA associations of all RVDs may improve the application of TAL effectors in bioengineering and precision therapy.

## Acknowledgments

We thank Xiaoyu Li (PKU) for technical advice regarding oligo synthesis and Yanyi Huang (PKU) for critical comments on the manuscript. This work was supported by the National Basic Research Program of China (2010CB911800), the National Natural Science Foundation of China (NSFC31070115, NSFC31170126), and the Peking-Tsinghua Center for Life Sciences.

Junjiao Yang<sup>1,\*</sup>, Yuan Zhang<sup>1,\*</sup>, Pengfei Yuan<sup>1</sup>,  
Yuexin Zhou<sup>1</sup>, Changzu Cai<sup>1</sup>, Qingpeng Ren<sup>1</sup>,  
Dingqiao Wen<sup>1,3</sup>, Coco Chu<sup>3</sup>, Hai Qi<sup>2</sup>, Wensheng Wei<sup>1</sup>

<sup>1</sup>State Key Laboratory of Protein and Plant Gene Research, College of Life Sciences, Peking University, Beijing 100871, China; <sup>2</sup>Tsinghua-Peking Center for Life Sciences, Laboratory of Dynamic Immunobiology, School of Medicine, Tsinghua University, Beijing 100084, China; <sup>3</sup>Current address: Department of Computer Science, Rice University, Houston, TX 77005, USA

\*These two authors contributed equally to this work.

Correspondence: Wensheng Wei

Tel: +86-10-62757227

E-mail: wswei@pku.edu.cn

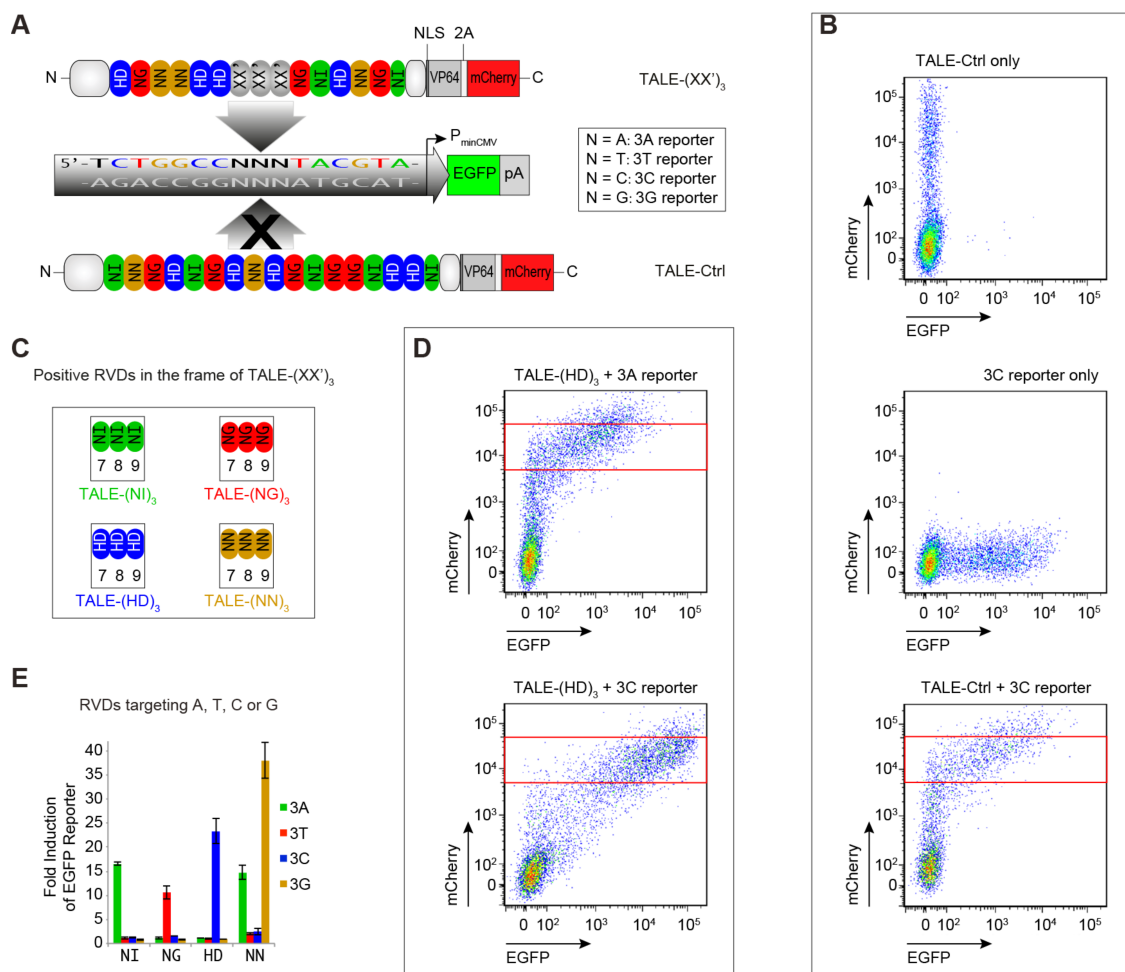
## References

- 1 Boch J, Bonas U. *Annu Rev Phytopathol* 2010; **48**:419-436.
- 2 Boch J, Scholze H, Schornack S, et al. *Science* 2009; **326**:1509-1512.
- 3 Miller JC, Tan S, Qiao G, et al. *Nat Biotech* 2011; **29**:143-148.
- 4 Bogdanove AJ, Voytas DF. *Science* 2011; **333**:1843-1846.
- 5 Moscou MJ, Bogdanove AJ. *Science* 2009; **326**:1501.
- 6 Streubel J, Blucher C, Landgraf A, et al. *Nat Biotech* 2012; **30**:593-595.
- 7 Cong L, Zhou R, Kuo YC, et al. *Nat Commun* 2012; **3**:968.
- 8 Yang J, Yuan P, Wen D, et al. *PLoS One* 2013; **8**:e75649.
- 9 Deng D, Yan C, Pan X, et al. *Science* 2012; **335**:720-723.
- 10 Mak AN, Bradley P, Cernadas RA, et al. *Science* 2012; **335**:716-719.
- 11 Aouida M, Piatek MJ, Bangarusamy DK, et al. *Curr Genet* 2013 Oct 1. doi: 10.1007/s00294-013-0412-z

(Supplementary information is linked to the online version of the paper on the *Cell Research* website.)

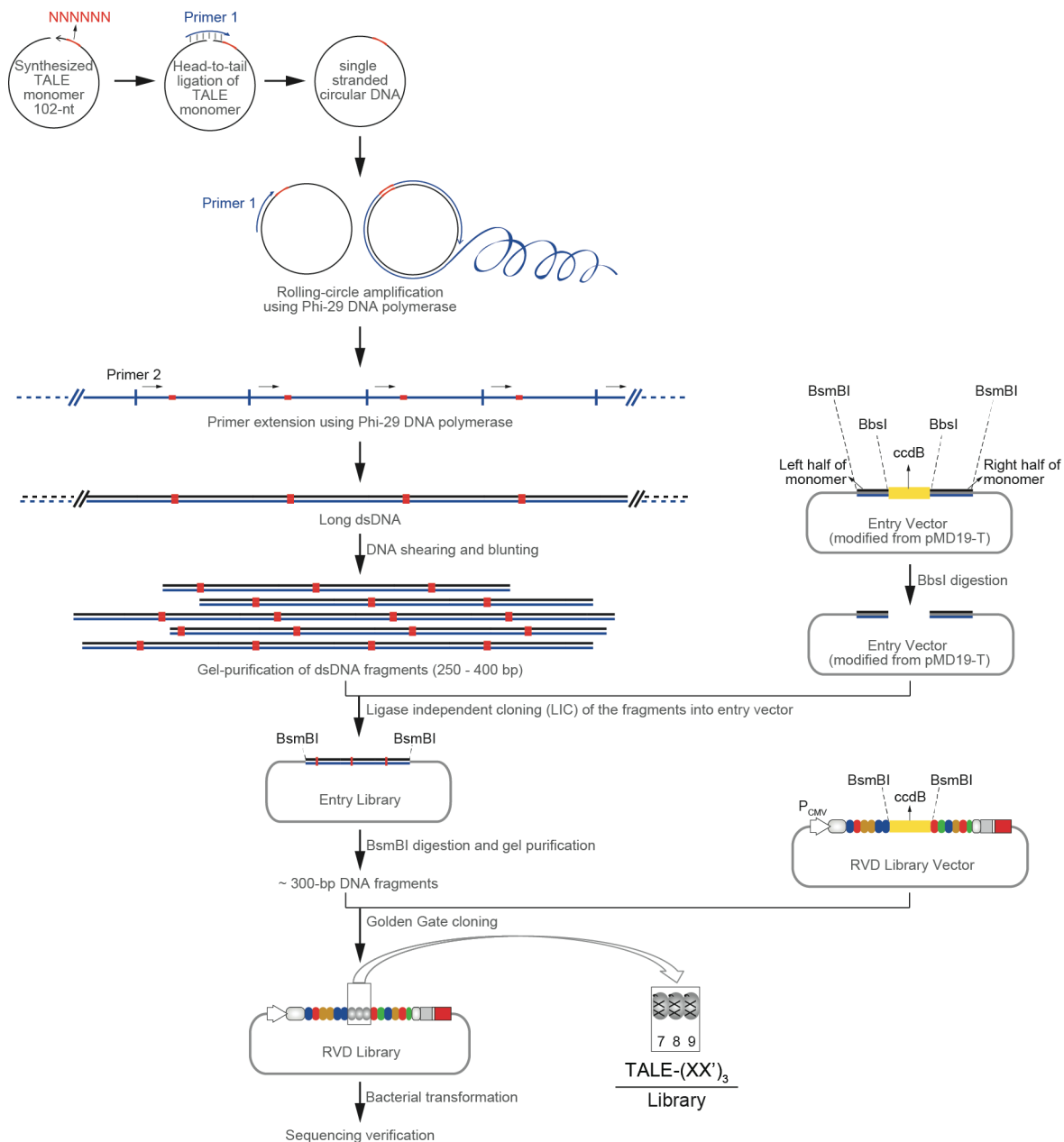


This work is licensed under the Creative Commons Attribution-NonCommercial-Share Alike Works 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0>



**Supplementary information, Figure S1** Screening system for the assessment of TALE RVD efficiencies and specificities. **(A)** Design of the screening system for novel TALE RVDs. The customized TALEs used for RVD screening contained 14.5 repeats fused with the VP64 trans-activation domain and 2A peptide-linked mCherry. The variable diresidues (XX') for testing were placed in the 7<sup>th</sup> - 9<sup>th</sup> repeat modules, and the customized TALE was designated as TALE-(XX')<sub>3</sub>. X and X' represents the 12<sup>th</sup> and 13<sup>th</sup> amino acids in the 7<sup>th</sup> - 9<sup>th</sup> repeat modules, respectively. To determine the DNA recognition specificity of variable RVDs, four reporters were constructed, consisting of TALE-(XX')<sub>3</sub> binding sites with three consecutive nucleotides (A, T, C or G) substituted at positions 7 - 9 in front of a minimal CMV promoter (P<sub>minCMV</sub>) and its downstream EGFP gene. Construct encoding TALE-Ctrl has the identical backbone as TALE-(XX')<sub>3</sub> except that its TALE repeat region is different

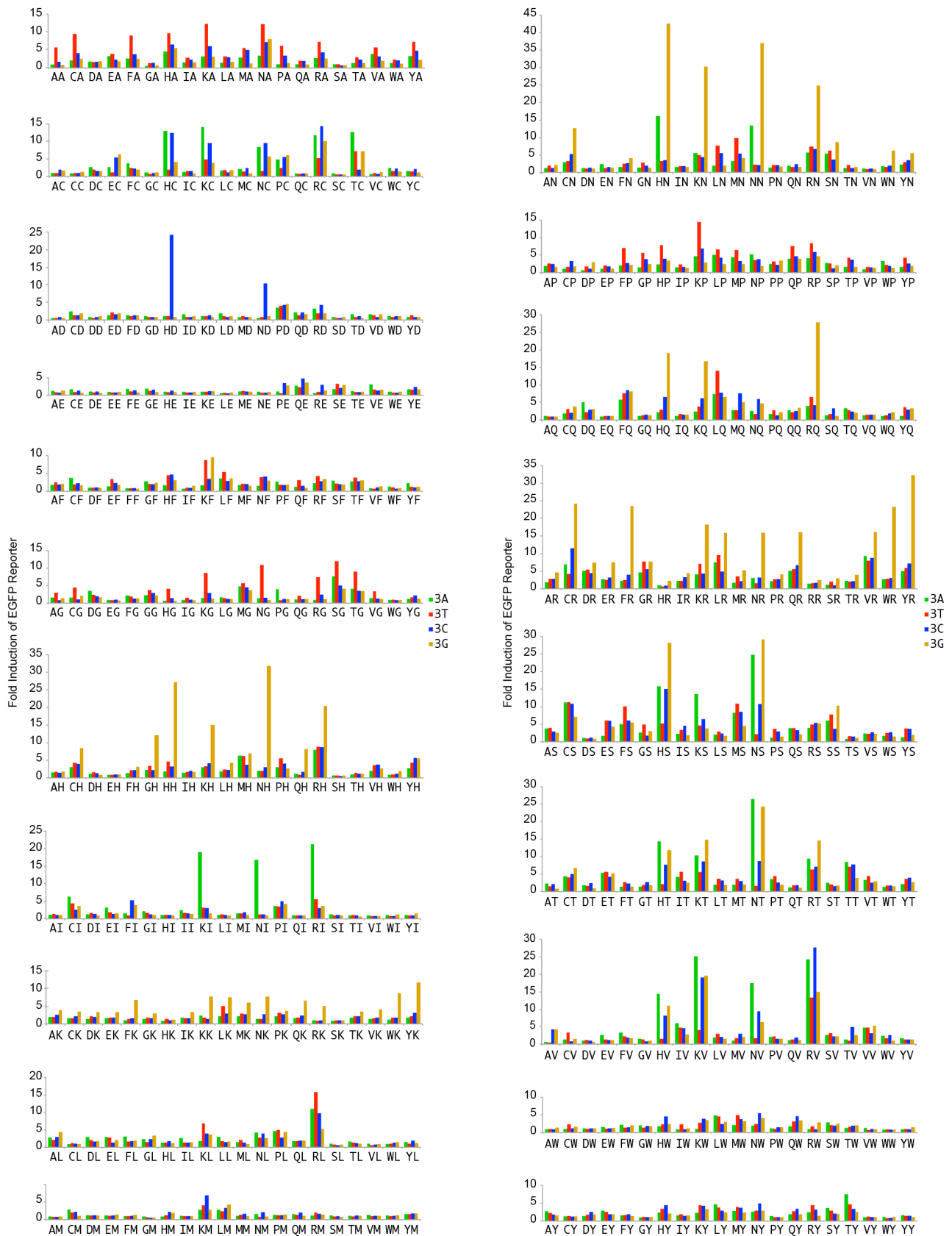
as indicated, which does not match with any reporters. **(B)** FACS analysis of HEK293T cells transfected with TALE-Ctrl only, 3C reporter only and TALE-Ctrl plus 3C reporter (from top to bottom). Red box indicates the region of data collection. **(C)** Customized TALE-(XX')<sub>3</sub> for testing the DNA binding activity of commonly used RVDs (NI, NG, HD and NN). **(D)** Representative FACS analysis of HEK293T cells co-transfected with TALE-(HD)<sub>3</sub> and the 3A (top) or 3C reporter (bottom). Red boxes indicate the region of data collection. **(E)** Base binding activity of the common RVDs in **(C)**. The horizontal axis labels indicate the variable RVD (XX') for testing in TALE-(XX')<sub>3</sub>. The color-coded bars represent the fold induction levels of the different reporters (A, green; T, red; C, blue; and G, yellow) in this and the subsequent figures. Data are means ± s.d., n = 3.



**Supplementary information, Figure S2** Schematic of TALE-(XX')<sub>3</sub> library construction for novel RVD screening. A 102-nt monomer encoding a standard TALE repeat unit, containing six random nucleotides at the RVD-encoding region, was synthesized and subsequently cyclized. Rolling-circle amplification of these single-strand circular DNA templates was conducted using phi29 DNA polymerase and primer 1, and dsDNA fragments were obtained from primer extension using primer 2. After ultrasonic shearing followed by DNA blunting, 250-400 bp DNA fragments were isolated. After gel purification,

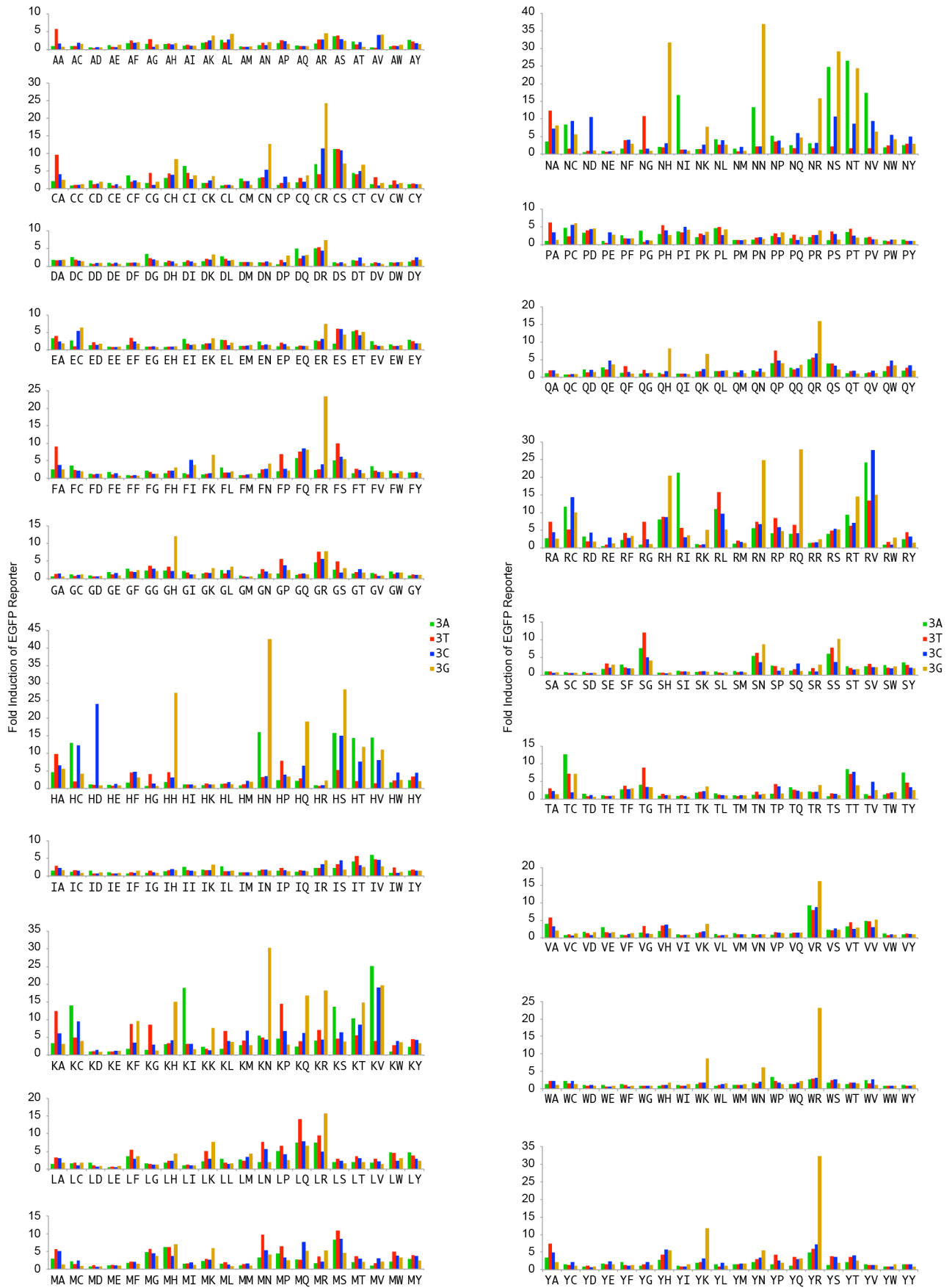
these DNA fragments were cloned into a pre-made entry vector through the LIC method. BsmBI digestion of clones in the entry library produced ~300 bp DNA fragments, which were subsequently ligated into a pre-made RVD library vector. After bacterial transformation and sequencing validation, we were able to obtain 400 types of plasmids encoding customized TALEs with three repeat modules in the middle (7<sup>th</sup> to 9<sup>th</sup>) carrying the variable RVDs for testing. A detailed protocol is provided in the Supplementary Methods.



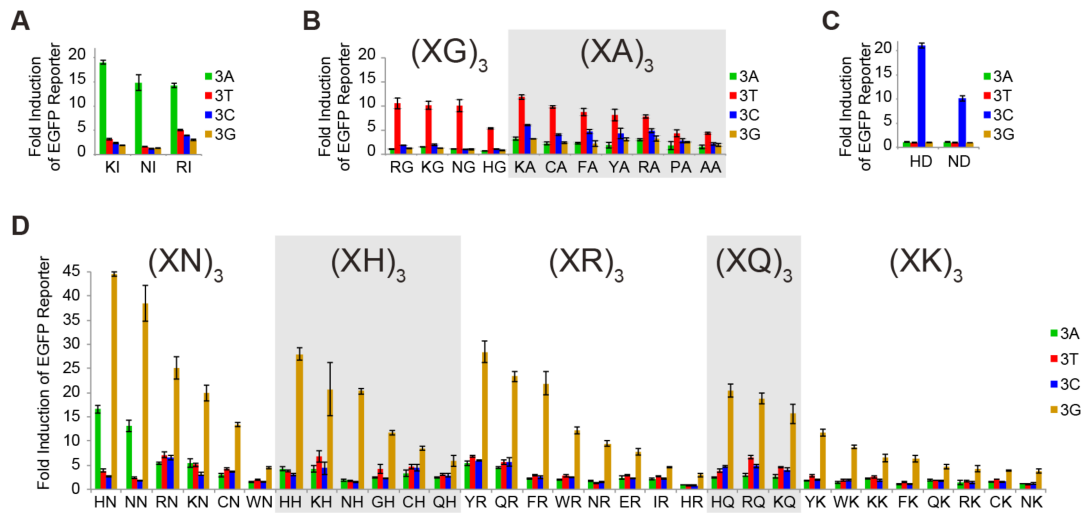


**Supplementary information, Figure S3** Base preference of 400 RVDs from the screening of the TALE-(XX')<sub>3</sub> library. The binding efficiencies and specificities of the

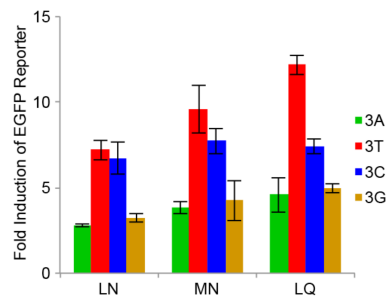
variable RVDs (XX') in each customized TALE were assayed based on the fold induction of EGFP reporters, by comparing with the basal level of EGFP in HEK293T cells transfected with the reporter plasmid and the customized TALE plasmid containing unmatched TALE repeats. The EGFP fluorescence intensity assayed from FACS analysis was normalized to the corresponding mCherry fluorescence intensity. The fold induction of the EGFP reporter of all 400 types of TALE-(XX')<sub>3</sub> (corresponding to four reporters) was categorized by the 13<sup>th</sup> residue (X'), and the data are listed in alphabetical order according to the 12<sup>th</sup> residue (X) of the variable RVD-carrying TALE. The x-axis labels indicate the variable RVDs (XX') tested in TALE-(XX')<sub>3</sub>. The color-coded bars represent the fold induction of the different reporters (A, green; T, red; C, blue; and G, yellow) in this and the subsequent figures.



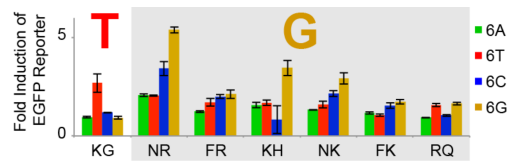
**Supplementary information, Figure S4** Base preference of 400 RVDs from the screening of the TALE-(XX')<sub>3</sub> library. The data points are the same as those in Figure 1B and Figure S3. The fold induction of the EGFP reporter of all 400 types of TALE-(XX')<sub>3</sub> (corresponding to the four reporters) was categorized by the 12<sup>th</sup> residue (X), and the data are listed in alphabetical order according to the 13<sup>th</sup> residue (X') of the variable RVD-carrying TALE. The x-axis labels indicate the variable RVDs (XX') tested in TALE-(XX')<sub>3</sub>.



**Supplementary information, Figure S5** Base preferences of RVDs. RVDs were clustered by position 13 (X') of the variable RVDs and ranked by fold induction of EGFP in the corresponding reporter construct. The x-axis labels indicate the variable RVDs tested. Data are means  $\pm$  s.d., n = 3. **(A)** Adenine-targeting RVDs. **(B)** Thymine-targeting RVDs. **(C)** Cytosine-targeting RVDs. **(D)** Guanine-targeting RVDs.



**Supplementary information, Figure S6** Base preference of RVDs in TALE-(XX')<sub>3</sub>. The x-axis labels indicate the variable RVDs tested in TALE-(XX')<sub>3</sub>. Data are means ± s.d., n = 3.



**Supplementary information, Figure S7** Base preference of RVDs in TALE-(XX')<sub>6</sub>. The results of the fold induction of TALE-(XX')<sub>6</sub> (corresponding to the four reporters) are shown on the y-axis, and the RVDs are clustered by base/reporter preference and ranked by fold induction of EGFP of the corresponding reporter. The x-axis labels indicate the variable RVDs tested in TALE-(XX')<sub>6</sub>. Data are means  $\pm$  s.d., n = 3.









**Supplementary information, Table S2.** Sequences of amino acid diresidues in RVDlibraries (TALE-(XX')<sub>3</sub>, TALE-(XX')<sub>6</sub> and TALE-(XX')<sub>12</sub>).

| Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) |
|------------------|--------------|------------------|--------------|------------------|--------------|------------------|--------------|------------------|--------------|
| AA               | gctgcg       | CA               | tgcgcc       | DA               | gacgct       | EA               | gaggcg       | FA               | ttgccc       |
| AC               | gcctgc       | CC               | tgctgc       | DC               | gattgc       | EC               | gaatgt       | FC               | ttctgc       |
| AD               | gccgac       | CD               | tgcgac       | DD               | gacgac       | ED               | gaagac       | FD               | ttcgac       |
| AE               | gccgaa       | CE               | tgcgag       | DE               | gacgaa       | EE               | gaggag       | FE               | ttcgag       |
| AF               | gccttc       | CF               | tgcttt       | DF               | gatttc       | EF               | gagttc       | FF               | ttcttc       |
| AG               | gctggt       | CG               | tgcggt       | DG               | gatggc       | EG               | gagggc       | FG               | ttcgga       |
| AH               | gcccac       | CH               | tgccat       | DH               | gatcac       | EH               | gagcac       | FH               | tttcat       |
| AI               | gccatc       | CI               | tgtatc       | DI               | gacatt       | EI               | gaaatc       | FI               | tttatt       |
| AK               | gccaag       | CK               | tgcaag       | DK               | gacaag       | EK               | gaaaag       | FK               | ttcaag       |
| AL               | gcccta       | CL               | tgcttg       | DL               | gacctc       | EL               | gaactc       | FL               | ttcttg       |
| AM               | gccatg       | CM               | tgtatg       | DM               | gatatg       | EM               | gaaatg       | FM               | tttatg       |
| AN               | gccaac       | CN               | tgcaac       | DN               | gacaac       | EN               | gagaac       | FN               | ttcaat       |
| AP               | gcccc        | CP               | tgccca       | DP               | gacccc       | EP               | gagcct       | FP               | ttcccc       |
| AQ               | gcccga       | CQ               | tgccag       | DQ               | gaccga       | EQ               | gagcag       | FQ               | ttccga       |
| AR               | gctcgt       | CR               | tgccgg       | DR               | gatcgt       | ER               | gaacga       | FR               | ttccgg       |
| AS               | gcatcg       | CS               | tgctcg       | DS               | gactct       | ES               | gaatcc       | FS               | ttctct       |
| AT               | gccacc       | CT               | tgactc       | DT               | gacact       | ET               | gaaacc       | FT               | ttcacc       |
| AV               | gctgtc       | CV               | tgtgtc       | DV               | gacggt       | EV               | gaggtc       | FV               | ttcgtg       |
| AW               | gcgtgg       | CW               | tgttgg       | DW               | gactgg       | EW               | gaatgg       | FW               | ttctgg       |
| AY               | gcgtat       | CY               | tgctat       | DY               | gactat       | EY               | gaatac       | FY               | ttctat       |
| GA               | ggggca       | HA               | cacgcc       | IA               | atcgcg       | KA               | aaggcc       | LA               | ttggct       |
| GC               | ggctgc       | HC               | cactgc       | IC               | atctgt       | KC               | aagtgc       | LC               | ctatgc       |
| GD               | ggcgac       | HD               | cacgac       | ID               | atcgac       | KD               | aaggat       | LD               | ctcgac       |
| GE               | ggcgag       | HE               | cacgag       | IE               | atcgaa       | KE               | aaggaa       | LE               | ttggag       |
| GF               | gggttc       | HF               | cacttc       | IF               | atcttt       | KF               | aagttt       | LF               | ctgttc       |
| GG               | ggtggc       | HG               | cacggc       | IG               | atcggc       | KG               | aagggc       | LG               | ctcggc       |
| GH               | ggccac       | HH               | caccac       | IH               | atccat       | KH               | aaacat       | LH               | cttcac       |
| GI               | ggcatc       | HI               | cacatc       | II               | attatc       | KI               | aagatc       | LI               | ctcatt       |
| GK               | ggcaag       | HK               | cacaag       | IK               | attaag       | KK               | aaaaag       | LK               | ctcaaa       |
| GL               | ggcctc       | HL               | cacctg       | IL               | atcttg       | KL               | aagctc       | LL               | cttctc       |
| GM               | ggtatg       | HM               | cacatg       | IM               | atcatg       | KM               | aagatg       | LM               | ctcatg       |
| GN               | ggcaac       | HN               | cacaac       | IN               | atcaac       | KN               | aagaac       | LN               | ctcaat       |
| GP               | ggtccc       | HP               | cacccc       | IP               | atcccc       | KP               | aagcct       | LP               | cttccc       |
| GQ               | ggccag       | HQ               | caccag       | IQ               | atccag       | KQ               | aagcag       | LQ               | ctccag       |
| GR               | ggtcgc       | HR               | cacaga       | IR               | atccgt       | KR               | aaacgc       | LR               | ttgcgc       |
| GS               | ggctcc       | HS               | cacagc       | IS               | atcagc       | KS               | aagtct       | LS               | ctctct       |
| GT               | ggcacc       | HT               | cacacc       | IT               | ataacc       | KT               | aagacc       | LT               | ttgact       |
| GV               | ggcgtc       | HV               | cacgtg       | IV               | atcgtc       | KV               | aaggtt       | LV               | ctgggt       |
| GW               | ggttgg       | HW               | cactgg       | IW               | atctgg       | KW               | aagtgg       | LW               | atttgg       |
| GY               | ggctac       | HY               | cactac       | IY               | atctat       | KY               | aagtac       | LY               | ctctac       |

| Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) | Amino Acid (XX') | DNA (NNNNNN) |
|------------------|--------------|------------------|--------------|------------------|--------------|------------------|--------------|------------------|--------------|
| MA               | atggcg       | NA               | aacgcc       | PA               | ccagcc       | QA               | caggct       | RA               | agagct       |
| MC               | atgtgc       | NC               | aactgc       | PC               | ccctgt       | QC               | cagtgc       | RC               | cgctgc       |
| MD               | atggat       | ND               | aacgac       | PD               | cctgac       | QD               | caggac       | RD               | agggac       |
| ME               | atggaa       | NE               | aacgag       | PE               | ccagaa       | QE               | caggag       | RE               | agagag       |
| MF               | atgttt       | NF               | aacttc       | PF               | ccattt       | QF               | caattt       | RF               | aggttt       |
| MG               | atgggc       | NG               | aacggc       | PG               | ccgggg       | QG               | cagggt       | RG               | cgcggc       |
| MH               | atgcac       | NH               | aaccac       | PH               | ccgcat       | QH               | caacat       | RH               | agacac       |
| MI               | atgata       | NI               | aacatc       | PI               | ccaatc       | QI               | cagatc       | RI               | agaatc       |
| MK               | atgaag       | NK               | aacaag       | PK               | cctaag       | QK               | cagaag       | RK               | cgcaag       |
| ML               | atgtta       | NL               | aacctg       | PL               | ccgctc       | QL               | caactt       | RL               | cgctc        |
| MM               | atgatg       | NM               | aacatg       | PM               | cccatg       | QM               | cagatg       | RM               | agaatg       |
| MN               | atgaac       | NN               | aacaac       | PN               | cctaac       | QN               | cagaac       | RN               | cgtaac       |
| MP               | atgccc       | NP               | aacccc       | PP               | ccccc        | QP               | caacct       | RP               | cgccc        |
| MQ               | atgcag       | NQ               | aaccag       | PQ               | ccgcag       | QQ               | cagcaa       | RQ               | agacag       |
| MR               | atgcgg       | NR               | aacaga       | PR               | cctcgt       | QR               | caacgt       | RR               | cgcgga       |
| MS               | atgtcc       | NS               | aacagc       | PS               | ccttcg       | QS               | caatcc       | RS               | cgctcc       |
| MT               | atgact       | NT               | aacacc       | PT               | ccgacg       | QT               | cagact       | RT               | cgcacc       |
| MV               | atggtc       | NV               | aacgtg       | PV               | cctgtg       | QV               | caggtg       | RV               | cggtt        |
| MW               | atgtgg       | NW               | aactgg       | PW               | ccgtgg       | QW               | cagtgg       | RW               | cgctgg       |
| MY               | atgtat       | NY               | aactac       | PY               | ccatac       | QY               | cagtat       | RY               | agatac       |
| SA               | tccgcc       | TA               | actgcc       | VA               | gtcgcc       | WA               | tgggcc       | YA               | tacgcg       |
| SC               | tcctgc       | TC               | acctgc       | VC               | gtatgc       | WC               | tggcgc       | YC               | tactgt       |
| SD               | tcggat       | TD               | accgac       | VD               | gtggac       | WD               | tgggat       | YD               | tacgat       |
| SE               | tccgaa       | TE               | accgag       | VE               | gttgag       | WE               | tgggaa       | YE               | tatgaa       |
| SF               | agcttt       | TF               | acgttc       | VF               | gtgttc       | WF               | tggttt       | YF               | tacttc       |
| SG               | tcgggg       | TG               | actggc       | VG               | gtcgga       | WG               | tggggc       | YG               | tatggc       |
| SH               | tcacac       | TH               | aaccac       | VH               | gtgcat       | WH               | tggcat       | YH               | tatcac       |
| SI               | agcatc       | TI               | accatc       | VI               | gtcatt       | WI               | tggatc       | YI               | tacatt       |
| SK               | tccaag       | TK               | accaag       | VK               | gtcaag       | WK               | tggaag       | YK               | tacaag       |
| SL               | tcctc        | TL               | actctc       | VL               | gtcttg       | WL               | tggttg       | YL               | tactta       |
| SM               | agcatg       | TM               | acatg        | VM               | gtcatg       | WM               | tggatg       | YM               | tatatg       |
| SN               | tccaac       | TN               | accaac       | VN               | gtgaac       | WN               | tggaac       | YN               | tacaac       |
| SP               | tccccg       | TP               | acacca       | VP               | gtccct       | WP               | tggccg       | YP               | tatccg       |
| SQ               | tcgcaa       | TQ               | acccaa       | VQ               | gtgcag       | WQ               | tggcag       | YQ               | taccag       |
| SR               | tcgagg       | TR               | accgac       | VR               | gtgcgg       | WR               | tggcgc       | YR               | tatcga       |
| SS               | tctagc       | TS               | acatcc       | VS               | gtgtcc       | WS               | tggcgc       | YS               | tattcc       |
| ST               | tctacc       | TT               | actacg       | VT               | gtgacc       | WT               | tggact       | YT               | tatact       |
| SV               | tccgtg       | TV               | acggtc       | VV               | gttgtt       | WV               | tgggtt       | YV               | tacgtt       |
| SW               | tcctgg       | TW               | acttgg       | VW               | gtgtgg       | WW               | tgggtg       | YW               | tactgg       |
| SY               | agttac       | TY               | acttac       | VY               | gtctac       | WY               | tgggat       | YY               | tactat       |

**Supplementary information, Table S3.** Base binding specificity of previously reported RVDs obtained from TALE-(XX')<sub>3</sub> library screening.

| Category | RVDs | A    | T  | C    | G    |
|----------|------|------|----|------|------|
| NX'      | NA   | -    | ++ | +    | +    |
|          | NC   | +    | -  | +    | -    |
|          | ND   | -    | -  | +    | -    |
|          | NG   | -    | +  | -    | -    |
|          | NH   | -    | -  | -    | ++++ |
|          | NI   | ++   | -  | -    | -    |
|          | NK   | -    | -  | -    | +    |
|          | NN   | ++   | -  | -    | ++++ |
|          | NP   | -    | -  | -    | -    |
|          | NQ   | -    | -  | +    | -    |
|          | NS   | ++++ | -  | +    | ++++ |
|          | NT   | ++++ | -  | +    | ++++ |
|          | NV   | ++   | -  | +    | +    |
| HX'      | HA   | -    | +  | +    | -    |
|          | HD   | -    | -  | ++++ | -    |
|          | HG   | -    | -  | -    | -    |
|          | HH   | -    | -  | -    | ++++ |
|          | HI   | -    | -  | -    | -    |
|          | HN   | ++   | -  | -    | ++++ |
| IX'      | IG   | -    | -  | -    | -    |
|          | IS   | -    | -  | -    | -    |
| SX'      | SH   | -    | -  | -    | -    |
|          | SN   | -    | +  | -    | +    |
|          | SS   | +    | +  | -    | +    |
| YX'      | YG   | -    | -  | -    | -    |

Note:

The ranges of fold induction of EGFP reporter for RVDs in TALE-(XX')<sub>3</sub> are indicated as follows:

- < 6  
 + 6 - 12  
 ++ 12 - 18  
 +++ 18 - 24  
 ++++ ≥ 24

## Supplementary information, Table S4. Efficiencies of TALENs-mediated indels with novel RVDs.

| Targeted Gene     | Type of TALENs      | TALEN <sup>L</sup>  | TALEN <sup>R</sup>   | Mean Indels (%)<br>± s.d. (n = 3) |
|-------------------|---------------------|---|--|-----------------------------------|
| LRP1 <sup>a</sup> | TALEN <sub>SN</sub> | NI NI <b>NN</b> NI HD NG NG <b>NN</b> HD NI <b>NN</b> HD HD HD HD                             | NI HD NI <b>NN</b> <b>NN</b> NG NG NI NG NG NG <b>NN</b> NI NG HD                      | 77.1 ± 4.3                        |
|                   | TALEN <sub>NH</sub> | NI NI <b>NH</b> NI HD NG NG <b>NH</b> HD NI <b>NH</b> HD HD HD HD                             | NI HD NI <b>NH</b> <b>NH</b> NG NG NI NG NG NG <b>NH</b> NI NG HD                      | 50.7 ± 4.9                        |
|                   | TALEN <sub>KN</sub> | NI NI <b>KN</b> NI HD NG NG <b>KN</b> HD NI <b>KN</b> HD HD HD HD                             | NI HD NI <b>KN</b> <b>KN</b> NG NG NI NG NG NG <b>KN</b> NI NG HD                      | 78.1 ± 3.4                        |
| DKK1 <sup>b</sup> | TALEN <sub>SN</sub> | NG HD HD NI NI HD <b>NN</b> HD NG NI NG HD NI NI <b>NN</b>                                    | <b>NN</b> HD HD HD HD <b>NN</b> HD NI <b>NN</b> HD <b>NN</b> HD HD <b>NN</b> HD        | 54.7 ± 4.8                        |
|                   | TALEN <sub>NH</sub> | NG HD HD NI NI HD <b>NH</b> HD NG NI NG HD NI NI <b>NH</b>                                    | <b>NH</b> HD HD HD HD <b>NH</b> HD NI <b>NH</b> HD <b>NH</b> HD HD <b>NH</b> HD        | 39.4 ± 7.8                        |
|                   | TALEN <sub>KN</sub> | NG HD HD NI NI HD <b>KN</b> HD NG NI NG HD NI NI <b>KN</b>                                    | <b>KN</b> HD HD HD HD <b>KN</b> HD NI <b>KN</b> HD <b>KN</b> HD HD <b>KN</b> HD        | 49.9 ± 2.7                        |
| LRP1 <sup>a</sup> | TALEN <sub>NG</sub> | NI NI NN NI HD <b>NG</b> <b>NG</b> NN HD NI NN HD HD HD HD                                    | NI HD NI NN NN <b>NG</b> <b>NG</b> NI <b>NG</b> <b>NG</b> <b>NG</b> NN NI <b>NG</b> HD | 77.1 ± 4.3                        |
|                   | TALEN <sub>RG</sub> | NI NI NN NI HD <b>RG</b> <b>RG</b> NN HD NI NN HD HD HD HD                                    | NI HD NI NN NN <b>RG</b> <b>RG</b> NI <b>RG</b> <b>RG</b> <b>RG</b> NN NI <b>RG</b> HD | 40.2 ± 10.3                       |
| ATG5 <sup>c</sup> | TALEN <sub>NG</sub> | <b>NN</b> <b>NG</b> <b>NG</b> <b>NG</b> HD NI HD NN HD <b>NG</b> NI <b>NG</b> NI <b>NG</b> HD | NI <b>NG</b> NN NN <b>NG</b> <b>NG</b> HD <b>NG</b> NN HD <b>NG</b> <b>NG</b> HD HD HD | 11.5 ± 0.6                        |
|                   | TALEN <sub>RG</sub> | <b>NN</b> <b>RG</b> <b>RG</b> <b>RG</b> HD NI HD NN HD <b>RG</b> NI <b>RG</b> NI <b>RG</b> HD | NI <b>RG</b> NN NN <b>RG</b> <b>RG</b> HD <b>RG</b> NN HD <b>RG</b> <b>RG</b> HD HD HD | 6.5 ± 0.5                         |

<sup>a</sup> Sequence for TALENs targeting is 5'-TAAGACTTGCAGCCCCAAGCAGTTTGCCTGCAGAGATCAAATAACCTGTA-3', and TALEN<sup>L</sup> and TALEN<sup>R</sup> targeted regions are labeled in grey shades. Sequences of primers used to amplify the TALENs targeting region for the assessment of indels are: 5'-AGCCTAGAGGACTTGGAGGG-3' and 5'-ACCTGGGCAATCATCAAG-3'.

<sup>b</sup> Sequence for TALENs targeting is 5'-TTCCAACGCTATCAAGAACCTGCCCCACCGCTGGGCGGCGCTGCGGGGCA-3', and TALEN<sup>L</sup> and TALEN<sup>R</sup> targeted regions are labeled in grey shades. Sequences of primers used to amplify the TALENs targeting region for the assessment of indels are: 5'-GACCCAGGCTTGCAAAGTGA-3' and 5'-CAGGCGAGACAGATTTGCAC-3'.

<sup>c</sup> Sequence for TALENs targeting is 5'-TGTTTCACGCTATATCAGGATGAGATAACTGAAAGGGAAGCAGAACCATA-3', and TALEN<sup>L</sup> and TALEN<sup>R</sup> targeted regions are labeled in grey shades. Sequences of primers used to amplify the TALENs targeting region for the assessment of indels are: 5'-TAGCAGGTTTCATTCATCGTTGCAAGGA-3' and 5'-ATGTAAGGAAAACAAAGTCCAGAACGC-3'.

**Supplementary information, Table S5.**

RVDs targeting multiple bases.

## RVDs targeting two bases

| Category | RVDs | A  | T  | C  | G    |
|----------|------|----|----|----|------|
| A/T      | SG   | +  | +  | -  | -    |
| A/C      | HC   | ++ | -  | ++ | -    |
|          | KC   | ++ | -  | ++ | -    |
|          | NC   | +  | -  | +  | -    |
|          | KS   | ++ | -  | +  | -    |
| A/G      | HN   | ++ | -  | -  | ++++ |
|          | NN   | ++ | -  | -  | ++++ |
| T/C      | HA   | -  | +  | +  | -    |
|          | KA   | -  | ++ | +  | -    |
|          | KP   | -  | ++ | +  | -    |
|          | FS   | -  | +  | +  | -    |
| T/G      | KF   | -  | +  | -  | +    |
|          | SN   | -  | +  | -  | +    |
|          | RQ   | -  | +  | -  | ++++ |
|          | GR   | -  | +  | -  | +    |
| C/G      | KR   | -  | +  | -  | +++  |
|          | HQ   | -  | -  | +  | +++  |
|          | KQ   | -  | -  | +  | ++   |
|          | QR   | -  | -  | +  | ++   |
|          | YR   | -  | -  | +  | ++++ |

## RVDs targeting four bases

| Category | RVDs | A   | T  | C   | G   |
|----------|------|-----|----|-----|-----|
| A/T/C/G  | RH   | +   | +  | +   | +++ |
|          | LQ   | +   | ++ | +   | +   |
|          | VR   | +   | +  | +   | ++  |
|          | CS   | +   | +  | +   | +   |
|          | RT   | +   | +  | +   | ++  |
|          | RV   | +++ | ++ | +++ | ++  |

Note:

The ranges of fold induction of EGFP reporter for RVDs in TALE-(XX')<sub>3</sub> are indicated as follows:

|      |         |
|------|---------|
| -    | < 6     |
| +    | 6 - 12  |
| ++   | 12 - 18 |
| +++  | 18 - 24 |
| ++++ | ≥ 24    |

## RVDs targeting three bases

| Category | RVDs | A    | T  | C   | G    |
|----------|------|------|----|-----|------|
| A/T/C    | RL   | +    | ++ | +   | -    |
|          | MS   | +    | +  | +   | -    |
|          | TT   | +    | +  | +   | -    |
| A/T/G    | TC   | ++   | +  | -   | +    |
|          | MH   | +    | +  | -   | +    |
|          | LR   | +    | +  | -   | ++   |
|          | SS   | +    | +  | -   | +    |
| A/C/G    | RC   | +    | -  | ++  | +    |
|          | CR   | +    | -  | +   | +++  |
|          | HS   | ++   | -  | ++  | ++++ |
|          | NS   | ++++ | -  | +   | ++++ |
|          | HT   | ++   | -  | +   | +    |
|          | KT   | +    | -  | +   | ++   |
|          | NT   | ++++ | -  | +   | ++++ |
|          | HV   | ++   | -  | +   | +    |
|          | KV   | ++++ | -  | +++ | +++  |
|          | NV   | ++   | -  | +   | +    |
| T/C/G    | FQ   | -    | +  | +   | +    |
|          | NA   | -    | ++ | +   | +    |
|          | RN   | -    | +  | +   | ++++ |

## SUPPLEMENTARY REFERENCES

1. Zhang, F. et al. Efficient construction of sequence-specific TAL effectors for modulating mammalian transcription. *Nature biotechnology* **29**, 149-153 (2011).
2. Yang, J. et al. ULtiMATE system for rapid assembly of customized TAL effectors. *PLoS One* **8**, e75649 (2013).
3. Boussif, O. et al. A versatile vector for gene and oligonucleotide transfer into cells in culture and in vivo: polyethylenimine. *Proc Natl Acad Sci U S A* **92**, 7297-7301 (1995).
4. Mussolino, C. et al. A novel TALE nuclease scaffold enables high genome editing activity in combination with low toxicity. *Nucleic acids research* (2011).



## Supplementary information, Data S1 Sequences and material and methods

## TALE-Ctrl

1 GTCGACGGATCGGGAGATCTCCCATCCCCTATGGTGCACCTCTCAGTACAATCTGCTCTG  
 61 ATGCCGCATAGTTAAGCCAGTATCTGCTCCCCTGCTTGTGTGTTGGAGGTCGCTGAGTAGT  
 121 GCGCGAGCAAAATTTAAGCTACAACAAGGCAAGGCTTGACCGACAATTGCATGAAGAATC  
 181 TGCTTAGGGTTAGGCGTTTTGCGCTGCTTCGCGATGTACGGGCCAGATATACGCGTTGAC  
 241 ATTGATTATTGACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCAT  
 301 ATATGGAGTTCCGCGTTACATAACTTACGGTAAATGGCCCCGCTGGCTGACCGCCCAACG  
 361 ACCCCCGCCCATTTGACGTCAATAATGACGTATGTTCCCATAGTAACGCCAATAGGGACTT  
 421 TCCATTGACGTCAATGGGTGGAGTATTTACGGTAAACTGCCCACTTGGCAGTACATCAAG  
 481 TGTATCATATGCCAAGTACGCCCCCTATTGACGTCAATGACGGTAAATGGCCCCGCTGGC  
 541 ATTATGCCCAGTACATGACCTTATGGGACTTTTCCCTACTTGGCAGTACATCTACGTATTAG  
 601 TCATCGCTATTACCATGGTGATGCGGTTTTTGGCAGTACATCAATGGGCGTGGATAGCGGT  
 661 TTGACTCACGGGGATTTCCAAGTCTCCACCCCATTTGACGTCAATGGGAGTTTTGTTTTGGC  
 721 ACCAAAATCAACGGGACTTTTCCAAAATGTCGTAACAACCTCCGCCCCATTGACGCAAAATGG  
 781 GCGGTAGGCGTGTACGGTGGGAGGTCATATAAAGCAGCGCGTTTTGCCTGTACTGGGTCT  
 841 CTCTGGTTAGACCAGATCTGAGCCTGGGAGCTCTCTGGCTAACTAGGGAACCCACTGCTT  
 901 AAGCTCAATAAAGCTTGCCCTGAGTGCCTCAAGTAGTGTGTGCCGCTCTGTTGTGTGAC  
 961 TCTGGTAACTAGAGATCCCTCAGACCCTTTTAGTCAGTGTGAAAAATCTCTAGCAGTGGC  
 1021 GCCCGAACAGGGACTTGAAAGCGAAAGGGAAACCAGAGGAGCTCTCTCGACGCAGGACTC  
 1081 GGCTTGCATGAAGCGCGCACGGCAAGAGGCGAGGGGCGCGACTGGTGAGTACGCCAAAAA  
 1141 TTTTACTAGCGGAGGCTAGAAGGAGAGAGATGGGTGCGAGAGCGTCAGTATTAAGCGGG  
 1201 GGAGAAATTAGATCGCGATGGGAAAAAATTCGGTTAAGGCCAGGGGGAAAGAAAAAATATA  
 1261 AATTAAAACATATAGTATGGGCAAGCAGGGAGCTAGAACGATTCGCAGTTAATCCTGGCC  
 1321 TGTTAGAAACATCAGAAGGCTGTAGACAAATACTGGGACAGCTACAACCATCCCTTCAGA  
 1381 CAGGATCAGAAGAACTTAGATCATTATATAATACAGTAGCAACCCTCTATTGTGTGCATC  
 1441 AAAGGATAGAGATAAAAGACACCAAGGAAGCTTTAGACAAGATAGAGGAAGAGCAAAAACA  
 1501 AAAGTAAGACCACCGCACAGCAAGCGGCCGGCCGCGCTGATCTTCAGACCTGGAGGAGGA  
 1561 GATATGAGGAGACAAATGGAGAAGTGAATATATAAAATATAAAGTAGTAAAAATGAACCA  
 1621 TTAGGATGAGCACCCCAAGGCAAGGAGAAGAGAGGTTGTCAGAGAGAAAAAAGACAGTG  
 1681 GGAATAGGAGCTTTTGTTCCTTGGGTTCTTGGGAGCAGCAGGAAGCACTATGGGCGCAGCG  
 1741 TCAATGACGCTGACGGTACAGGCCAGACAATTAATTGTCTGGTATAGTGCAGCAGCAGAAC  
 1801 AATTTGCTGAGGGCTATTGAGGCGCAACAGCATCTGTTGCAACTCACAGTCTGGGGCATC  
 1861 AAGCAGCTCCAGGCAAGAATCCTGGCTGTGGAAAGATACCTAAAGGATCAACAGCTCCTG  
 1921 GGGATTTGGGGTTGCTCTGGAAAACCTATTTGCACCCTGCTGTGCCTTGGAAATGCTAGT  
 1981 TGGAGTAATAAATCTCTGGAACAGATTTGGAATCACACGACCTGGATGGAGTGGGACAGA  
 2041 GAAATTAACAATTACACAAGCTTAATACACTCCTTAATTAAGAATCGCAAAACCAGCAA  
 2101 GAAAAGAAATGAACAAGAATTATTGGAATTAGATAAATGGGCAAGTTTTGTGGAATTGGTTT  
 2161 AACATAACAAATTTGGCTGTGGTATATAAAATTAATCATAATGATAGTAGGAGGCTTGGTA  
 2221 GGTTTAAGAATAGTTTTTGTGTACTTTCTATAGTGAATAGAGTTAGGCAGGGATATTCA  
 2281 CCATTATCGTTTTAGACCCACCTCCCAACCCCGAGGGGACCCGACAGGCCCGAAGGAATA  
 2341 GAAGAAGAAGGTTGAGAGAGAGACAGACAGATCCATTTCGATTGAGTGAACGGATCGGCA  
 2401 CTGCGTGCGCCAATTTCTGCAGACAAATGGCAGTATTCATCCACAATTTTAAAGAAAAGG  
 2461 GGGGATTTGGGGGTACAGTGCAGGGGAAAGAATAGTAGACATAATAGCAACAGACATACA  
 2521 AACTAAAGAAATTACAAAAACAAATTACAAAAATTCAAAATTTTCGGGTTTATTACAGGGA  
 2581 CAGCAGAGATCCAGTTTGGTTAGTACCGGGCCCGCTCTAGACGATGTACGGGCCAGATAT  
 2641 ACGCGTTGACATTGATTATTGACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTT  
 2701 CATAGCCCATATATGGAGTTCCGCGTTACATAACTTACGGTAAATGGCCCCGCTGGCTGA  
 2761 CCGCCCAACGACCCCGCCCATTTGACGTCAATAATGACGTATGTTCCCATAGTAACGCCA  
 2821 ATAGGGACTTTCCATTGACGTCAATGGGTGGAGTATTTACGGTAAACTGCCCACTTGGCA  
 2881 GTACATCAAGTGTATCATATGCCAAGTACGCCCCCTATTGACGTCAATGACGGTAAATGG  
 2941 CCCGCCGTCATTATGCCCAGTACATGACCTTATGGGACTTTTCCCTACTTGGCAGTACATC  
 3001 TACGTATTAGTCATCGCTATTACCATTGGTGTAGTGCCTTTTGGCAGTACATCAATGGGCGT  
 3061 GGATAGCGGTTTGGATCACGGGATTTCCAAGTCTCCACCCATTGACGTCAATGGGAGT  
 3121 TTGTTTTGGCACCAAAATCAACGGGACTTTTCCAAAATGTCGTAACAACCTCCGCCCCATTG  
 3181 ACGCAAATGGGCGGTAGGCGTGTACGGTGGGAGGTCATATAAAGCAGAGCTCTCTGGCTA  
 3241 ACTAGAGAACCCACTGCTTACTGGCTTATCGAAATTAATACGACTCACTATAGGGAGACC  
 3301 CAAGCTGGCTAGCGAAGTTCCATTTCTCTAGAAAGTATAGGAACTTCATGAGGACCAGGC  
 3361 TGCCATCTCCCCCTGCCCTTCCCCCGCATTTAGCGCTGGGAGCTTTAGCGACCTGCTTA  
 3421 GGCAGTTCGACCCAGCTTGTTCACACCAGCCTGTTTTGACAGCCTGCCTCCCTTCGGAG

3481 CGCACCACACCGAAGCCGCCACCGGCGAGTGGGACGAGGTGCAAAGCGGCCTGAGGGCAG  
3541 CGGACGCTCCTCCGCCAACCATGAGGGTGGCAGTGACAGCAGCTAGGCCCCCTCGGGCAA  
3601 AACCTGCACCCAGGAGAAGGGCTGCCAACCAGCGACGCGAGTCCAGCCGCACAGGTGG  
3661 ACCTCAGGACGCTGGGCTACAGCCAGCAACAGCAAGAGAAGATCAAGCCCAAAGTAAGGA  
3721 GCACCGTGGCCAGCACCACGAGGCCCTGGTGGGTACGGCTTACCCACGCGCATATCG  
3781 TTGCTCTGAGCCAACATCCCGCAGCTCTGGGTACCCTTGCAGTGAAGTATCAGGACATGA  
3841 TCGCGGCACCTGCCCTGAAGCTACACACGAAGCCATAGTGGGCGTTGGCAAGCAGTGGAGCG  
3901 GTGCCAGAGCGCTTGAGGCACCTGTTGACGGTGGCTGGCGAGCTGAGGGGACCGCCACTGC  
3961 AACTGGACACCGGCCAACCTGCTGAAGATCGCCAAGAGGGGAGGCGTGACGGCGGTGGAGG  
4021 CCGTGCATGCCCTGGAGGAATGCCCTGACCGGCGCGCCCTGAACCTGACACCAGAGCAAG  
4081 TAGTGGCTATTGCAAGTAACATCGGTGGCAAACAAGCGCTGGAGACCGTGCAGAGGCTCC  
4141 TTCCGGTGCCTTGCCAAGCACACGGTTTGGACCCCGAACAGGTTGTAGCCATAGCTTCTA  
4201 ACAACGGAGGTAAGCAGGCACCTGGAAACCGTGCAGCGCCTGCTCCCAGTACTGTGTCAGG  
4261 CTCATGGGCTCACTCCGGAACAGGTGGTGCGAATCGCGAGCAACGGCGGGCGCAAGCAAG  
4321 CCTTGGAGACAGTCCAAAGACTTTTTGCCCTGTCTTTGTGTCAGGCGCATGGCCTTACGCCTG  
4381 AGCAAGTCGTTGCGATCGCCTCCACGACGGCGGAAAACAGGCTTTGGAAACCGTGCAGC  
4441 GGTGCTGCCGTTTTTGTGCCAAGCCCACGGATTGACCCCGAACAGGTTGTAGCCATAG  
4501 CTTCTAACATCGGAGGTAAGCAGGCACCTGGAAACCGTGCAGCGCCTGCTCCCAGTACTGT  
4561 GTCAGGCTCATGGGCTCACTCCGGAACAGGTGGTGCGAATCGCGAGCAACGGCGGGCGGCA  
4621 AGCAAGCCCTTGAGACAGTCCAAAGACTTTTTGCCCTGTCTTTGTGTCAGGCGCATGGCCTTA  
4681 CGCTGAGCAAGTCGTTGCGATCGCCTCCACGACGGCGGAAAACAGGCTTTGGAAACCG  
4741 TGCAGCGGTTGCTGCCGTTTTTGTGCCAAGCCCACGGATTGACCCCGAACAGGTTGTAG  
4801 CCATAGCTTCTAACAAACGGAGGTAAGCAGGCACCTGGAAACCGTGCAGCGCCTGCTCCAG  
4861 TACTGTGTCAGGCTCATGGGCTTACGCCTGAGCAAGTCGTTGCGATCGCCTCCACGACG  
4921 GCGGAAAACAGGCTTTGGAAACCGTGCAGCGGTTGCTGCCGTTTTTGTGCCAAGCCCACG  
4981 GACTCACTCCGGAACAGGTGGTGCGAATCGCGAGCAACGGCGGGCGCAAGCAAGCCCTTG  
5041 AGACAGTCCAAAGACTTTTTGCCCTGTCTTTGTGTCAGGCGCATGGCCTGACACCAGAGCAAG  
5101 TAGTGGCTATTGCAAGTAACATCGGTGGCAAACAAGCGCTGGAGACCGTGCAGAGGCTCC  
5161 TTCCGGTGCCTTGCCAAGCACACGGTTTTGACCCCGAACAGGTTGTAGCCATAGCTTCTA  
5221 ACGGCGGAGGTAAGCAGGCACCTGGAAACCGTGCAGCGCCTGCTCCCAGTACTGTGTCAGG  
5281 CTCATGGGCTCACTCCGGAACAGGTGGTGCGAATCGCGAGCAACGGCGGGCGGCAAGCAAG  
5341 CCCTTGAGACAGTCCAAAGACTTTTTGCCCTGTCTTTGTGTCAGGCGCATGGCCTGACACCAG  
5401 AGCAAGTAGTGGCTATTGCAAGTAACATCGGTGGCAAACAAGCGCTGGAGACCGTGCAGACA  
5461 GGCTCCTTCCGGTGCCTTGCCAAGCACACGGTCTCACTCCGGAACAGGTGGTGCGAATCG  
5521 CGAGCCACGACGGCGGCAAGCAAGCCCTTGAGACAGTCCAAAGACTTTTTGCCCTGTCTTTT  
5581 GTCAGGCGCATGGCCTTACGCCTGAGCAAGTCGTTGCGATCGCCTCCACGACGGCGGAA  
5641 AACAGGCTTTGGAAACCGTGCAGCGGTTGCTGCCGTTTTTGTGCCAAGCCCACGGACTGA  
5701 CACCAGAGCAAGTAGTGGCTATTGCAAGTAACATCGGTGGCAAACAAGCGCTGGAGAGCA  
5761 TCGTGGCCCAGCTGTCTCGGCCCGACCTGCCCTCGCCGCTCTGACCAACGACCACCTGG  
5821 TGGCCCTGGCTTGCCCTCGGGGGCAGGCCAGCTCTTGACGCCGTGAAGAAGGGCCTTCCCTC  
5881 ACGCCCCAGCCCTGATCAAGCGGACCAACAGAAGGATTCCCGAGAGGACATCACATCGAG  
5941 TGGCAGATCACGCGCAAGTGGTCCGCGTGTCTCGGATTCCTCCAGTGTCACTCCCACCCCG  
6001 CACAAGCGTTCGATGACGCCATGACTCAATTTGGTATGTGCGAGACACGGACTGTGTCAGC  
6061 TCTTTCGTAGAGTCGGTGTACAGAACTCGAGGCCCGCTCGGGCACACTGCCTCCGCCCT  
6121 CCCAGCGTGGGACAGGATTCCTCAAGCGAGCGGTATGAAACGCGCGAAGCCTTCCACCTA  
6181 CGTCAACTCAGACACCTGACCAGGCGAGCCTTCATGCGTTTCGAGACTCGCTGGAGAGGG  
6241 ATTTGGACGCGCCCTCGCCCATGCATGAAGGGGACCAAACCTCGCGCGTCACTAGCCCCA  
6301 AGAAGAAGAGAAAGGTGGAGGCCAGCGGTTCCGGACGGGCTGACGCATTGGACGATTTTG  
6361 ATCTGGATATGCTGGGAAGTGACGCCCTCGATGATTTTGACCTTGACATGCTTGGTTCGG  
6421 ATGCCCTTGATGACTTTGACCTCGACATGCTCGGCAGTGACGCCCTTGATGATTTTCGACC  
6481 TGGACATGCTGATTAACCTTAGAGGCAGTGGAGAGGGCAGAGGAAGTCTGCTAACATGCG  
6541 GTGACGTCGAGGAGAATCCTGGCCAGTGAGCAAGGGCGAGGAGGATAACATGGCCATCA  
6601 TCAAGGAGTTCATGCGCTTCAAGGTGCACATGGAGGGCTCCGTGAACGGCCACGAGTTCCG  
6661 AGATCGAGGGCGAGGGCGAGGGCCGCCCTTACGAGGGCACCCAGACCGCCAAGCTGAAGG  
6721 TGACCAAGGGTGGCCCCCTGCCCTTCGCCCTGGGACATCCTGTCCCCTCAGTTCATGTACG  
6781 GCTCCAAGGCCTACGTGAAGCACCCCGCCGACATCCCCGACTACTTGAAGCTGTCTTCC  
6841 CCGAGGGCTTCAAGTGGGAGCGCTGATGAACCTTCGAGGACGGCGGCTGGTGCACCGTGA  
6901 CCCAGACTCCCTCCCTGCGAGGACGGGAGTTCATCTACAAGGTGAAGCTGCGGACCCACA  
6961 ACTTCCCCTCCGACGGCCCCGTAATGCAGAAGAAGACCATGGGCTGGGAGGCCTCCTCCG  
7021 AGCGGATGTACCCCGAGGACGGCGCCCTGAAGGGCGAGATCAAGCAGAGGCTGAAGCTGA  
7081 AGGACGGCGGCCACTACGACGCTGAGGTCAAGACCACCTACAAGGCCAAGAAGCCCGTGC  
7141 AGCTGCCCGGCGCCTACAACGTCAACATCAAGTTGGACATCACCTCCCACAACGAGGACT  
7201 ACACCATCGTGGAAACAGTACGAACGCGCCGAGGGCCGCCACTCCACCGGCGGCATGGACG

7261 AGCTGTACAAGTAACATGTTTAAAGGGTCCGGTCCACTAGGTACAATTCGATATCAAGC  
 7321 TTATCGATAATCAACCTCTGGATTACAAAATTTGTGAAAGATTGACTGGTATTCTTAACT  
 7381 ATGTTGCTCCTTTTACGCTATGTGGATACGCTGCTTTAATGCCTTTGTATCATGCTATTG  
 7441 CTCCCCGTATGGCTTTCATTTTCTCCTCCTTGTATAAATCCTGGTTGCTGTCTCTTTATG  
 7501 AGGAGTTGTGGCCCGTTGTCCAGGCAACGTGGCGTGGTGTGCACTGTGTTTGTGCTGACGCAA  
 7561 CCCCCACTGGTTGGGGCATTGCCACCACCTGTCAGCTCCTTTCCGGGACTTTTCGCTTTCC  
 7621 CCTCCCTATTTGCCACGGCGGAACCTATCGCCGCCTGCCTTGCCCGCTGCTGGACAGGGG  
 7681 CTCGGCTGTTGGGCACGACAATCCCGTGGTGTGTGCGGGGAAATCATCGTCTCTTCCCTT  
 7741 GGCTGCTCGCCTGTGTTGCCACCTGGATTCTGCGCGGGACGTCCTTCTGCTACGTCCCTT  
 7801 CGGCCCTCAATCCAGCGGACCTTCCCTTCCCGCGCCTGCTGCCGGCTCTGCGGCCTCTTC  
 7861 CGCGTCTTCGCCTTCGCCCTCAGACGAGTCGGATCTCCCTTTGGGCCGCCTCCCCGCATC  
 7921 GATACCGTGCACCTCGATCGAGACCTAGAAAAACATGGAGCAATCACAAGTAGCAATACA  
 7981 GCAGCTACCAATGCTGATTGTGCCCTGGCTAGAAGCACAAGAGGAGGAGGAGGTGGGTTTT  
 8041 CCAGTCACACCTCAGGTACCTTTAAGACCAATGACTTACAAGGCAGCTGTAGATCTTAGC  
 8101 CACTTTTTTAAAAGAAAAGGGGGGACTGGAAGGGCTAATTCCTCCCAACGAAGACAAGAT  
 8161 ATCCTTGATCTGTGGATCTACCACACACAAGGCTACTTCCCTGATTGGCAGAACTACACA  
 8221 CCAGGGCCAGGGATCAGATATCCACTGACCTTTGGATGGTGTCTACAAGCTAGTACCAGTT  
 8281 GAGCAAGAGAAGGTAGAAGAAGCCAATGAAGGAGAGAACACCCCGCTTGTACACCCTGTG  
 8341 AGCCTGCATGGGATGGATGACCCGGAGAGAGAAGTATTAGAGTGGAGGTTTGACAGCCGC  
 8401 CTAGCATTTTCATGACATGGCCCCGAGACTGCATCCGGACTGTACTGGGTCTCTCTGGTTA  
 8461 GACCAGATCTGAGCCTGGGAGCTCTCTGGCTAACCTAGGGAAACCCACTGCTTAAGCTCAA  
 8521 TAAAGCTTGCCCTTGAGTGCCTCAAGTAGTGTGTGCCCGTCTGTTGTGTGACTCTGGTAAC  
 8581 TAGAGATCCCTCAGACCCTTTTAGTCAGTGTGGAATACTCTAGCAGCATGTGAGCAAAA  
 8641 GGCCAGCAAAAAGGCCAGGAACCGTAAAAAGGCCGCGTTGCTGGCGTTTTTCCATAGGCTC  
 8701 CGCCCCCTGACGAGCATCACAAAATCGACGCTCAAGTCAGAGGTGGCGAAACCCGACA  
 8761 GGACTATAAAGATAACCAGGCGTTTCCCCCTGGAAGCTCCCTCGTGCCTCTCCTGTTCCG  
 8821 ACCCTGCCGCTTACCGGATACCTGTCCGCCTTCTCCCTTCGGGAAGCGTGGCGCTTTCT  
 8881 CATAGCTCACGCTGTAGGTATCTCAGTTCGGTGTAGGTGCTTCGCTCCAAGCTGGGCTGT  
 8941 GTGCACGAACCCCCGTTTCCAGCCGACCGCTGCGCCTTATCCGGTAACTATCGTCTTGAG  
 9001 TCCAACCCGGTAAGACACGACTTATCGCCACTGGCAGCAGCCACTGGTAACAGGATTAGC  
 9061 AGAGCGAGGTATGTAGGCGGTGCTACAGAGTCTTGAAGTGGTGGCCTAACTACGGCTAC  
 9121 ACTAGAAGAACAGTATTTGGTATCTGCGCTCTGCTGAAGCCAGTTACCTTCGGAAAAGA  
 9181 GTTGGTAGCTCTTGTCCGGCAAACAACCACCGCTGAGGTAGCGGTGGTTTTTTTTTTTGC  
 9241 AAGCAGCAGATTACGCGCAGAAAAAAGGATCTCAAGAAGATCCTTTGATCTTTTCTACG  
 9301 GGGTCTGACGCTCAGTGGAACGAAAACCTACGTTAAGGGATTTTGGTCATGAGATTATCA  
 9361 AAAAGGATCTTCACCTAGATCCTTTTAAATTAATAATGAAGTTTTAAATCAATCTAAAGT  
 9421 ATATATGAGTAAACTTGGTCTGACAGTTACCAATGCTTAATCAGTGAGGCACCTATCTCA  
 9481 GCGATCTGTCTATTTCTGTTCCATCCATAGTTGCCCTGACTCCCCGTCGTGTAGATAACTACG  
 9541 ATACGGGAGGGCTTACCATCTGGCCCCAGTGCATGATAACCGCGAGACCCACGCTCA  
 9601 CCGGCTCCAGATTTATCAGCAATAAACCAGCCAGCCGGAAGGGCCGAGCGCAGAAGTGGT  
 9661 CCTGCAACTTTTATCCGCCTCCATCCAGTCTATTAATTGTTGCCGGGAAGCTAGAGTAAGT  
 9721 AGTTCGCCAGTTAATAGTTTGGCGAACGTTGTGTCATTGCTACAGGCATCGTGGTGTCA  
 9781 CGCTCGTCTGTTGGTATGGCTTTCATTCAGCTCCGGTTCCTCAACGATCAAGGCGAGTTACA  
 9841 TGATCCCCATGTTGTGCAAAAAGCGGTTAGCTCCTTCGGTCTCCGATCGTTGTGTCAGA  
 9901 AGTAAGTTGGCCGAGTGTATCACTCATGGTTATGGCAGCACTGCATAATTCTCTTACT  
 9961 GTCATGCCATCCGTAAGATGCTTTTCTGTGACTGGTGAGTACTCAACCAAGTCATTCTGA  
 10021 GAATAGTGTATGCGGCGACCGAGTTGCTCTTGGCCGCGTCAATACGGGATAATACCGCG  
 10081 CCACATAGCAGAACTTTAAAAGTGCTCATCATTTGAAAACGTTCTTCGGGGCGAAAACCTC  
 10141 TCAAGGATCTTACCGCTGTTGAGATCCAGTTCGATGTAACCCACTCGTGCACCCAACTGA  
 10201 TCTTCAGCATCTTTTACTTTTACCAGCGTTTCTGGGTGAGCAAAAACAGGAAGGCAAAAT  
 10261 GCCGCAAAAAGGGAATAAGGGCGACACGGAAATGTTGAATACTCATACTCTTCTTTTTT  
 10321 CAATATTAATTGAAGCATTTTATCAGGGTTATTGTCTCATGAGCGGATACATATTTGAATGT  
 10381 ATTTAGAAAAATAAACAAATAGGGGTTCCGCGCACATTTCCCCGAAAAGTGCCACCTGAC

2622-3276 bp: CMV promoter

3348-4064 bp: TALE N-terminus

4065-5756 bp: TALE repeat array (NI NN NG HD NI NG HD NN HD NG NI NG NG NI HD HD NI)

5757-6290 bp: TALE C-terminus

6291-6332 bp: NLS

6333-6497 bp: VP64

6504-6566 bp: 2A peptide

6567-7274 bp: mCherry

**RVD Library Entry Vector**

The following sequence was inserted into the TA-cloning site of pMD19-T vector (Takara, Inc.).

```

1   TAGCTATACGTCTCATTGACCCCCGAACAGGTTGTAGCCATAGCTTCTAAGTCTTCAGAG
61  ACGCTGGCTTATCGAAATTAATACGACTCACATATAGGGAGACCCAAGCTGGCTAGTTAAG
121 CTATCAACAAGTTTGTACAAAAAAGCTGAACGAGAAACGTAAAATGATATAAATATCAAT
181 ATATTTAAATTAGATTTTGCATAAAAAACAGACTACATAATACTGTAAAACACAACATATC
241 CAGTCACTATGAATCAACTACTTAGATGGTATTAGTGACCTGTAGTCGACCGACAGCCTT
301 CCAAATGTTCTTTCGGGTGATGCTGCCAACTTAGTCGACCGACAGCCTTCCAAATGTTCTT
361 CTCAAAACGGAATCGTCGTATCCAGCCTACTCGCTATTGTCTCAATGCCGTATTAAATCA
421 TAAAAAGAAAATAAGAAAAAGAGGTGCGAGCCTCTTTTTTGTGTGACAAAATAAAAAACATC
481 TACCTATTTCATATACGCTAGTGTATAGTCTGAAAATCATCTGCATCAAGAACAATTTT
541 ACAACTCTTATACTTTTCTCTTACAAGTCGTTTCGGCTTCATCTGGATTTTTCAGCCTCTAT
601 ACTTACTAAACGTGATAAAGTTTCTGTAATTTCTACTGTATCGACCTGCAGACTGGCTGT
661 GTATAAGGGAGCCTGACATTTATATTTCCCGAAGACATCAGGTTAATGGCGTTTTTGTATGT
721 CATTTTTCGCGGTGGCTGAGATCAGCCACTTCTTCCCGGATAACGGAGACCGGCACACTGG
781 CCATATCGGTGGTCATCATGCGCCAGCTTTCATCCCGGATATGCACCACCGGGTAAAGTT
841 CACGGGAGACTTTATCTGACAGCAGACGTGCACTGGCCAGGGGGATCACCATCCGTCGCC
901 CGGGCGTGTCAATAATATCACTCTGTACATCCACAAACAGACGATAACGGCTCTCTCTTT
961 TATAGGTGTA AACCTTAAACTGCATTTTACCAGCCCTGTTCTCGTCAGCAAAAAGAGCCG
1021 TTCATTTCAATAAACCGGGGACCTCAGCCATCCCTTCCCTGATTTTCCGCTTTCCAGCGT
1081 TCGGCACGCAGACGACGGGCTTTCATTTCTGCATGGTTGTGCTTACCAGACCGGAGATATTG
1141 ACATCATATATGCCTTGAGCAACTGATAGCTGTGCTGTCAACTGTCAGTGTAAATACGCT
1201 GCTTCATAGCATACTCTTTTTGACATACTTCGGGTATACATATCAGTATATATTCTTAT
1261 ACCGCAAAAATCAGCGCGCAAATACGCATACTGTTATCTGGCTTTTAGTAAGCCGGATCC
1321 ACGCGGCGTTTACGCCCCCTGCCACTCATCGCAGTACTGTTGTAATTCATTAAGCATT
1381 CTGCCGACATGGAAGCCATCACAACGGCATGATGAACCTGAATCGCCAGCGGCATCAGC
1441 ACCTTGTCGCCTTGCGTATAATATTTGCCCATGGTGAACCGGGGGCGAAGAAGTTGTCC
1501 ATATTGGCCACGTTTAAATCAAACCTGGTGAACCTCACCCAGGGATTGGCTGAGACGAAA
1561 AACATATTCTCAATAAACCCTTTAGGGAAATAGGCCAGGTTTTCACCGTAACACGCCACA
1621 TCTTGCGAATATATGTGTAGAAACTGCCGGAATCGTCGTGGTATTCACTCCAGAGCGAT
1681 GAAAACGTTTCAGTTTGCATGGAACCGGTGTAACAAGGGTGAACACTATCCCATATC
1741 ACCAGCTCACCGTCTTTTCAATGCCATACGGAATTCGGGATGAGCATTTCATCAGGCGGCA
1801 AGAATGTGAATAAAGCCGGATAAAAACCTTGTGCTTATTTTTCTTTACGGTCTTTAAAAAG
1861 GCCGTAATATCCAGCTGAACGGTCTGGTTATAGGTACATTGAGCAACTGACTGAAATGCC
1921 TCAAAAATGTTCTTTACGATGCCATTTGGGATATATCAACGGTGGTATATCCAGTGATTTTT
1981 TTCTCCATTTTAGCTTCCCTTAGCTCCTGAAAATCTCGATAACTCAAAAAATACGCCCGGT
2041 AGTGATCTTATTTTCAATATGGTGAAGTTGGAACCTCTTACGTGCCGATCAACGTCTCAT
2101 TTTCGCCAAAAGTTGGCCCAGGGCTTCCCGGTATCAACAGGGACACCAGGATTTATTTAT
2161 TCTGCGAAGTGATCTTCCGTACAGGTATTTATTCGGCGCAAAGTGCGTCGGGTGATGCT
2221 GCCAACTTAGTCTGACTACAGGTCAATAACCATCTAAGTAGTTGATTCATAGTGACTGG
2281 ATATGTTGTGTTTTACAGTATTATGTAGTCTGTTTTTTATGCAAAATCTAATTTAATATA
2341 TTGATATTTATATCATTTTACGTTTCTCGTTTCAGCTTCTTGTACAAAGTGGTTGATCTA
2401 GAGGGCCCCGGTTCCGACGTCCTGAAGACAAGGAGGTAAGCAGGCACTGGAACCGTG
2461 CAGCGCCTGCTCCAGTACTGTGTGTCAGGCTCATGGGTGAGACGTATAGCTA

```

16-48 bp: N-terminus of the TALE repeats

57-2425 bp: ccdB cassette

2434-2496 bp: C-terminus of the TALE repeats

**RVD Library Vector**

1 GTCGACGGATCGGGAGATCTCCCATCCCCATGGTGCCTCTCAGTACAATCTGCTCTG  
61 ATGCCGCATAGTTAAGCCAGTATCTGCTCCCTGCTTGTGTGTTGGAGGTCGCTGAGTAGT  
121 GCGCGAGCAAAATTTAAGCTACAACAAGGCAAGGCTTGACCCGACAATTGCATGAAGAATC  
181 TGCTTAGGGTTAGGCGTTTTGCGCTGCTTCGCGATGTACGGGCCAGATATACGCGTTGAC  
241 ATTGATTATTGACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCAT  
301 ATATGGAGTTCGCGTTACATAACTTACGGTAAATGGCCCGCTGGCTGACCGCCCAACG  
361 ACCCCCGCCATTGACGTCAATAATGACGTATGTTCCCATAGTAACGCCAATAGGGACTT  
421 TCCATTGACGTCAATGGGTGGAGTATTTACGGTAAACTGCCCACTTGGCAGTACATCAAG  
481 TGTATCATATGCCAAGTACGCCCCCTATTGACGTCAATGACGGTAAATGGCCCGCCTGGC  
541 ATTATGCCCAGTACATGACCTTATGGGACTTTCCCTACTTGGCAGTACATCTACGTATTAG  
601 TCATCGCTATTACCATGGTGTATGCGGTTTTGGCAGTACATCAATGGGCGTGGATAGCGGT  
661 TTGACTCACGGGGATTTCCAAGTCTCCACCCCATTTGACGTCAATGGGAGTTTTGTTTTGGC  
721 ACCAAAATCAACGGGACTTTTCCAAAATGTCGTAACAACCTCCGCCCATTTGACGCAAAATGG  
781 GCGGTAGGCGTGTACGGTGGGAGGTCATATAAGCAGCGCGTTTTGCTGTACTGGGTCT  
841 CTCTGGTTAGACCAGATCTGAGCCTGGGAGCTCTCTGGCTAACTAGGGAACCCACTGCTT  
901 AAGCCTCAATAAAGCTTGCCTTGAGTGTCTCAAGTAGTGTGTGCCCGTCTGTTGTGTGAC  
961 TCTGGTAACTAGAGATCCCTCAGACCTTTTTAGTCAAGTGTGGAAAATCTCTAGCAGTGGC  
1021 GCCCGAACAGGGACTTGAAAGCGAAAGGGAAACCAGAGGAGCTCTCTCGACGCAGGACTC  
1081 GGCTTGTGTAAGCGCGCACGGCAAGAGGCGAGGGGCGCGACTGGTGAGTACGCCAAAAA  
1141 TTTTACTAGCGGAGGCTAGAAGGAGAGAGATGGGTGCGAGAGCGTCAGTATTAAGCGGG  
1201 GGAGAAATTAGATCGCGATGGGAAAAAATTCGGTTAAGGCCAGGGGGAAAGAAAAAATATA  
1261 AATTAAAACATATAGTATGGGCAAGCAGGGAGCTAGAACGATTCGCAGTTAATCCTGGCC  
1321 TGTTAGAAACATCAGAAGGCTGTAGACAATACTGGGACAGCTACAACCATCCCTTCAGA  
1381 CAGGATCAGAAGAACTTAGATCATTATATAATACAGTAGCAACCCCTTATTGTGTGCATC  
1441 AAAGGATAGAGATAAAAGACACCAAGGAAGCTTTAGACAAGATAGAGGAAGAGCAAAACA  
1501 AAAGTAAGACCACCGCACAGCAAGCGGCCGGCCGCGCTGATCTTCAGACCTGGAGGAGGA  
1561 GATATGAGGGACAATTGGAGAAGTGAATTATATAAAATATAAAGTAGTAAAAATTGAACCA  
1621 TTAGGAGTAGCACCCACCAAGGCAAAGAGAAGAGTGGTGCAGAGAGAAAAAGAGCAGTG  
1681 GGAATAGGAGCTTTGTTCCTTGGGTTCTTGGGAGCAGCAGGAAGCACTATGGGCGCAGCG  
1741 TCAATGACGCTGACGGTACAGGCCAGACAATTATTGTCTGGTATAGTGCAGCAGCAGAAC  
1801 AATTTGCTGAGGGCTATTGAGGCGCAACAGCATCTGTTGCAACTCACAGTCTGGGGCATC  
1861 AAGCAGCTCCAGGCAAGAATCCTGGCTGTGGAAGATACCTAAAGGATCAACAGCTCCTG  
1921 GGGATTTGGGGTTGCTCTGGAAAACCTATTTCACCACCTGCTGTGCCTTGGAATGCTAGT  
1981 TGGAGTAATAAATCTCTGGAACAGATTTGGAATCACACGACCTGGATGGAGTGGGACAGA  
2041 GAAATTAACAATTACACAAGCTTAATACACTCCTTAATTGAAGAATCGCAAAACCAGCAA  
2101 GAAAAGAATGAACAAGAATTATTGGAATTAGATAAATGGGCAAGTTTGTGGAATTGGTTT  
2161 AACATAACAAAATTGGCTGTGGTATATAAAAATTATTATAATGATAGTAGGAGGTTGGTA  
2221 GGTTTAAGAATAGTTTTTGTGCTGACTTTCTATAGTGAATAGAGTTAGGCAGGGATATTCA  
2281 CCATTATCGTTTTCAGACCCACCTCCCAACCCGAGGGGACCCGACAGGCCCGAAGGAATA  
2341 GAAGAAGAAGGTGGAGAGAGAGACAGAGACAGATCCATTCGATTAGTGAACGGATCGGCA  
2401 CTGCGTGCGCCAATTCCTGCAGACAAATGGCAGTATTCATCCACAATTTTAAAAGAAAAGG  
2461 GGGGATTTGGGGGTACAGTGCAGGGGAAAGAATAGTAGACATAATAGCAACAGACATACA  
2521 AACTAAAGAATTACAAAAACAAATTACAAAAATTCAAAATTTTCGGGTTTTATTACAGGGA  
2581 CAGCAGAGATCCAGTTTGGTTAGTACCGGGCCCGCTCTAGACGATGTACGGGCCAGATAT  
2641 ACGCGTTGACATTGATTATTGACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTT  
2701 CATAGCCCATATATGGAGTTCCGCGTTACATAACTTACGGTAAATGGCCCGCTGGCTGA  
2761 CCGCCCAACGACCCCGCCCATTTGACGTCAATAATGACGTATGTTCCCATAGTAACGCCA  
2821 ATAGGGACTTTCCATTGACGTCAATGGGTGGAGTATTTACGGTAAACTGCCCACTTGGCA  
2881 GTACATCAAGTGTATCATATGCCAAGTACGCCCCCTATTGACGTCAATGACGGTAAATGG  
2941 CCGCCTGGCATTATGCCAGTACATGACCTTATGGGACTTTCCCTACTTGGCAGTACATC  
3001 TACGTATTAGTCATCGCTATTACCATGGTGTATGCGGTTTTGGCAGTACATCAATGGGCGT  
3061 GGATAGCGGTTTACTCACGGGATTTCCAAGTCTCCACCCCATTTGACGTCAATGGGAGT  
3121 TTGTTTTTGGCACCAAAATCAACGGGACTTTCCAAAATGTCGTAACAACCTCCGCCCATTTG  
3181 ACGCAAAATGGGCGTAGGCGTGTACGGTGGGAGGTCATATAAGCAGAGCTCTCTGGCTA  
3241 ACTAGAGAACCACCTGCTTACTGGCTTATCGAAATTAATACGACTCACTATAGGGAGACC  
3301 CAAGCTGGCTAGCGAAGTTCCTATTCTCTAGAAAAGTATAGGAACTTCATGAGGACCAGGC  
3361 TGCCATCTCCCCCTGCCCTTCCCCCGCATTTAGCGCTGGGAGCTTTAGCGACCTGCTTA  
3421 GGCAGTTTCGACCCAGCTTGTTCACACCAGCCTGTTTACAGCCTGCCTCCCTTCGGAG  
3481 CGCACCACACCGAAGCCGCCACCGGCGAGTGGGACGAGGTGCAAAGCGGCCTGAGGGCAG  
3541 CGGACGCTCCTCCGCCAACCATGAGGGTGGCAGTGACAGCAGCTAGGCCCCCTCGGGCAA

3601 AACCTGCACCCAGGAGAAGGGCTGCCAACCCAGCGACGCGAGTCCAGCCGCACAGGTGG  
3661 ACCTCAGGACGCTGGGCTACAGCCAGCAACAGCAAGAGAAGATCAAGCCCAAAGTAAGGA  
3721 GCACCGTGGCCAGCACCACGAGGCCCTGGTGGGTACGGCTTCACCCACGCGCATATCG  
3781 TTGCTCTGAGCCAACATCCCGCAGCTCTGGGTACCGTTGCGGTGAAGTATCAGGACATGA  
3841 TCGCGGCACTGCCTGAAGCTACACACGAAGCCATAGTGGGCGTTGGCAAGCAGTGGAGCG  
3901 GTGCCAGAGCGCTTGAGGCACTGTTGACGGTGGCTGGCGAGCTGAGGGGACCGCCACTGC  
3961 AACTGGACACCGCCAACTGCTGAAGATCGCCAAGAGGGGAGGCGTGACGGCGGTGGAGG  
4021 CCGTGCAATGCCCTGGAGGAATGCCCTGACCGGCGGCCCTGAACAGAGACGCTGGCTTAT  
4081 CGAAATTAATACGACTCACATATAGGGAGACCCAAGCTGGCTAGTTAAGCTATCAACAAGT  
4141 TTGTACAAAAAAGCTGAACGAGAAAACGTAATAATGATATAAATATCAATATATTAAATTAG  
4201 ATTTTGCATAAAAAACAGACTACATAATACTGTA AAAACACAACATATCCAGTCACTATGA  
4261 ATCAACTACTTAGATGGTATTAGTGACCTGTAGTTCGACCGACAGCCTTCCAAATGTTCTT  
4321 CGGGTGATGCTGCCAACTTAGTTCGACCGACAGCCTTCCAAATGTTCTTCTCAAACGGAAT  
4381 CGTCGTATCCAGCCTACTCGCTATTGTCTCAATGCCGTATTAATCATAAAAAAGAAATA  
4441 AGAAAAAGAGGTGCGAGCCTCTTTTTTGTGTGACAAAATAAAAAACATCTACCTATTCTATA  
4501 TACGCTAGTGTCTATAGTCCTGAAAATCATCTGCATCAAGAACAATTTCACTACTTTATA  
4561 CTTTTCTCTTACAAGTCTTCCGGCTTCATCTGGATTTTCAGCCTCTATACTTACTAAACG  
4621 TGATAAAGTTTCTGTAATTTCTACTGTATCGACCTGCAGACTGGCTGTGTATAAGGGAGC  
4681 CTGACATTTATATTTCCCCAGAACATCAGGTTAATGGCGTTTTTGATGTCATTTTCGCGGT  
4741 GGCTGAGATCAGCCACTTCTTCCCGGATAACGGGACCGGCACACTGGCCATATCGGTGG  
4801 TCATCATGCGCCAGCTTTCATCCCGATATGCACCACCGGGTAAAGTTACGGGAGACTT  
4861 TATCTGACAGCAGACGTGCACTGGCCAGGGGATCACCATCCGTCGCCCGGGCGTGTCAA  
4921 TAATATCACTCTGTACATCCACAAACAGACGATAACGGCTCTCTCTTTTTATAGGTGTAAA  
4981 CCTTAAACTGCATTTACCAGCCCTGTTCTCGTCAGCAAAAGAGCCGTTTCAATTTCAATA  
5041 AACCGGGCGACCTCAGCCATCCCTTCTGATTTTCCGCTTTCAGCGTTCGGCACGCAGA  
5101 CGACGGGCTTCAATCTGCATGGTGTGCTTACCAGACCGGAGATATTGACATCATATATG  
5161 CCTTGAGCAACTGATAGCTGTGCTGTCAACTGTCACTGTAATACGCTGCTTCATAGCAT  
5221 ACCTCTTTTTGACATACTTCGGGTATACATATCAGTATATATTCTTATACCGCAAAAATC  
5281 AGCGCGCAAATACGCATACTGTTATCTGGCTTTTAGTAAGCCGGATCCACGCGGGCTTTA  
5341 CGCCCCCTGCCACTCATCGCAGTACTGTTGTAATTCATTAAGCATTCTGCCGACATGG  
5401 AAGCCATCACAAACGGCATGATGAACCTGAATCGCCAGCGGCATCAGCACCTTGTGCGCT  
5461 TGCGTATAAATTTGCCCATGGTGA AAACGGGGGCGAAGAAGTTGTCCATATTGGCCACG  
5521 TTTAAATCAAACCTGGTGA AACCTACCCAGGGATTGGCTGAGACGAAAAACATTTCTCA  
5581 ATAAACCTTTTAGGGAATAGGCCAGGTTTTTACCCTAACACGCCACATCTTGCAATAT  
5641 ATGTGTAGAACTGCCGGAATCGTCTGTTGATTTCACTCCAGAGCGATGAAAACGTTTCA  
5701 GTTTGCTCATGGAAAACGGTGTAAACAAGGGTGAACACTATCCCATATCACCAGCTCACCG  
5761 TCTTTCATTGCCATACGGAATTCGGGATGAGCATTCATCAGGCGGGCAAGAATGTGAATA  
5821 AAGCCGGATAAAACTTGTGCTTATTTTTCTTTACGGTCTTTAAAAAGGCCGTAATATCC  
5881 AGCTGAACGGTCTGGTTATAGGTACATTGAGCAACTGACTGAAATGCCTCAAATGTTCT  
5941 TTACGATGCCATTGGGATATATCAACGGTGGTATATCCAGTGATTTTTTTCTCCATTTTA  
6001 GCTTCTTAGCTCCTGAAAATCTCGATAACTCAAAAATACGCCCGGTAGTGATCTTATT  
6061 TCATTTATGGTGAAGTTGGAACCTCTTACGTGCCGATCAACGTCTCATTTTCGCCAAAAG  
6121 TTGGCCAGGGCTTCCCGGTATCAACAGGGACACCAGGATTTATTTATTCTGCGAAGTGA  
6181 TCTTCCGTACAGGTATTTATTCGGCGCAAAGTGCCTCGGGTGTGATGCCAACTTATGTC  
6241 GACTACAGTCACTAATACCATCTAAGTAGTTGATTCATAGTGACTGGATATGTTGTGTT  
6301 TTACAGTATTATGTAGTCTGTTTTTTTATGCAAAAATCTAATTTAATATATTGATATTTATA  
6361 TCATTTTACGTTTCTCGTTCAGCTTTCTTGTACAAAGTGGTTGATCTAGAGGGCCCGCGG  
6421 TTCGAACGTCCTTAGCATCGTGGCCAGCTGTCTCGGCCCGACCCTGCCCTCGCCGCTCT  
6481 GACCAACGACCACCTGGTGGCCCTGGCTTGCCCTCGGGGGCAGGCCAGCTCTTGACGCCGT  
6541 GAAGAAGGGCTTCCCTACGCCCCAGCCCTGATCAAGCGGACCAACAGAAGGATTCCCGA  
6601 GAGGACATCACATCGAGTGGCAGATCACGCGCAAGTGGTCCGCGTGCTCGGATTCTTCCA  
6661 GTGTCAC TCCCACCCGCACAAGCGTTCGATGACGCCATGACTCAATTTGGTATGTGCGAG  
6721 ACACGGACTGCTGCAGCTCTTTTCGTAGAGTCCGGTGTACAGAACTCGAGGCCCGCTCGGG  
6781 CACACTGCCCTCCCGCTCCAGCGGTGGGACAGGATTTCTCCAAGCGAGCGGTATGAAACG  
6841 CGCGAAGCCTTACCCTACGTCAACTCAGACACCTGACCAGGCGAGCCTTCATGCGTTTCG  
6901 AGACTCGCTGGAGAGGGATTTGGACGCGCCCTCGCCATGCATGAAGGGGACCAAACTCG  
6961 CGCGTACGATGCCCAAAGAAGAAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG  
7021 CGCATTTGACGATTTTGTATCTGGATATGCTGGGAAGTGCAGCCCTCGATGATTTTGACCT  
7081 TGACATGCTTGGTTCCGATGCCCTTGTATGACTTTGACCTCGACATGCTCGGCAGTGACGC  
7141 CCTTGATGATTTTCGACCTGGACATGCTGATTAACCTTAGAGGCAGTGGAGAGGGCAGAGG  
7201 AAGTCTGCTAACATGCGGTGACGTCGAGGAGAATCCTGGCCAGTGAGCAAGGGCGAGGA  
7261 GGATAACATGGCCATCATCAAGGAGTTTATGCGCTTCAAGGTGCACATGGAGGGCTCCGT  
7321 GAACGGCCACGAGTTCGAGATCGAGGGCGAGGGCGAGGGCCGCCCTACGAGGGCACCCA

7381 GACCGCCAAGCTGAAGGTGACCAAGGGTGGCCCCCTGCCCTTCGCCTGGGACATCCTGTG  
7441 CCCTCAGTTCATGTACGGCTCCAAGGCCTACGTGAAGCACCCCGCCGACATCCCCGACTA  
7501 CTTGAAGCTGTCCCTCCCGAGGGCTTCAAGTGGGAGCGCGTGATGAACTTCGAGGACGG  
7561 CGGCGTGGTGACCGTGACCCAGGACTCCTCCCTGCAGGACGGCGAGTTCATCTACAAGGT  
7621 GAAGCTGCGCGGCACCAACTTCCCTCCGACGGCCCCGTAATGCAGAAGAAGACCATGGG  
7681 CTGGGAGGCCCTCCTCCGAGCGGATGTACCCCGAGGACGGCGCCCTGAAGGGCGAGATCAA  
7741 GCAGAGGCTGAAGCTGAAGGACGGCGGCCACTACGACGCTGAGGTCAAGACCACCTACAA  
7801 GGCCAAGAAGCCCGTGCAGCTGCCCGCGCCCTACAACGTCAACATCAAGTTGGACATCAC  
7861 CCCCCACAACGAGGACTACACCATCGTGGAACAGTACGAACGCGCCGAGGGCCGCCACTC  
7921 CACCGGCGGCATGGACGAGCTGTACAAGTAACATGTTTAAGGGTTCGGTTCCTACTAGGT  
7981 ACAATTTCGATATCAAGCTTATCGATAATCAACCTCTGGATTACAAAATTTGTGAAAGATT  
8041 GACTGGTATTTCTTAACATATGTTGCTCCTTTTACGCTATGTGGATACGCTGCTTTAATGCC  
8101 TTTGTATCATGCTATTTGCTTCCCGTATGGCTTTCATTTTCTCCTTGTATAAATCCTG  
8161 GTTGTGTCTCTTTATGAGGAGTTGTGGCCCGTTGTCAGGCAACGTGGCGTGGTGTGCAC  
8221 TGTGTTTGTGACGCAACCCCCACTGGTGGGGCATTGCCACCACCTGTGAGCTCCTTTC  
8281 CGGGACTTTCGCTTTCCCCCCTCCCTATTGCCACGGCGGAACCTATCGCCGCCTGCCTTGC  
8341 CCGCTGCTGGACAGGGGCTCGGCTGTTGGGCACTGACAATTCGGTGGTGTGTCGGGGAA  
8401 ATCATCGTCCCTTCCCTGGCTGCTCGCCTGTGTTGCCACCTGGATTCTGCGCGGGACGTC  
8461 CTTCCTGCTACGTCCCTTCGGCCCTCAATCCAGCGGACCTTCCCTCCCGCGGCCCTGCTGCC  
8521 GGCTCTGCGCCCTTCCGCGTCTTCGCTTCGCTTCGCTTCGCTTCGCTTCGCTTCGCTTCG  
8581 GGCCGCTCCCGCTCATGATACCGTACCGTTCGCTTCGCTTCGCTTCGCTTCGCTTCGCTTCG  
8641 TCACAAGTAGCAATACAGCAGCTACCAATGCTGATTGTGCCTGGCTAGAAGCACAAGAGG  
8701 AGGAGGAGGTGGGTTTTCCAGTACACCTCAGGTACCTTTAAGACCAATGACTTACAAGG  
8761 CAGCTGTAGATCTTAGCCACTTTTTAAAAGAAAAGGGGGACTGGAAGGGCTAATTCACT  
8821 CCCAACGAAGACAAGATATCCTTGATCTGTGGATCTACCACACACAAGGCTACTTCCCTG  
8881 ATTTGGCAGAACTACACACCAGGGCCAGGGATCAGATATCCACTGACCTTTGGATGGTGTCT  
8941 ACAAGCTAGTACCAGTTGAGCAAGAGAAGGTAGAAGAAGCCAATGAAGGAGAGAACACCC  
9001 GCTTGTACACCTGTGAGCCTGCATGGGATGGATGACCCGGAGAGAGAAGTATTAGAGT  
9061 GGAGGTTTTGACAGCCGCTTAGCATTTTCATCACATGGCCCCGAGAGCTGCATCCGGACTGTA  
9121 CTGGGTCTCTCTGGTTAGACCAGATCTGAGCCTGGGAGCTCTCTGGCTAACTAGGGAAACC  
9181 CACTGCTTAAGCCTCAATAAAGCTTGCCCTGAGTGTCTCAAGTAGTGTGTGCCCGTCTGT  
9241 TGTGTGACTCTGGTAACTAGAGATCCTCAGACCTTTTAGTCAGTGTGGAAAATCTCTA  
9301 GCAGCATCTGAGCAAAAGGCCAGCAAAAAGGCCAGGAACCGTAAAAAGGCCGCGCTTGTG  
9361 CGTTTTTCCATAGGCTCCGCCCCCTGACGAGCATCACAAAAATCGACGCTCAAGTCAGA  
9421 GGTGGCGAAACCCGACAGGACTATAAAGATACCAGGCGTTTTCCCCCTGGAAGCTCCCTCG  
9481 TGCGCTCTCCTGTTCCGACCTTGCCGCTTACCAGGATACCTGTCCGCTTTTCTCCCTTCGG  
9541 GAAGCGTGGCGCTTTCCTCATAGCTCACGCTGTAGGTATCTCAGTTCGGTGTAGGTCGTTT  
9601 GCTCCAAGCTGGGCTGTGTGCACGAACCCCCGTTCCAGCCGACCGCTGCGCCTTATCCG  
9661 GTAACATATCGTCTTTGAGTCCAACCCGGTAAGACACGACTTATCGCCACTGGCAGCAGCCA  
9721 CTGGTAACAGGATTAGCAGAGCGAGGTATGTAGGCGGTGCTACAGAGTTCTTGAAGTGGT  
9781 GGCCTAACACGGCTACACTAGAAGAACAGTATTTGGTATCTGCGCTCTGCTGAAGCCAG  
9841 TTACCTTCGGAAAAAGAGTTGGTAGCTCTTGATCCGGCAAACAAACCACCGCTGGTAGCG  
9901 GTGGTTTTTTTTGTTTTGCAAGCAGCAGATTACGCGCAGAAAAAAGGATCTCAAGAAGATC  
9961 CTTTGTATCTTTTCTACGGGGTCTGACGCTCAGTGGAACGAAAACCTCACGTTAAGGGATTT  
10021 TGGTCATGAGATTATCAAAAAGGATCTTACCCTAGATCCTTTTAAATTAAAAATGAAGTT  
10081 TTAAATCAATCTAAAGTATATATAGTAAACTTGGTCTGACAGTTACCAATGCTTAATCA  
10141 GTGAGGCACCTATCTCAGCGATCTGTCTATTTGTTTCATCCATAGTTGCCTGACTCCCCG  
10201 TCGTGTAGATAACTACGATACGGGAGGGCTTACCATCTGGCCCCAGTGCTGCAATGATAC  
10261 CGCGAGACCCACGCTCACCGGCTCCAGATTTATCAGCAATAAACCAGCCAGCCGGAAGGG  
10321 CCGAGCGCAGAAGTGGTCCTGCAACTTTATCCGCTCCATCCAGTCTATTAATTGTTGCC  
10381 GGGAAGCTAGAGTAAGTAGTTCCGCCAGTTAATAGTTTGCGCAACGTTGTTGCCATTGCTA  
10441 CAGGCATCGTGGTGTACGCTCGTCTGTTGGTATGGCTTCATTCAGCTCCGTTCCCAAC  
10501 GATCAAGGCGAGTTACATGATCCCCATGTTGTGCAAAAAAGCGGTTAGCTCCTTCGGTCT  
10561 CTCCGATCGTTGTGCAAGTAAGTTGGCCGCACTGTTATCACTCATGGTTATGGCAGCAC  
10621 TGCATAATTCCTTACTGTGATGCCATCCGTAAGATGCTTTTCTGTGACTGGTGTGACTACT  
10681 CAACCAAGTCATTCGAGAATAGTGTATGCGGCGACCGAGTTGCTCTTGGCCGGCGCTCAA  
10741 TACGGGATAAATACCAGCCACATAGCAGAACTTTAAAAGTGTCTCATCTTGGAAAACGTT  
10801 CTTCCGGGCGAAAACTCAAGGATCTTACCCTGTTGAGATCCAGTTTCGATGTAACCCA  
10861 CTGCTGCACCCAACTGATCTTACGATCTTTTACTTTTACCAGCGTTTTTGGGTGAGCAA  
10921 AAACAGGAAGGCAAAATGCCGCAAAAAGGGAATAAGGGCGACACGGAATGTTGAATAC  
10981 TCATACTCTTCTTTTCAATATTTATGAAGCATTTATCAGGGTTATTGTCTCATGAGCG  
11041 GATACATATTTGAATGTATTTAGAAAAATAAACAAATAGGGGTTCCGCGCACATTTCCCC  
11101 GAAAAGTGCCACCTGAC

2622-3276 bp: CMV promoter  
3348-4064 bp: TALE N-terminus  
4065-6433 bp: ccdB cassette  
6434-6967 bp: TALE C-terminus  
6968-7009 bp: NLS  
7010-7174 bp: VP64  
7181-7243 bp: 2A peptide  
7244-7951 bp: mCherry

## SUPPLEMENTARY METHODS

### *Artificial system for RVD screening*

The artificial screening system was composed of four reporters and a TALE-VP64 expression library in which the RVDs of three consecutive monomers in the middle of an artificial TALE array were encoded by the same 6 randomly synthesized nucleotides (TALE-(XX')<sub>3</sub>). TALE-(XX')<sub>3</sub> contained 14.5 repeats fused with the VP64 trans-activation domain and 2A peptide-linked mCherry. The variable diresidues (XX') for testing were placed in the 7<sup>th</sup> - 9<sup>th</sup> repeat modules, and the variable RVD-carrying TALE modules were purposely designed as triplets to augment the DNA binding capability. X and X' represents the 12<sup>th</sup> and 13<sup>th</sup> amino acids in the 7<sup>th</sup> - 9<sup>th</sup> repeat modules, respectively. Four reporters consist of TALE-(XX')<sub>3</sub> binding sites CTGGCCNNNTACGTA, in which N represents A, T, C or G, was located immediately upstream of a minimal CMV promoter (P<sub>minCMV</sub>) and its downstream EGFP gene. TALE-Ctrl was constructed to have the identical backbone as TALE-(XX')<sub>3</sub> except that its TALE repeats (16.5-mers) are different, not matching with any reporters. The mCherry-normalized EGFP level of TALE-Ctrl co-transfected with reporter served as the corresponding basal level. For each sample, EGFP fluorescence intensity was normalized to mCherry intensity. Fold induction is calculated as the result of normalized sample EGFP intensity divided by normalized basal level shown as follows.



$$\text{Fold Induction} = \frac{(\prod_{i=1}^n \text{Exp}_{\text{EGFP}_i})^{1/n}}{(\prod_{i=1}^n \text{Ctrl}_{\text{EGFP}_i})^{1/n}} \bigg/ \frac{(\prod_{i=1}^n \text{Exp}_{\text{mCherry}_i})^{1/n}}{(\prod_{i=1}^n \text{Ctrl}_{\text{mCherry}_i})^{1/n}}$$

### Notes:

$(\prod_{i=1}^n \text{Exp}_{\text{EGFP}_i})^{1/n}$ : Geometric mean of  $\text{Exp}_{\text{EGFP}}$  of cells from FACS ( $n$  = number of cells)

$\text{Exp}_{\text{EGFP}}$ : EGFP intensity of HEK293T cells co-transfected with TALE-(XX')<sub>m</sub> plus reporter ( $m$  = 3, 6 and 12, corresponding to TALE-(XX')<sub>3</sub>, TALE-(XX')<sub>6</sub> and TALE-(XX')<sub>12</sub>, respectively)

$\text{Exp}_{\text{mCherry}}$ : mCherry intensity of HEK293T cells co-transfected with TALE-(XX')<sub>m</sub> plus reporter

$\text{Ctrl}_{\text{EGFP}}$ : EGFP intensity of HEK293T cells co-transfected with TALE-Ctrl plus reporter

$\text{Ctrl}_{\text{mCherry}}$ : mCherry intensity of HEK293T cells co-transfected with TALE-Ctrl plus reporter

### Construction of the TALE-(XX')<sub>3</sub> library

A 102-nt monomer encoding a standard TALE repeat unit was synthesized with the following sequences: 5'-GCTATGGCTACAACCTGTTTCGGGGGTCAACCCATGAGCCT GACACAGTACTGGGAGCAGGCGCTGCACGGTTTCCAGTGCCTGCTTACCTCCNNNN NNAGAA-3'. The six random nucleotides (underlined Ns) corresponding to the RVD-encoding region were purposely placed near the 3' end to ensure unbiased oligonucleotide synthesis. This single-stranded DNA was cyclized using linker primer 1 (5'-GAACAGGTT GTAGCCATAGCTTCT-3') and T4 DNA ligase (NEB). Using the agarose gel-purified single stranded circular DNA as template, rolling-circle amplification was conducted with phi29 DNA polymerase (NEB) and primer 1 (30°C, 90 min) followed by primer extension using primer 2 (5'-AGGTTGTAGCCATAGCT-3') and phi29 (30°C, 90 min) to acquire long dsDNA. After ultrasonic shearing (270 W, work 10 s, pause 10 s, 10 cycles) and T4 DNA polymerase (NEB) treatment to blunt the ends of the DNA fragments (12°C, 15 min, with 400 μM dNTP mix), 250-400 bp DNA fragments were harvested by gel purification. These DNA fragments were then cloned into a pre-made entry vector using the ligase-independent cloning (LIC) method to create an entry library. BsmBI digestion of entry

library clones produced ~300 bp DNA fragments that were subsequently cloned, via the Golden Gate approach<sup>1</sup>, into a pre-made RVD library vector, which was constructed using the ULtiMATE protocol previously developed by our group<sup>2</sup>. Each plasmid in the final RVD library was verified through sequencing analysis. About 350 kinds of TALE-(XX')<sub>3</sub> constructs were obtained from this approach.

#### *Individual construction of TALE-(XX')<sub>3</sub>*

A complementary primer (5'-aaCGTCTCaGTTTCGGGGGTCAACCCATGAGCCTGACACAGTACTGGGAGCAGGGCGCTGCACGGTTTCCAGTGCCTGCTT-3') and a specific primer (5'-tCGTCTCaGAACAGGTTGTAGCCATAGCTTCTNNNNNNGGAGGTAAGCAGGCACTGGAA-3'; NNNNNN indicates the RVD codons) were annealed and PCR extended to generate a 102 bp monomer with BsmBI sites at each end. The monomer was ligated via a 6 cycles of the Golden Gate method<sup>1</sup> to generate repeats. The repeat product was PCR amplified using the primers G-lib-F (5'-TAGCTATACGTCTCATTGACCCCCGAACAGGTTGTAGCC-3') and G-lib-R (5'-TAGCTATACGTCTCACCCATGAGCCTGACACAGTACTGGGAGCA-3') and Taq Hifi (Transgen, Inc.). The 3-repeat fragment was gel purified and subsequently cloned into a pre-made RVD library vector via the Golden Gate method<sup>1</sup>. Trans1-T1 competent cells (Transgen, Inc.) were used for bacterial transformation. About 50 kinds of TALE-(XX')<sub>3</sub> constructs were obtained from this approach.

#### *Design and construction of TALE-(XX')<sub>6</sub> and TALE-(XX')<sub>12</sub>*

For TALE-(XX')<sub>6</sub>, the customized TALEs (17.5-mer) are the same as TALE-(XX')<sub>3</sub> except that there are six identical repeats containing the variable RVDs. Accordingly, four reporters were constructed, consisting of TALE-(XX')<sub>6</sub> binding sites with six consecutive nucleotides (A, T, C or G) substituted at positions 7 - 12 in front of a minimal CMV promoter and its downstream EGFP gene. The 7<sup>th</sup> - 9<sup>th</sup> repeats in TALE-(XX')<sub>3</sub> were PCR

amplified using the primers G-lib-seq-F (5'-TCTAGGTACCAAGCCCACGGATTGA-3') and G-lib-seq-R (5'-ATCGATCGTCCGGAGTGAGCCCA-3'). The 300 bp PCR product was purified and further PCR amplified using the primers G-lib-F and G-lib-R. The product of the 6-repeat fragment was gel purified and cloned into a pre-made RVD library vector via the Golden Gate method<sup>1</sup> to construct the final TALE-(XX')<sub>6</sub> plasmid.

For TALE-(XX')<sub>12</sub>, the customized TALEs (15.5-mer) contained four ACTC-targeting TALE repeat modules followed by 11.5 repeat units with variable RVDs to be tested. Accordingly, four reporters were constructed, consisting of TALE-(XX')<sub>12</sub> binding sites with 12 consecutive nucleotides (A, T, C or G) substituted at positions 5 - 16 in front of a minimal CMV promoter and its downstream EGFP gene. TALE-(XX')<sub>12</sub> plasmids were constructed with a similar method using TALE-(XX')<sub>6</sub> as a PCR template and applying the ULtiMATE protocol<sup>2</sup>.

|                           |             |
|---------------------------|-------------|
| Backbone plasmid          | 10 ng       |
| 6-repeat fragment         | 5 ng        |
| 10X Tango buffer (Thermo) | 1 µl        |
| 10 mM ATP                 | 1 µl        |
| BsmBI (Thermo)            | 0.75 µl     |
| T4 ligase (NEB)           | 0.25 µl     |
| ddH <sub>2</sub> O        | up to 10 µl |

|      |       |           |
|------|-------|-----------|
| 37°C | 5 min | \         |
|      |       | 10 cycles |
| 16°C | 5 min | /         |
| 37°C | 5 min |           |

#### *Cell culture, transfection and flow cytometric analysis*

HEK293T cells (from Stanley Cohen lab at Stanford University) were cultured in DMEM medium with 10% FBS and 1% penicillin-streptomycin at 37°C and 5% CO<sub>2</sub>. Cells were

seeded 24 h prior to transfection in 24-well plates at a density of  $10^5$  cells per well. The cells in each well were co-transfected with 0.2  $\mu\text{g}$  TALE-(XX')<sub>3</sub> plasmid and 0.4  $\mu\text{g}$  reporter plasmid using polyethylenimine (PEI)<sup>3</sup>. At 48 h post-transfection, the cells were collected and analyzed on a BD LSR II flow cytometer (BD Biosciences). Lasers with wavelengths of 488 nm and 561 nm were used to quantify EGFP and mCherry protein expression, respectively. At least 10,000 events were collected from every sample to obtain sufficient data for analysis. Majority of cells showed mCherry fluorescence intensity of  $5 \times 10^4 - 5 \times 10^5$ , thus were gated for further analysis. The binding efficiencies and specificities of the variable RVDs (XX') in each customized TALE were assayed by comparing the fold induction of EGFP reporters with the basal level of EGFP in HEK293T cells transfected with the reporter plasmid and a customized TALE plasmid containing unmatched TALE repeats (Supplementary Sequences). The EGFP fluorescence intensity assayed from FACS analysis was normalized to the corresponding mCherry fluorescence intensity.

#### *Generation of a heat map illustrating the base-preference of RVDs*

The heat map was generated from library screening of TALE-(XX')<sub>3</sub> using four reporters (3A, 3T, 3C, and 3G), thereby reflecting the base preference of 400 RVDs. EGFP activities from different reporters are coded by different colors representing the reporter identities (3A, green; 3T, red; 3C, blue; and 3G, yellow). The brightness of the colors indicates the fold induction of the reporters by TALE-(XX')<sub>3</sub> compared to their basal levels.

#### *Criteria for selection of novel RVDs from TALE-(XX')<sub>3</sub> for intensive study*

The standards are as follows: (1) the highest fold induction among four reporters is at least equal to NK for G recognition; and (2) the second highest fold induction is lower than half of the highest one. In addition to these simple rules, we have also included groups of

RVDs that displayed some unique patterns, i.e. seven RVDs ended with Ala (KA, CA, FA, YA, RA, PA, and AA) that showed preference for T-recognition, and RVDs recognizing both A and G (NN and HN).

### *Statistical analysis*

For comparison between the induction level of TALE-(XX')<sub>12</sub> and the basal level (Fig. 1f), two-sample, one-tailed t-test was performed assuming unequal variance. \* $P < 0.05$ , \*\* $P < 0.01$ , and \*\*\* $P < 0.005$ .

### *Construction of TALENs*

TALENs with natural RVDs (i.e. NI, NG, HD and NN) were constructed using ULtiMATE system as previously described<sup>2</sup>. For TALEN repeats using novel RVDs KN (in place of NN), NH (in place of NN) or RG (in place of NG), TALE monomers containing new RVDs were individually synthesized. The final assembly of these TALENs constructs was conducted using the same ULtiMATE protocol as above.

### *Assessment of TALENs-mediated indels*

HEK293T cells were seeded in 6-well plates at a density of  $3 \times 10^5$  cells per well and incubated at 37°C with 5% CO<sub>2</sub>. For each well, a pair of TALEN plasmids and pmaxGFP (Lonza Group Ltd.) were co-transfected at a ratio of 9:9:2 (0.9 µg : 0.9 µg : 0.2 µg) using PEI method. The transfected cells were cultured at 37°C for one day followed by 3 days of incubation at 30°C (cold shock) before flow cytometric sorting for GFP positive cells. TALENs-targeting regions were PCR-amplified from the genome DNA of the isolated GFP positive cells. The TALENs-mediated indels were analyzed by mismatch-sensitive T7 endonuclease I (T7E1; New England Biolabs) as described previously<sup>4</sup>.