



저작자표시-비영리-변경금지 2.0 대한민국

이용자는 아래의 조건을 따르는 경우에 한하여 자유롭게

- 이 저작물을 복제, 배포, 전송, 전시, 공연 및 방송할 수 있습니다.

다음과 같은 조건을 따라야 합니다:



저작자표시. 귀하는 원저작자를 표시하여야 합니다.



비영리. 귀하는 이 저작물을 영리 목적으로 이용할 수 없습니다.



변경금지. 귀하는 이 저작물을 개작, 변형 또는 가공할 수 없습니다.

- 귀하는, 이 저작물의 재이용이나 배포의 경우, 이 저작물에 적용된 이용허락조건을 명확하게 나타내어야 합니다.
- 저작권자로부터 별도의 허가를 받으면 이러한 조건들은 적용되지 않습니다.

저작권법에 따른 이용자의 권리는 위의 내용에 의하여 영향을 받지 않습니다.

이것은 [이용허락규약\(Legal Code\)](#)을 이해하기 쉽게 요약한 것입니다.

[Disclaimer](#)

이 학박사 학위논문

A genome-wide library of
TAL effector nucleases

2014년 2월

서울대학교 대학원

화학부 생화학 전공

김 용 섭

A genome-wide library of TAL effector nucleases

지도교수 김진수

이 논문을 이학박사학위논문으로 제출함
2013년 11월

서울대학교 대학원
화학과 생화학전공
김 용 섭

김용섭의 박사학위논문을 인준함
2013년 11월

위 원 장	_____	(인)
부 위 원 장	_____	(인)
위 원	_____	(인)
위 원	_____	(인)
위 원	_____	(인)

Abstract

A genome-wide library of TAL effector nucleases

Yongsub Kim

Department of Chemistry

The Graduate School

Seoul National University

Transcription activator-like (TAL) effector nucleases (TALENs) are newly developed programmable nucleases which use a simple ‘protein–DNA code’ that relates modular DNA-binding TALE repeat domains to individual bases in a target binding site. Unlike homing endonucleases and zinc finger nucleases (ZFNs), they can be readily engineered to bind specific genomic loci, enabling the introduction of precise genetic modifications such as gene disruptions, additions,

corrections and genome rearrangements. In this thesis, I developed improved TALEN architectures to avoid unwanted mutations in genome. Then, I carefully chose genome-wide TALEN target sites that did not have highly similar sequences elsewhere in the genome and assembled TALEN pairs using high throughput Golden Gate cloning system. A pilot test including over a hundred pairs of TALENs showed that all TALENs were active and disrupted their target genes at high frequencies, although two of these TALENs became active only after their target sites were partially demethylated using a DNA methyltransferase inhibitor. I used the TALEN library to generate single- and double-gene knockout cells in which NF- κ B signaling pathways were disrupted. Compared with cells treated with short interfering RNAs (siRNA), these cells showed unambiguous suppression of the signal transduction. Furthermore, I developed the TALEN library for targeting every exon in protein coding genes of several organisms including human. The TALEN library reported here will be broadly useful for research and drug discovery.

Keywords : TAL effector nuclease (TALEN), Double-strand breaks (DSB), Non-homologous end-joining (NHEJ), Genome engineering.

Student Number : 2008-22719

Table of Contents

Abstract	i
Table of Contents	iii
List of Figures	v
List of Tables	viii
List of Abbreviations	ix
I. Introduction	1
II. Materials and Methods	
1. Dual fluorescent reporter plasmid construction	4
2. Cell culture and transfection	4
3. High-throughput assembly of TALENs	4
4. Flow cytometry	6
5. T7E1 assay for mutation detection	6
6. PCR analysis to detect genomic mutations and sequencing	7
7. Analysis and rescue of genome-inactive TALENs	8
8. Digital PCR to estimate mutation frequency	8
9. Gene-knockout cell lines	8
10. Fluorescent PCR analysis	9
11. Episomal reporter assay to detect NF- κ B signaling	9
12. Western blotting	10
III. Results	
A. Optimization of TAL Effector Nucleases (TALENs)	

1. Dual fluorescent reporter system	11
2. Design of prototype TALENs	15
B. Human genome-wide TALEN library	
1. One-step Golden-Gate cloning system	18
2. Design of human genome-wide TALENs	29
3. TALEN-mediated genome editing activity	36
4. Rescue of inactive TALENs	46
5. Undetectable off-target mutations	52
6. TALEN induced genome rearrangement	56
C. TALEN-mediated knockout cell lines	
1. Establishment of knockout cell lines for NF- κ B pathway study	63
2. Episomal reporter assay	73
D. Expansion of TALEN library	
1. Design of expanded TALEN library	84
2. TALEN library for several organisms	91
E. Comparison of ZFNs and TALENs	96
IV. Discussion	105
V. Appendix	
1. List of TALEN activity in reporter assay	110
2. List of TALEN activity in the T7E1 assay	120
VI. References	123
Abstract in Korean	136

List of Figures

Figure 1. Scheme of dual fluorescent reporter-based assay	13
Figure 2. Improved results of dual fluorescent reporter assay for measuring activities of engineered nucleases	14
Figure 3. Optimization of TALEN architectures	16
Figure 4. Scheme of one-step Golden-Gate cloning system	20
Figure 5. High-throughput synthesis of TALENs	23
Figure 6. Validation of TALENs constructed by Golden-Gate assembly	25
Figure 7. Pilot test of 15 TALENs	28
Figure 8. Reporter-based assay for detecting TALEN activities	37
Figure 9. 103 of TALENs test using the T7E1 assay	39
Figure 10. Distribution of TALEN activities	41
Figure 11. DNA sequences of indels induced by TALENs	42
Figure 12. Comparison between episomal reporter assay and T7E1 assay	45
Figure 13. Episomal reporter assays of two genome-inactive TALENs	47
Figure 14. TALEN-driven mutations in drug-treated cells	48
Figure 15. Targeted mutagenesis using alternative sets of TALENs	51
Figure 16. Undetectable off-target mutations with TALENs	55
Figure 17. TALEN-mediated targeted genomic deletions	58
Figure 18. TALEN-mediated targeted genomic inversions	60

Figure 19. TALEN-mediated targeted genomic duplications	61
Figure 20. TALEN-mediated targeted genomic translocations	62
Figure 21. Validation of knockout cell lines	65
Figure 22. Undetectable off-target mutations in TALEN-mediated knockout cell lines	71
Figure 23. Undetectable gene expression levels in gene knockout cells	72
Figure 24. Schematic of reporter assay related NF- κ B signaling	75
Figure 25. Functional assay in knockout cells using reporter	76
Figure 26. Cytokine-independent NF- κ B activation by mitoxantrone	78
Figure 27. Validation of knockout cells related siRNA screening	80
Figure 28. Invalidation of a false-positive gene identified in a genome-wide siRNA screen: <i>NR1H4</i>	82
Figure 29. Invalidation of a false-positive gene identified in a genome-wide siRNA screen: <i>SMEK</i>	83
Figure 30. TALEN-mediated mutations in each exons of <i>NRAS</i>	88
Figure 31. TALEN-mediated homology-dependent repair (HDR)	90
Figure 32. TALEN-mediated gene disruption in zebrafish	94
Figure 33. Sequencing validation of TALEN-induced indels in zebrafish	95
Figure 34. ZFN and TALEN mutation patterns reported in the literature	98
Figure 35. Overlapping target sites of ZFNs and TALENs	100
Figure 36. Comparison of ZFNs and TALENs targeting overlapping sites	101

List of Tables

Table 1. List of 15 TALEN target sites in pilot test	27
Table 2. Pilot test of computational strategy for TALENs design	31
Table 3. Summary of human genome-wide TALEN library	34
Table 4. Analysis of TALEN target sites in the human genome	35
Table 5. Potential off-target sites of highly active TALENs in the human genome	54
Table 6. Potential off-target sites of TALENs used for knockout cell lines	70
Table 7. Summary of TALEN target sites for exons in human protein- coding genes	86
Table 8. List of TALEN target sites for targeting <i>NRAS</i>	87
Table 9. Summary of TALEN library for several organisms	92
Table 10. List of TALEN target sites for zebrafish	93

List of Abbreviations

DSB	Double-strand break
FACS	Fluorescence-activated cell sorting
HEK	Human embryonic kidney
HR	Homologous recombination
KO	Knockout
NHEJ	Non-homologous end-joining
PCR	Polymerase chain reaction
T7E1	T7 Endonuclease 1
ZFN	Zinc-finger nuclease
TALEN	Transcription activator-like effector nuclease
TALE	TAL Effector
RNAi	RNA interference
siRNA	short interfering RNA
RGEN	RNA-Guided Endonuclease

I. Introduction

In the post-genome era, although human genome has been sequenced (Lander et al. 2001; Venter et al. 2001), the precise function of most genes is still unknown. Hence, one of the most challenging tasks in biological research is investigation the function of these genes. Important contributions to our understanding of gene functions often derive from knockout studies that a specific gene is silenced in cells or organisms achieved by gene targeting through homologous recombination (HR) (Smithies et al. 1985), an endogenous DNA double strand break (DSB) repair mechanism. Although gene targeting is widely used in mouse embryonic stem cells to create gene knockout animals, the efficiency of HR is extremely low in mammalian and other higher eukaryotic cells, ranging from 10^{-6} to 10^{-7} (Deng and Capecchi 1992).

In a decade ago, attractive alternative to the above described techniques was established, called RNA interference (RNAi). RNAi, discovered in the nematode *Caenorhabditis elegans*, is used to suppress the expression of a gene of interest by artificial short interfering RNAs (siRNAs) (Fire et al. 1998). Despite of the broad utility of this approach in basic research and drug discovery, siRNAs are limited by many factors: 1) siRNAs induce gene knockdown rather than gene knockout. Thus, small fraction of the activities of target genes remains after siRNA treatment (Krueger et al. 2007). 2) siRNAs are not specific, displaying sequence dependent off target effects (Jackson et al.

2003). Only several nucleotide matches between an siRNA seed sequence and the 3' untranslated region (UTR) of an off target gene can trigger suppression of gene expression (Birmingham et al. 2006). 3) siRNAs can broadly affect cellular physiology by activating an innate immune response and competing with endogenous microRNAs (miRNAs) for proteins such as Dicer and the RNA-induced silencing complex (RISC), which are critical for miRNA function (Sledz et al. 2003; Khan et al. 2009). 4) Furthermore, many genes are refractory to inhibition via siRNAs (Krueger et al. 2007).

Engineered nucleases such as zinc finger nucleases (ZFNs) (Bibikova et al. 2003; Urnov et al. 2005; Doyon et al. 2008; Maeder et al. 2008), TAL effector nucleases (TALENs) (Cermak et al. 2011; Miller et al. 2011; Sander et al. 2011) and CRISPR/Cas9-based RNA-guided endonucleases (RGENs) (Cho et al. 2013; Cong et al. 2013; Hwang et al. 2013; Mali et al. 2013), enable the targeted alteration precisely at a predetermined locus. These programmable nucleases induce site specific DSBs in a genome, which are then repaired by endogenous mechanisms that can be exploited to create sequence alterations at the cleavage site. ZFNs and TALENs, although they share the same FokI-derived nuclease domain, differ in that they employ distinctive DNA binding arrays: ZFNs use zinc finger arrays (Kim et al. 1996) and TALENs use TAL effector (TALE) repeat arrays (Miller et al. 2011). Because these arrays recognize target DNA sequences in a modular fashion, tailor-made DNA-binding arrays with desired specificities can be constructed by mixing-and-matching pre-characterized

modules (Kim et al. 2011).

In this thesis, I described novel TALEN architectures for genome engineering using improved fluorescent reporter system (Kim et al. 2013b). This TALEN architecture showed high activity at target site with 12- to 14-bp spacers comparable with TALEN constructs in previous study. And then, to prepare a collection of TALENs that is designed to target every protein coding gene in the human genome, I developed one-step Golden Gate cloning system to assemble TALEN plasmids and computational strategy to search for unique TALEN target sites in each gene to avoid potential off target mutations. Nearly 100% of TALENs induced site specific mutations at high frequencies. Only two TALENs failed to show any measurable mutagenesis activities. Each of two TALEN sites was turned out to be heavily methylated. With this TALEN library, I generated single- and double-gene knockout cells in which NF- κ B signaling pathways were disrupted. Furthermore, I designed TALEN library for several organisms targeting every exon of protein coding genes and showed that almost TALENs were highly active. Finally, I investigated difference of mutation signature between TALENs and ZFNs. I found that ZFN-induced insertion frequency was much higher than the TALEN-induced insertion frequency. From these result, the TALEN library resources will be broadly useful for research and drug discovery.

II. Materials and Methods

1. Dual fluorescent reporter plasmid construction

The amplified product of encoding EGFP sequence from pEGFP-N1 using Primer A and B was cloned into the NotI site of the pRGS plasmid (Kim et al. 2011). This vector was named pRG2S. Oligonucleotides that contained each target sites were synthesized (Macrogen, Seoul, South Korea) and annealed *in vitro*. The annealed oligonucleotides were ligated into the vector (pRG2S) digested with EcoR1 and BamH1. As previous pRGS plasmid, note that an in-frame stop codon should be avoided at the target site either by altering the frame or by changing the orientation of the target site.

2. Cell culture and transfection

HEK293T/17 (ATCC, CRL-11268) and HeLa cells (ATCC, CCL-2) were maintained in Dulbecco's modified Eagle's medium (DMEM) supplemented with 100 units/mL penicillin, 100 µg/mL streptomycin, 0.1 mM nonessential amino acids, and 10% FBS (FBS). We transfected 200,000 HEK293 cells using 3 µl of polyethylenimine and 1 µg of plasmid DNA in 24-well plates. We transfected 100,000 HeLa cells with Lipofectamine 2000 (Invitrogen) according to the manufacturer's protocol.

3. High-throughput assembly of TALENs

All steps in TALEN assembly were performed in 96-well plates. In each plate, 47 pairs of TALENs were assembled and one pair of negative control (FokI vector alone) was included. Our one-step Golden-Gate system consisted of 424 TALE array plasmids (6×64 tripartite arrays, 2×16 bipartite arrays, and 2×4 monopartite arrays). For convenience, we numbered each TALE array as follows. We used these numbers to choose appropriate arrays to assemble TALEN plasmids.

For example, the half-site sequence, "5'-TGGGGGAGGTGGCG AGGAAC", can be divided into 8 parts (the first T, GGG, GGA, GGT, GGC, GAC, GAA, and the last C). The first T and last C are not recognized by TALE arrays. To assemble a TALEN subunit specific to this sequence, we chose the following arrays: position1-GGG + position2-GGA + position3-GGT + position4-GGC + position5-GAC + position6-GAA + the FokI expression vector that contains C-specific half-repeat. A detailed protocol is shown below:

1) Six TALE array plasmids and a FokI expression vector are mixed in each well as follows:

1.0 μ l	TALE array vectors (50ng/ μ l each)
0.5 μ l	FokI expressing vector (50ng/ μ l)
0.5 μ l	BsaI (New England BioLabs, 10U/ μ l)
2.0 μ l	$10\times$ T4 DNA Ligase Reaction Buffer
0.1 μ l	T4 DNA Ligase (New England BioLabs, 2000U/ μ l)
10.9 μ l	ddH ₂ O

In a 20 μ l restriction-ligation reaction.

2) The restriction-ligation reaction is carried out in a thermocycler as follows:

37°C 5 min] → 20 cycles

16°C 5 min

50°C 15 min

80°C 5 min

3) After the thermocycling reaction, the reaction mixture (6 µl) from each well is used for transformation of chemically-competent DH5 α cells (30 µl). Subsequently, cells are inoculated in Flat-Bottom Blocks (Qiagen) that are filled with LB medium (800 µl) containing ampicillin (50µg/ml). The transformants in 96-well blocks are incubated overnight at 37°C with vigorous shaking.

4) Two sets of *E. coli* stocks are prepared by mixing the LB culture (50 µl) with 60% glycerol (150 µl); they are separately stored at -80°C.

4. Flow cytometry (FACS)

Transfected HEK293T/17 and HeLa cells were trypsinized and resuspended in growth media. Single-cell suspensions were analyzed and sorted using the FACSCanto (BD Biosciences). Untransfected cells and cells transfected with reporters alone were used as controls.

5. T7E1 assay for mutation detection

HEK293T/17 cells (2×10^5) pre-cultured in a 24 well plate were

transfected with two plasmids encoding a TALEN pair (500ng each). After 72h of incubation, genomic DNA was extracted from the transfected cells using the G-dexTM Genomic DNA Extraction Kit (iNtRON BIOTECH NOLOGY, Seongnam, S.Korea). Purified genomic DNA samples were subjected to the T7E1 assay as described previously (Kim et al. 2009). Briefly, genomic region around TALEN target site was amplified, melted, and annealed to form heteroduplex DNA. The annealed DNA was treated with 5 units of T7E1 endonuclease I (T7E1) for 20 min at 37°C. DNA was analyzed by agarose gel electrophoresis. PCR primers used for the T7E1 assay in previous study (Kim et al. 2013b).

6. PCR analysis to detect genomic mutations and sequencing

Genomic DNA (50 ng per reaction) was subjected to PCR analysis using Taq DNA polymerase (GeneAll Biotech, Seoul, Korea) and appropriate primers as described in previous study (Kim et al. 2013b). For sequencing analysis, PCR products corresponding to genomic rearrangements were purified using the QIAquick Gel Extraction Kit (Qiagen, Valencia, CA) and cloned into the T-Blunt vector using the T-Blunt PCR Cloning Kit (SolGent, Daejeon, S. Korea). Cloned PCR products were sequenced using the M13 primer or primers used for PCR amplification.

7. Analysis and rescue of genome-inactive TALENs

One day before TALEN plasmid transfection, HEK293 cells (two million) were pretreated with 0.2 μ M 5-aza-dC or 100 ng/ml trichostatin A in 24-well plates. After 3 days of incubation, genomic DNA was isolated from transfected cells and subjected to the T7E1 assay. To determine whether the TALEN sites were methylated, genomic DNA was treated with bisulfite using the EpiTect Bisulfite kit (QIAGEN) according to the manufacturer's protocol. The bisulfite-converted DNA was then amplified using bisulfite-specific primers, and the amplified products were subcloned into the T-Blunt vector and sequenced.

8. Digital PCR to estimate mutation frequency

Digital PCR analysis was performed based on previous study (Lee et al. 2010). Genomic DNA samples were quantified and serially diluted with Tris-EDTA buffer. The serially diluted genomic DNA were amplified using appropriate primers. The frequencies of mutation were calculated as previous study.

9. Gene-knockout cell lines

HEK293T/17 and HeLa cells were co-transfected with TALEN plasmids and surrogate reporter plasmids that contain the TALEN target site. Gene-knockout cells were enriched by selection as described (Kim

et al. 2013a) and cloned by limiting dilution in 96-well plates. Typically, cells were maintained for 2 weeks in 96-well plates to isolate single clones. These clones were analyzed using T7E1, fPCR and dideoxy sequencing.

10. Fluorescent PCR analysis

Genomic DNA was extracted from each single clones and subjected to fluorescence PCR using 5'-carboxyfluorescein-labeled primers. Each PCR amplicons were analyzed using an ABI 3730xl DNA analyzer. The positions of peaks indicate the lengths of PCR products.

11. Episomal reporter assay to detect NF- κ B signaling

The NF- κ B-dependent firefly luciferase reporter was constructed by placing three tandem copies of the NF- κ B recognition element (TGGGGACTTTCCGC) (Duan et al. 2005) in front of a synthetic promoter that consists of the TATA-box and the initiator element. Gene-knockout or wild-type cells were co-transfected with the luciferase reporter plasmid and the *Renilla* luciferase plasmid. After 24 h of incubation, cells were treated with TNF α (1 ng/ml) or IL-1 β (25 ng/ml) and incubated for 15 h. Cells were lysed in 1 \times lysis buffer (100 μ l) (Promega), and the dual luciferase assays were done according to the manufacturer's protocol.

12. Western blotting

HEK293 knockout clones were lysed and the lysates were electrophoresed on a 7% SDS-PAGE gel. Primary antibodies specific for *TNFR1* (1:200, Santa cruz biotechnology) or *GAPDH* (1:200, Santa cruz biotechnology) and anti-mouse secondary antibody (1:1000, Santa cruz biotechnology) were used. Immunoreactive bands were visualized using the ECL method.

III. Results

A. Optimization of TAL Effector Nucleases (TALENs)

1. Dual fluorescent reporter system

First of all, to measure genome-editing activity of engineered nucleases, I developed dual fluorescent reporter system (Figure 1). In the previous study, surrogate reporter systems were developed for detection of engineered nuclease's activity and enrichment of mutant cells (Kim et al. 2011). The reporter plasmid contains nuclease's target site between the RFP- and GFP- encoding DNA sequences. Because the GFP sequence fused to the RFP sequence is out of frame, the cells transfected with the reporter plasmid and inactive nuclease-encoding plasmid express RFP only. In contrast, functional GFP is expressed only when nucleases induce DSBs at target site, whose repair via error-prone NHEJ lead to indels that often result in frameshifting mutations. Although this system is valuable resource for nuclease mediated genome editing, it is limited that only one third of the reporter plasmids, which introduced +1 or +2 frameshift, can produce in-frame frameshifting. Actually, some active TALENs did not generate functional GFP proteins in this reporter system (Figure 2). In this regard, I developed improved fluorescence reporter for detecting nuclease's activity by fusing another GFP sequence at the behind of

GFP sequence as other frame. Nuclease-mediated indels formation generate two different frame shift. The +1 frameshifting mutation makes the former GFP protein functional, The +2 frameshifting mutation makes latter GFP protein functional. In this dual fluorescent reporter assay, TALENs which were inactive in the previous reporter assay were highly active (Figure 2). This result showed that dual fluorescent reporter system is valuable for detecting engineered nucleases activity.

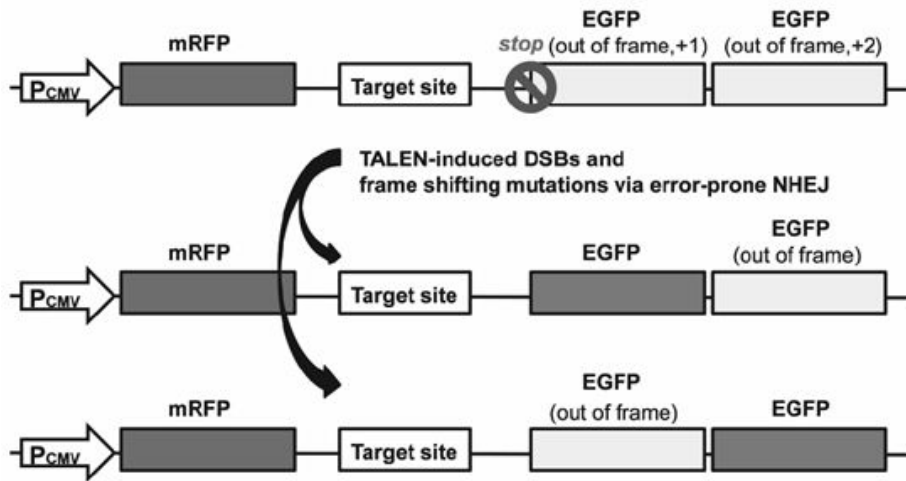
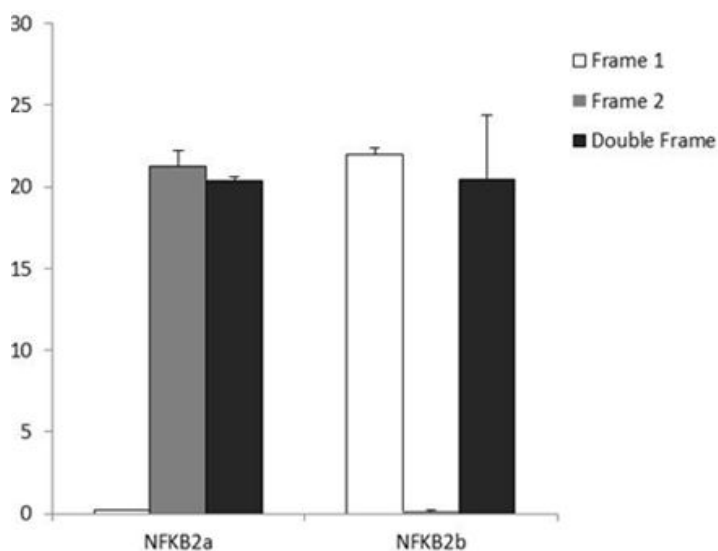


Figure 1. Scheme of dual fluorescent reporter-based assay. The reporter consists of the mRFP gene, the engineered nuclease's target sequence and the two EGFP genes. mRFP is constitutively expressed by the CMV promoter, whereas functional EGFP is not expressed because it is out of frame. Nuclease-mediated indel formation in target sites generate two different frame shifts. The +1 frameshifting mutation makes the former GFP protein functional, The +2 frameshifting mutation makes the latter GFP protein linked with fused non-functional amino acid functional.

a



b

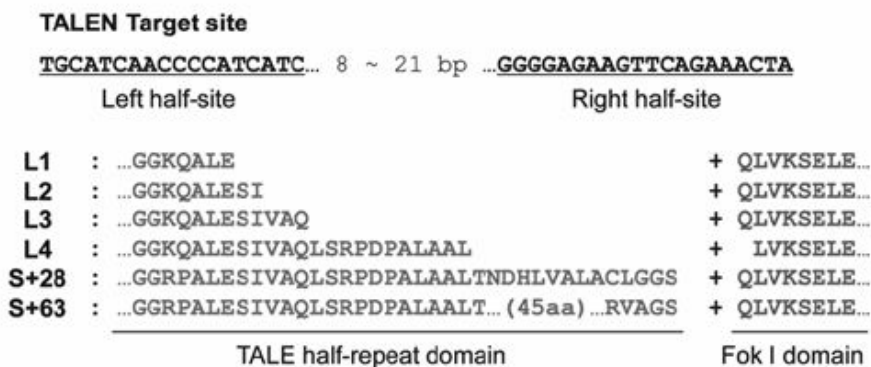
Label	Target sequence
NFkB2a	<u>tcgggggtggctccacatgggtggaggctctgggggtgcagccgggggcta</u>
NFkB2b	<u>tcgactacggcgctaccgcggacgcgcgcgcgtgctggcgggacagcgcca</u>

Figure 2. Improved results of dual fluorescent reporter assay for measuring activities of engineered nucleases. (a) Two TALEN pairs were not active in RGS reporter assay (Kim et al. 2011). But in other frame reporters and RG2S reporters, these TALENs were highly active. (b) Target sequence of TALEN pairs in this assay.

2. Design of prototype TALENs

In the previous study, TAL effectors fused FokI nuclease domain can induce DSBs in human cells and leads to gene disruption as ZFNs. Although TALENs are highly active than ZFNs, their broad spacer between two TALEN monomer which could induce unwanted mutation in genome is barrier for gene therapy or other applications (Christian et al. 2010; Li et al. 2011; Miller et al. 2011). To optimize the fusion junction between a TALE domain and the FokI nuclease domain, I prepared a series of TALE-FokI fusions with different junctions by linking each TALE to various amino acid residues in the appropriate region of the FokI nuclease domain (Figure 3a). I tested these TALE-FokI fusions using the dual reporter system with variable spacers between two TALEN monomer binding sites. As the optimized TALENs, I chose L4 that showed high activity at target sites with 12- to 14-bp spacers but no or negligible activity at sites with < 12-bp or > 14-bp spacers (Figure 3b). Compared to the two original TALEN constructs that contain additional amino acid residues between the TAL effector array and the FokI sequence (S+28 and S+63) (Miller et al. 2011), L4 TALEN constructs showed nuclease activities at sites with a narrow range of spacers, a desirable property for high specificity.

a



b

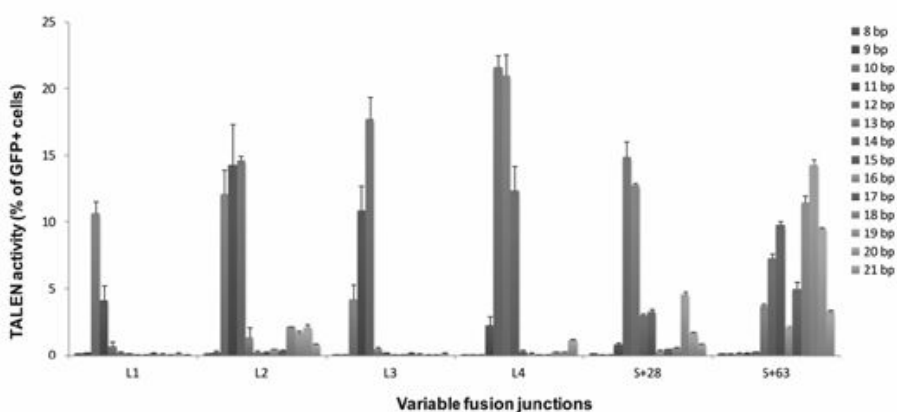


Figure 3. Optimization of TALEN architectures. (a) TALEN target site and amino-acid sequences in the fusion junctions that connect the TAL effector array to the FokI cleavage domain. (b) Comparison of TALEN gene-editing activities. Reporter plasmids that contain target sites with variable spacers and TALEN plasmids were co-transfected into HEK293 cells, and GFP+ cells were counted by flow cytometry. S+28 and S+63

are two prototype TALEN architectures previously reported (Miller et al. 2011). Error bars indicated s.e.m. from at least three independent experiments.

B. Human genome-wide TALEN library

1. One-step Golden-Gate cloning system

I developed a one-step Golden-Gate cloning system to assemble TALEN plasmids with variable lengths in a high-throughput manner (Figure 4). Although Golden-Gate cloning methods have been used for the assembly of TALEN plasmids previously (Cermak et al. 2011; Li et al. 2011; Morbitzer et al. 2011; Weber et al. 2011; Zhang et al. 2011), these methods either rely on the use of PCR, or require gel isolation of DNA segments, or require at least two rounds of subcloning steps. First of all, I used four TAL effector repeat domains termed NI, NN, NG and HD repeat variable di-residues (RVDs), each specific to one of the four bases (A, G, T and C, respectively), to make the cloning system (Boch et al. 2009; Moscou and Bogdanove 2009). These repeat domains consisted of 34 amino acid residues with similar sequences; the RVDs at positions 12 and 13 determine the base specificities. Next, I designed 64 tripartite-, 16 bipartite- and 4 monopartite- TAL effector repeat domain arrays with minimum sequence similarity as following criteria; i) removed human's rare codon ($< 10\%$) and generated ii) 192 (64×3) monopartite with minimum similarity (max. 81.25%), iii) 64 tripartite, 16 bipartite with minimum similarity in single arrays (max. 72%). iv) Then, these arrays were subcloned to 6 subgroups via PCR-based mutagenesis. Therefore, I prepared a total of 424 TAL effector array plasmids (6×64 tripartite arrays, 2×16 bipartite arrays,

and 2×4 monopartite arrays). Then, to make engineered nucleases, I prepared 8 obligatory heterodimeric FokI-encoding plasmids (DAS/RR) with L4 fusion junction in Figure 3 (Guo et al. 2010).



Figure 4. Scheme of one-step Golden-Gate cloning system. A total of 424 TAL effector array plasmids ($64 \times 6 + 16 \times 2 + 4 \times 2$) (Kan^R) and 8 FokI plasmids (Amp^R) are used. One member in each of the six positions is chosen; the six arrays are combined with each other for subcloning into one of the FokI expression plasmids. This system allows construction of TALEN plasmids that contain at least 14.5 (4 tripartite arrays + 2 monopartite arrays) and up to 18.5 (6 tripartite arrays) RVD modules in a single Golden-Gate reaction.

The TAL effector array plasmids are separated into six subgroups in accordance with their positions (Figure 4). When digested with BsaI restriction enzyme, each array in a given position generates the same four-base pair overhang sequences. TAL effector arrays in different subgroups produce different four-base overhang sequences. One member in each of the six subgroups is chosen; the six arrays are combined with each other for subcloning into one of the FokI expression plasmids. This system allows construction of TALEN plasmids that contain at least 14.5 (4 tripartite arrays + 2 monopartite arrays) and up to 18.5 (6 tripartite arrays) RVDs in a single Golden-Gate reaction. The last half-repeat (0.5) is encoded in the FokI plasmids. These TALENs recognize DNA sequences of 16–20 bps in length, including a conserved base T at the 5' end. Because TALENs function as dimers, these TALEN pairs recognize 32- to 40-bp DNA sequences, which consist of two half-sites separated by 12- to 14-bp spacers because of L4 fusion junction. This one-step Golden-Gate cloning system enables us to TALEN synthesis high-throughput manner (Figure 5).

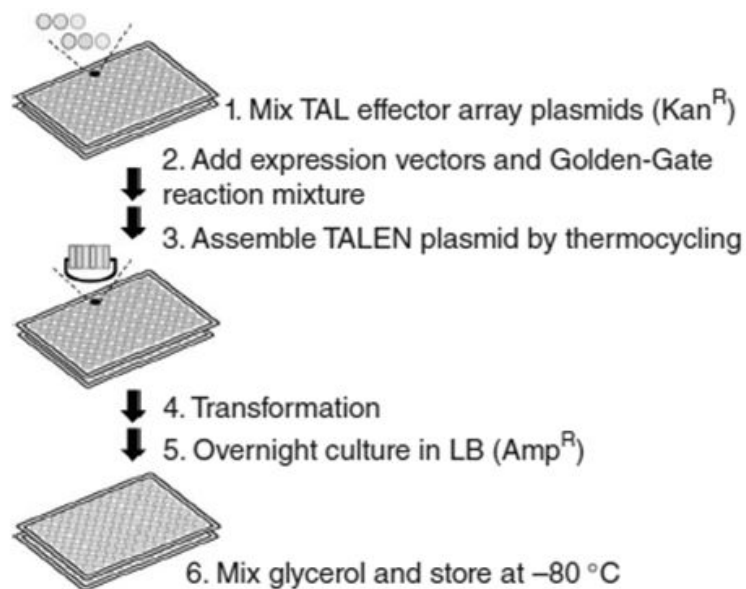


Figure 5. High-throughput synthesis of TALENs. High-throughput Golden-Gate cloning could process in 96-well plates. Six TAL effector array plasmids and one FokI plasmid are mixed in each well. BsaI releases the TAL effector arrays and allows an ordered assembly of six TAL effector arrays into the FokI plasmid.

I investigated whether TALEN pairs made by one-step Golden-Gate system are functional as TALEN pairs generated by modular assembly method. I generated TALEN pair targeting *CCR5* gene as previous study (Kim H.J. 2011) (Figure 6a), and confirmed protein expression levels and mutagenesis activity. TALEN pair that synthesized by Golden-Gate cloning method (GG-TALEN) was highly expressed and induced similarly mutagenesis with TALENs synthesized by modular assembly method (MA-TALEN) (Figure 6b, c). And then, I constructed 15 TALEN pairs targeting each human genes (Table 1). Because longer TALENs are more active (Reyon et al. 2012), I designed these TALEN pairs with 18.5 RVD modules. I measured the mutagenesis activities of these TALEN pairs in HEK293 cells using T7E1 assay. HEK293 cells were transfected with TALEN-encoding plasmids, and PCR amplicons from genomic DNA were analyzed by T7E1. Mutation frequencies were detected from the intensities of cleaved bands relative to intact bands. Mutations were estimated at all 15 sites at frequencies of 3.9-43% (Figure 7). These pilot tests demonstrate that both the new TALEN construct and the one-step Golden-Gate cloning system are robust enough to allow genome-scale construction of TALENs.

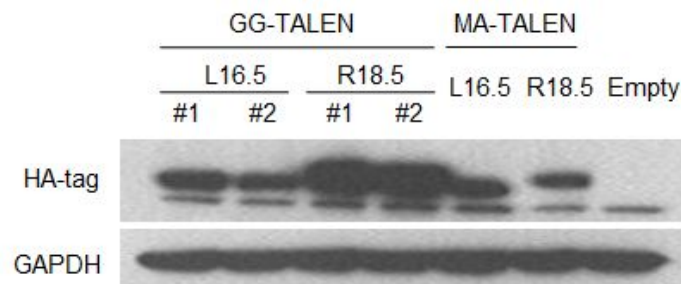
a

Human CCR5 (891 site)

5' -ATGACGCACTGCTGCATCAACCCCATCATCTATGCCTTTGTCGGGGAGAAGTTCAGAACTACCTCTTAGTCTTC-3'
 3' -TACTGCGTGACGACGTAGTTGGGGTAGTAGATACGGAACAGCCCTCTTCAAGTCTTTGATGGAGAATCAGAAG-5'

L16.5 TGCATCAACCCCATCATC CCCCTCTTCAAGTCTTTGAT R18.5

b



c

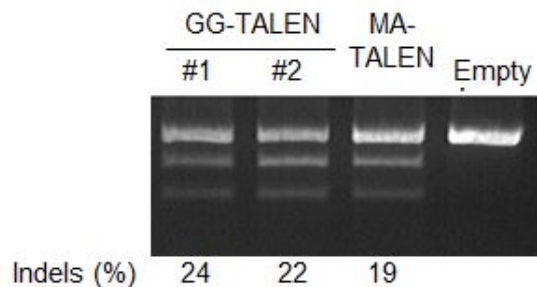


Figure 6. Validation of TALENs constructed by Golden-Gate assembly.

(a) TALENs targeted CCR5 in previous study (Kim H.J. 2011) were synthesized. (b) Protein expression levels were detected by Western blotting. Protein expression of TALEN pairs synthesized by Golden-Gate cloning method (GG-TALEN) were comparable with TALEN pairs constructed by modular assembly (MA-TALEN). Each numbers are

independently synthesized TALENs. (c) T7E1 assay were detected endogenous mutagenesis activities of TALEN pairs. GG-TALENs could induce mutations in endogenous target sites with high frequency. The numbers at the bottom of the gels indicate mutation frequencies measured by band intensities

Table 1. List of 15 TALEN target sequences in pilot test

	Gene name	TALEN binding site	Number of CpG sites			Mutation frequency (%)
			Left site	Right site	Total	
Pilot study	ACAT1	TGCTTGAAGCAAGAGAAaaggtcttgccaGTATTGCAATGGAGGAGGA	0	0	0	9.5
	ANGPT1	TGGGATAAAGTGGCACTACTtcaaagggcccaGTTACTCCTTACGTTCCACA	0	1	1	24
	AIRE	TCTACAAGCACCTGCGGCTcogcettctgtcaGCCCCGCTGCCAGGGCTGGA	1	1	2	43
	CYP11A1	TGGTGCAAGTGGCCATCTATgctctgggcccgaGAGCCACCTTCTTCTTGA	0	1	1	11
	F8	TTACTGCTTCATCCTACTTTaccaatagtgttGCCACCTGGTCTCCTTCAAA	0	0	0	26
	FOXF1	TCAAQCCATGGCTCTCTTccatgcactcgGCGGGGGGCTCCTACTA	2	2	4	18
	GFAP	TGGAGGAAGAGGGGCAGAGCctcaaggacgagATGGCCGCCACTTGCAGGA	0	1	1	11
	AGT	TGGCCGAGGCTGAGCTGCGGccattctgtcacACGAGCTGAACCTGCAAAA	2	1	3	4.2
	KRT23	TGTTATCGAGACCAAGTCTcggtactctgtcAAGCTCCAGGACATGCAAGA	1	0	1	3.9
	BBS9	TATTTTGAAAAACAGGGAGTcaaagattttgcATGTCTTTTTGCGGATCTA	0	1	1	13
	TRIM45	TCACCCCTGCTTTGTAAAGATgcgcagagagaaATCATGGGCAGGGAGGAGA	0	0	0	15
	FCRL3	TCCAGGAACAGAACAGGCCTtaogetggggAATCACGGGGCTGGTGCTCA	0	1	1	32
	BEND2	TGTGCTAGATACCTTATTcagaactcttcACAAAAGATGTCCTGGTCCA	1	0	1	21
	EME1	TACACACAGAGCCGAGCTcaaattgtgcagAGCTGGAAAGAGCTGGCGGA	0	1	1	12
	MOSPD2	TAAGTTAAAGACAATGCTTtcaatagtgcagATAAAACCAAGTGAAGATATA	1	0	1	7.9

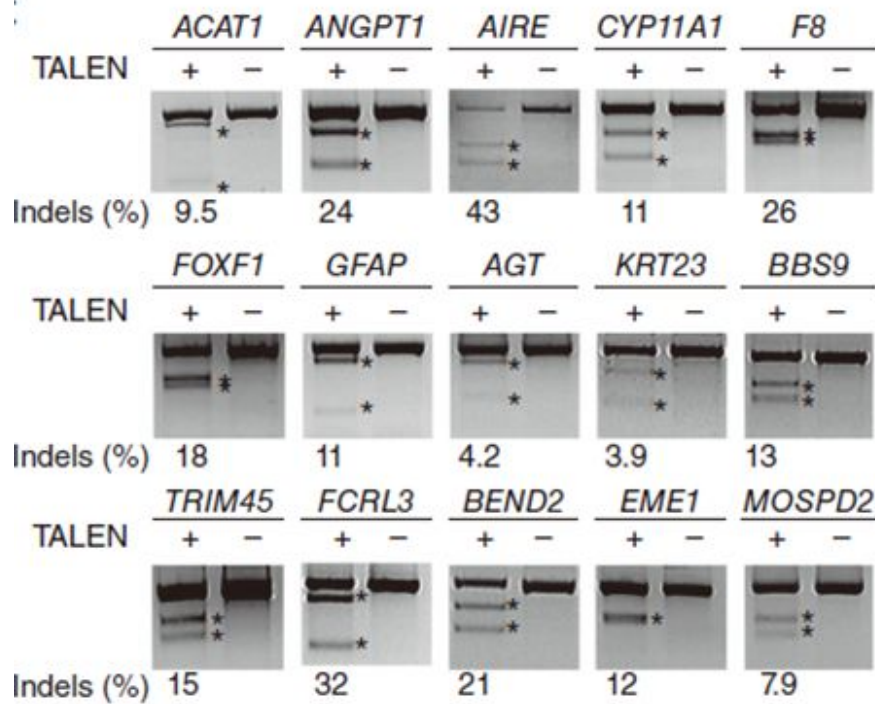


Figure 7. Pilot test of 15 TALENs. TALEN-mediated mutation frequencies were detected using the T7E1 assay. Asterisks indicate the expected positions of DNA bands cleaved by T7E1. The numbers at the bottom of the gels indicate mutation frequencies measured by band intensities.

2. Design of human genome-wide TALENs

To construct a human TALEN library - a genome-scale collection of human gene-targeting TALENs - I obtained the sequences of 18,742 protein-coding genes from HUGO Gene Nomenclature Committee(Seal et al. 2011) (www.genenames.org) on March 2011 and Refseq mRNA database from the National Center for Biotechnology Information (Pruitt et al. 2012) (www.ncbi.nlm.nih.gov) on November 2011. From the coding sequences (CDS) in Refseq mRNA database, we searched for candidate TALEN target sites that consist of two half-sites of 20-bp in length (18.5RVDs, total 40-bp) separated by 12- or 13-bp spacers and that start with the base T and end with A. To identify appropriate TALEN target sites for functionally gene disruption, I sequentially filtered these sites in the following steps:

- 1) I removed TALEN target sites that correspond to exon-exon junctions.
- 2) I selected target sites that reside in common exons in cases in which genes are expressed in two or more splicing variants
- 3) I scored the relative positions of target sites in CDS from the longest mRNA variant and excluded sites that resided within the downstream 30% of each coding sequence.
- 4) I searched for unique target sequences with the minimum number of potential off-target sites in the following manner:
 - i) I screened for potential off-target sites in the human genome sequence (hg19) by allowing 3-bp mismatches to

the left and right half-target sites using Bowtie, an ultrafast memory-efficient short read aligner (Langmead et al. 2009).

- ii) All sequences were aligned at chromosome loci positions for calculating distance between two half-target sites
- iii) I excluded candidate sites that are associated with potential off-target sites that carry 6 or fewer mismatches. I defined potential off-target sites as any homodimeric or heterodimeric half-sites separated by 12- to 14-bp spacers. As a result, the most homologous potential off-target sites carried at least 7-base mismatches with the target site of choice.

These criteria were stringent enough to avoid poor sites that were not appropriate for site-specific gene knockout but were flexible enough to identify multiple sites in most genes (Table 2).

Table 2. Pilot test of computational strategy for TALENs design

Gene ID	Symbol	Initial found sites	off-target filtering		Position in CDS (< 70%) filtering		Common exon filtering	
			numbers	percents	numbers	percents	numbers	percents
672	<i>BRCA1</i>	944	777	82.3	628	66.5	56	5.9
1029	<i>CDKN2A</i>	38	27	71.1	20	52.6	0	0.0
1234	<i>CCR5</i>	140	83	59.3	54	38.6	54	38.6
3043	<i>HBB</i>	39	7	17.9	4	10.3	4	10.3
3651	<i>PDX1</i>	34	25	73.5	23	67.6	23	67.6
4791	<i>NFKB2</i>	195	91	46.7	82	42.1	82	42.1
5728	<i>PTEN</i>	222	59	26.6	56	25.2	56	25.2
5888	<i>RAD51</i>	167	84	50.3	59	35.3	40	24.0
6657	<i>SOX2</i>	56	25	44.6	16	28.6	16	28.6
7157	<i>TP53</i>	121	77	63.6	59	48.8	32	26.4
10661	<i>KLF1</i>	48	40	83.3	32	66.7	32	66.7
29102	<i>DROSHA</i>	548	353	64.4	245	44.7	228	41.6
Total		2552	1648	64.6	1278	50.1	623	24.4

Genome-wide search identified at least one TALEN site that satisfied all of the criteria described in the text in 17,120 out of 18,742 genes (91%) (Group A in Table 3). To identify additional target sites, I loosened the criteria as follows: 1) I considered TALEN pairs that consisted of at least one monomer with 14.5 to 17.5 RVDs (16- to 19-bp, total 32- to 39-bp) and found 2,842 such sites (Group B) in 1,361 genes including 162 additional genes. 2) I searched for target sites whose potential off-target sites elsewhere in the genome carry at least 2-bp mismatches at each of the two half-sites: I identified 706 sites (Group C) in 393 genes including 338 additional genes. 3) In cases in which I failed to find sites that resided within the upstream 70% of a coding sequence, I searched for TALEN sites that resided within the upstream 90%: A total of 270 sites in 154 genes were identified (Group D). 4) 60 Genes encode splicing variants that do not share a common protein-coding exon. For these genes, I searched for sites in exons shared by at least two variants (Group E). 5) Finally, I identified 4,582 sites in the remaining 922 genes after I further loosened the off-target criterion (Group F). I designed to choose all of these additional sites (Groups B to F) preferentially recognized by 18.5 RVD TALENs with minimum off-target effects. As a result, I identified a total of 169,362 TALEN target sites in 18,740 protein-coding genes. I couldn't identify target sites in the rest 2 genes because potential TALEN target sites were not exist in their coding sequences. The TALEN target sites were designed 9.0 target sites per gene on average. The vast majority (98%) of these sites are targeted by 18.5/18.5 RVD

TALENs. Of these sites 95% (160,712/169,362) do not have any homologous sites with >85% sequence identity (that is, 6-base mismatches/ \leq 40-bp sequence) in the human genome (Table 4).

Table 3. Summary of human genome-wide TALEN Library

Group	Target-site length (bp) ^a	TALEN length (RVD)	Position in coding sequences	Number of genes	Number of genes (running total)	Number of target sites	Number of target sites (running total)	Minimum number of mismatches at potential off-target sites
A	40	18.5	upstream 70%	17,120	17,120	160,712	160,712	>6 (=3 + 3) ^b
B	32~39	14.5~18.5	upstream 70%	1,361	17,282	2,842	163,554	>6 (=3 + 3) ^b
C	32~40	14.5~18.5	upstream 70%	393	17,620	706	164,260	>4 (=2 + 2) ^c
D	32~40	14.5~18.5	upstream 90%	154	17,758	270	164,530	>6 (=3 + 3) ^b → >4 (=2 + 2) ^c
E	32~40	14.5~18.5	upstream 70%	60	17,818	250	165,780	>6 (=3 + 3) ^b → >4 (=2 + 2) ^c
F	32~40	14.5~18.5	upstream 70%	922	18,740	4,582	169,362	≤4 (=2 + 2) ^b

^a TALEN target sequence excluding 12- or 13-bp spacers

^b At least 3-base mismatches at each half-site

^c At least 2-base mismatches at each half-site

Table 4. Analysis of TALEN target sites in the human genome.

Number of target sites	Numer of genes	Position in CDS	Total sites	Number of CpG sites	Total sites
1	41	upstream 10%	41,276	0	60,266
2	685	upstream 10% to 20%	33,323	1	48,353
3	245	upstream 20% to 30%	26,808	2	27,525
4	202	upstream 30% to 40%	21,877	3	15,102
5	1,112	upstream 40% to 50%	18,699	4	8,445
6	258	upstream 50% to 60%	14,610	5	4,785
7	329	upstream 60% to 70%	12,529	6	2,692
8	491	upstream 70% to 80%	104	7	1,306
9	701	upstream 80% to 90%	166	8	616
10	14,676	upstream 90% to 100%	0	>8	272
Total genes	18,740	Total	169,392	Total	169,362
Total target sites	169,362				

3. TALEN-mediated genome editing activity

To validate success rate of TAL effector nucleases, I randomly chose 478 TALEN pairs that targeted different genes to validate their mutagenesis activities in HEK293 cells using dual fluorescent reporter assay. As previous study, I measured GFP positive ratio in each TALENs by FACS (Kim et al. 2011). Although all of cells transfected with reporter plasmid were only <1% GFP positive populations, 472 of 478 cells (98.7%) transfected with reporter and TALEN-encoding plasmids were active (detection limit, ~1%) (Figure 8 and Appendix 1).

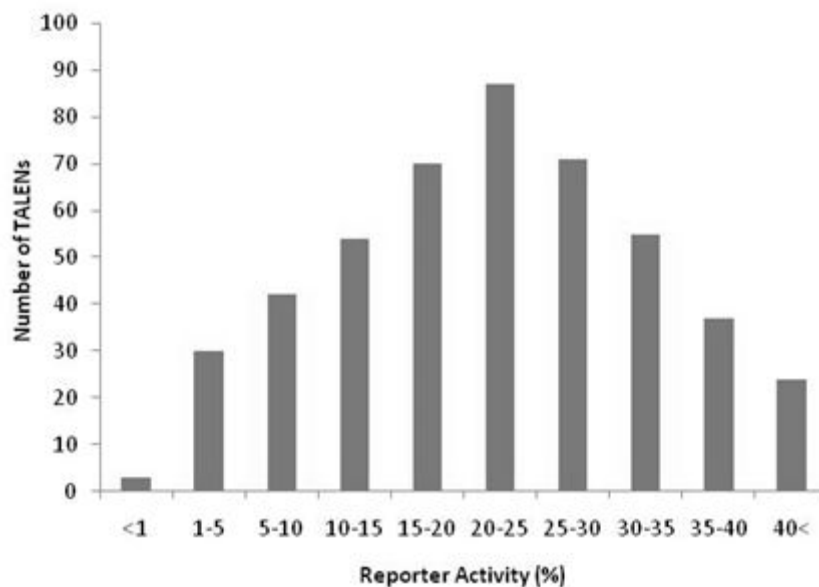
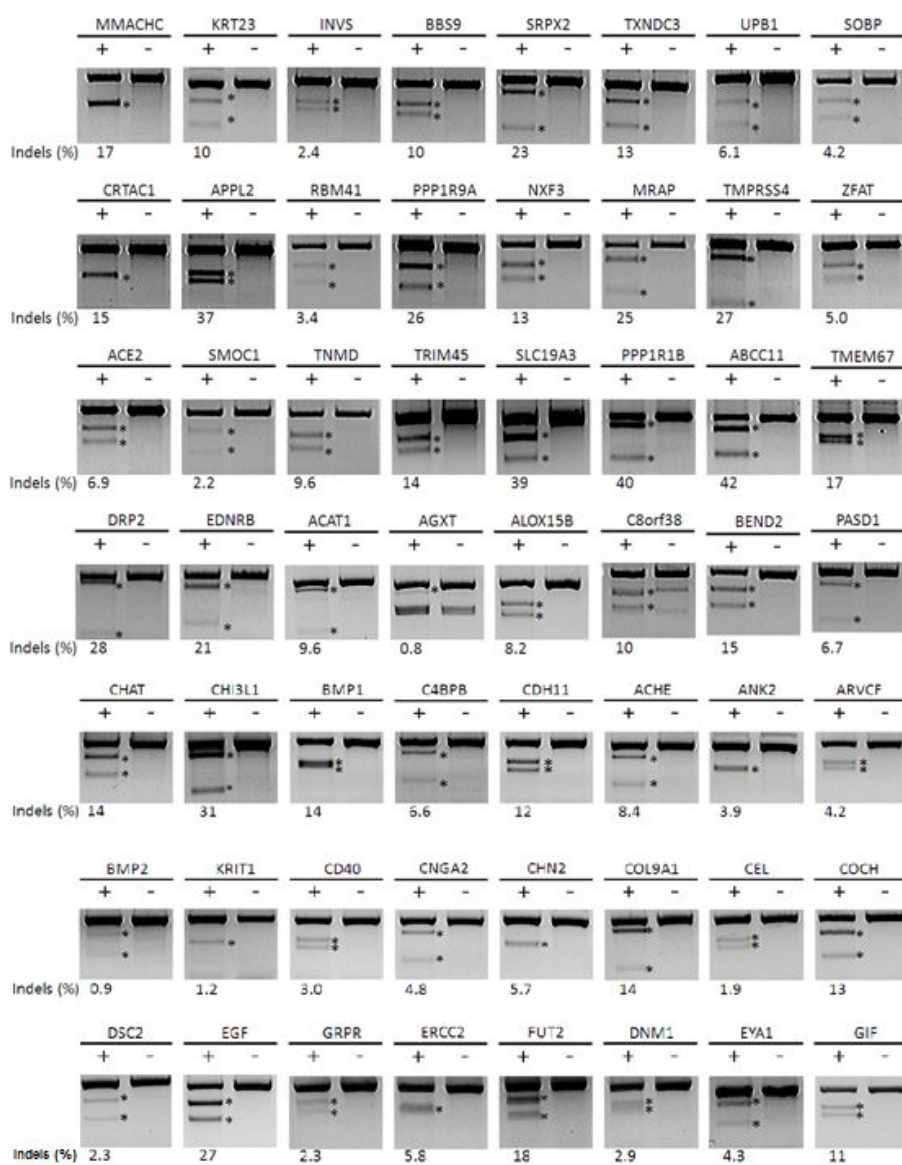


Figure 8. Reporter-based assay for detecting TALEN activities. HEK293 cells were co-transfected with the reporter plasmid including appropriate target site and TALEN-encoding plasmis. Two days after transfection, GFP+ cells were counted by FACS. Almost TALENs were active in this reporter assay.

I investigated whether these active TALENs could, indeed, induce mutagenesis in endogenous target sites highly successful as reporter assay in HEK293 cells. I selected 104 TALEN pairs in 479 TALENs and validated their genome editing activities in HEK293 cells using T7E1 assay (Figure 9 and Appendix 2). Mutations were detected at 101 out of 103 target sites that were successfully PCR-amplified (assay sensitivity, ~0.5%). Accordingly, the success rate of our TALENs was 98.1%. These TALENs were highly active: 78/103 (76%) of TALENs were associated with mutation frequencies of >5%, and 57/103 (55%) of TALENs showed frequencies of >10% (Figure 10). The average mutation frequency was 16%. The two best-performing TALENs, which were specific to the *FKTN* and *OPN1SW* genes, induced mutations at frequencies of 54% and 46%, respectively. Mutations induced by ten different TALENs were confirmed by DNA sequencing (Figure 11). Indels, signatures of error-prone NHEJ repair of DSBs, were detected at the target sites.



(continued)

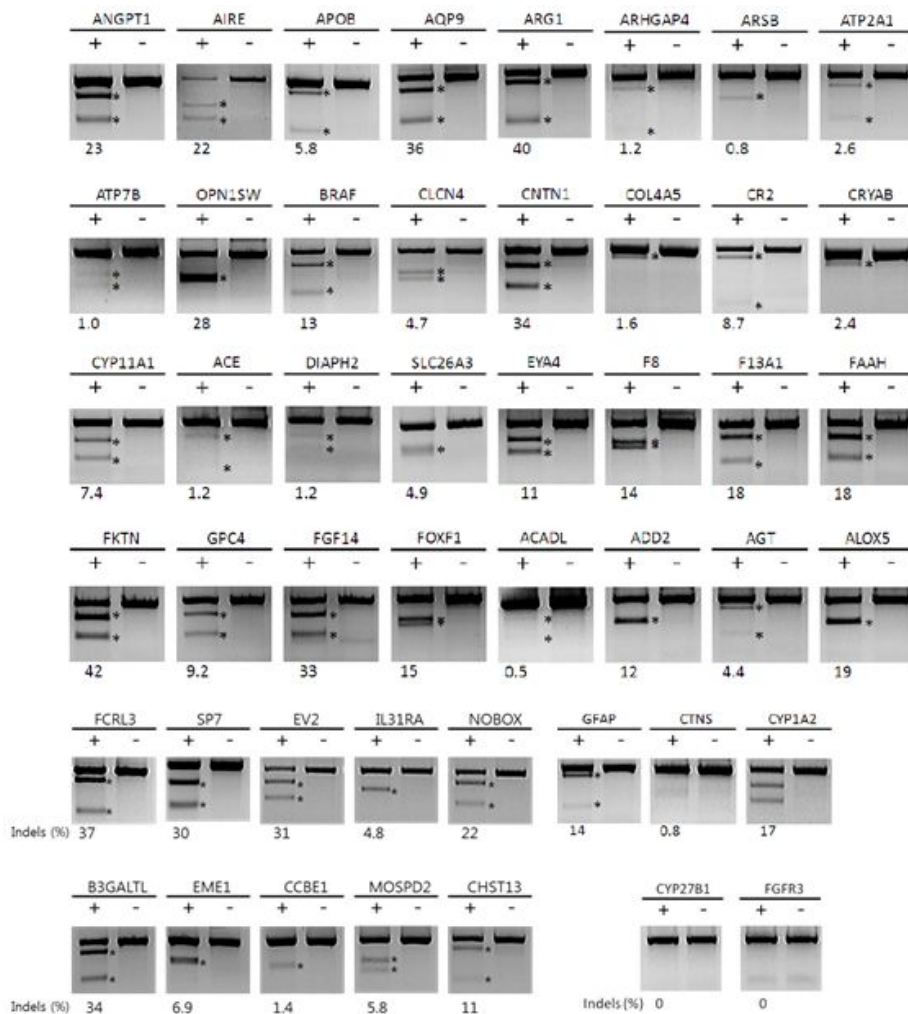


Figure 9. 103 of TALENs test using the T7E1 assay. The genome-editing activities of 103 TALENs were measured by the T7E1 assay. Asterisks indicate the expected positions of DNA bands cleaved by T7E1. The numbers at the bottom of the gels indicate mutation frequencies measured by the intensities of cleaved bands. + or - indicates the presence or absence of TALENs in transfected cells.

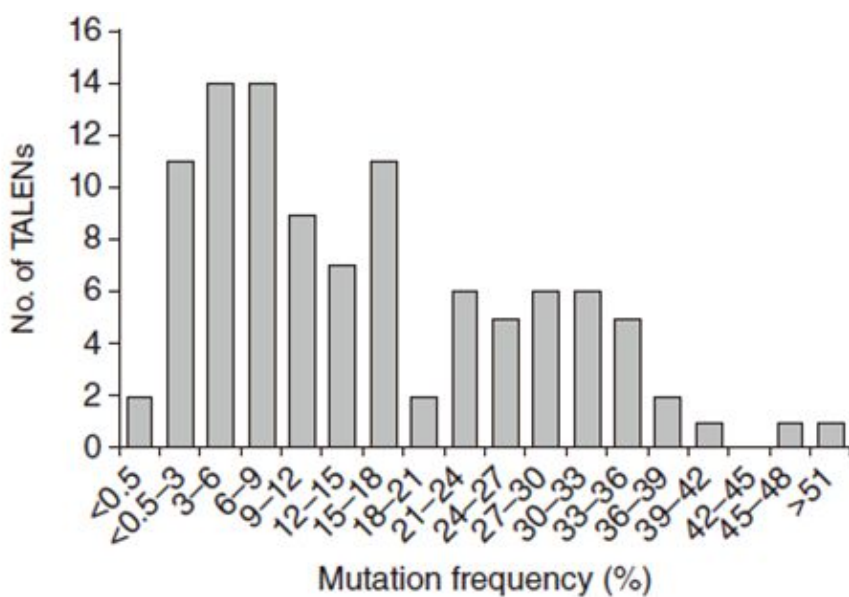
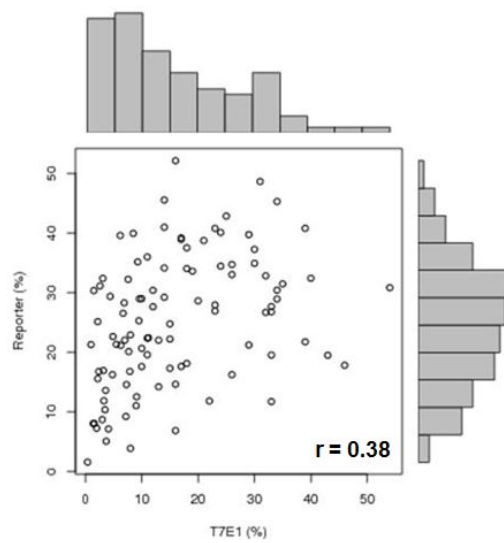


Figure 10. Distribution of TALEN activities. Detected mutation frequencies by the T7E1 assay were analyzed. These TALENs were highly active: 76% (78/103) of TALENs were associated with mutation frequencies of >5% (or indel %), and 55% (57/103) of TALENs showed frequencies of >10%

Figure 11. DNA sequences of indels induced by TALENs. Endogenous mutations induced by 10 different TALENs were confirmed by dideoxy DNA sequencing. The numbers of inserted or deleted bases are shown on the right side of each mutant sequence. WT, wild-type sequence.

Furthermore, to know what is the main factors that influenced on genome editing frequencies of nucleases, I investigated correlation between nuclease mediated episomal reporter activity and endogenous mutation activity of 103 TALEN pairs. The estimated values of linear (Pearson's r) and non-linear (Spearman, ρ) correlation coefficients between endogenous activities and episomal reporter activities were 0.38 and 0.44, respectively (Figure 12a, b). This result showed that nuclease's genome editing activity is little correlate with episomal reporter based activity. In other words, some endogenous factors such as chromatin accessibility or DNA methylation influenced on TALEN's gene disruption frequency.

a



b

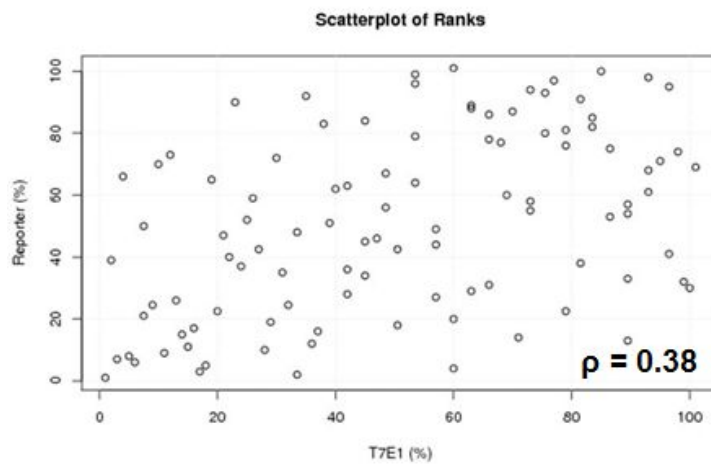


Figure 12. Comparison between episomal reporter assay and the T7E1 assay. The estimated values of linear (Pearson's r) and non-linear (Spearman, ρ) correlation coefficients between endogenous activities and episomal reporter activities were 0.38 and 0.44, respectively.

4. Rescue of inactive TALENs

I investigated why two of the 118 TALENs (including 15 TALENs used in the pilot study), which were designed to target *CYP27B1* and *FGFR3*, failed to show any genome-editing activity in the T7E1 assay in contrast with highly active in reporter assay (Figure 13 and Appendix 1). First, I sequenced the two target sites in HEK293 cells and found no poly-morphisms or mutations at these sites. Thus, the target sites do not contain sequence variations that prevent the TALENs from binding. To resolve this unexpected discrepancy in cleavage of episomal and chromosomal DNA, I hypothesized that these TALENs could not access the endogenous sites owing to the local chromatin structure or DNA methylation. To test this idea, I treated cells with either trichostatin A, an inhibitor of histone deacetylase (HDAC), or 5-aza-2-deoxycytidine (5-aza-dC), an inhibitor of DNA methyltransferase. The two TALENs both induced mutations when cells were pretreated with the inhibitor of DNA methylation but not with the HDAC inhibitor, as shown by the T7E1 assay and by DNA sequencing (Figure 14a, c). In addition, bisulfite DNA sequencing analysis revealed that all the six CpG dinucleotides in each of the two target sites were, indeed, methylated (Figure 14b). These results showed that TALENs cannot recognize heavily methylated DNA and that 5-aza-dC rescued the TALENs that did not cleave chromosomal DNA initially, consistent with recent biochemical and structural studies (Bultmann et al. 2012; Deng et al. 2012; Valton et al. 2012).

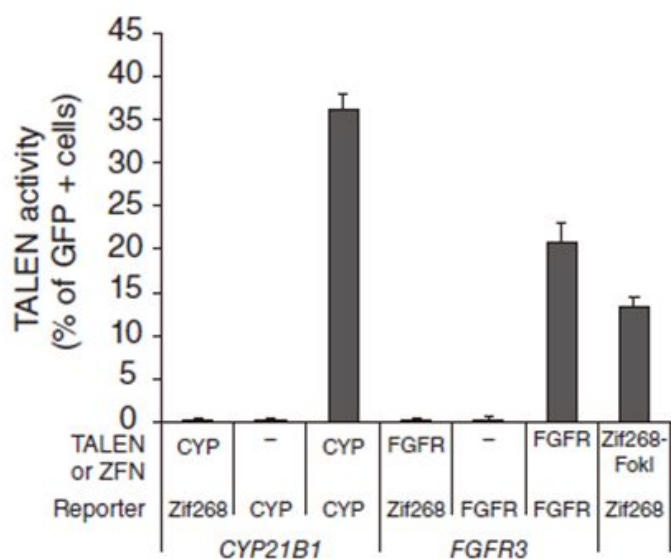
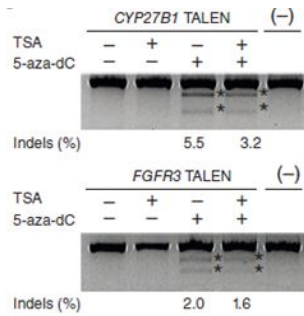


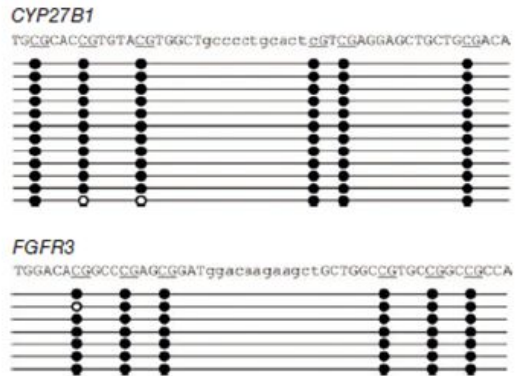
Figure 13. Episomal reporter assays of two genome-inactive TALENs.

Dual fluorescent reporter assays of two TALENs that could not detected mutations by the T7E1 assay. Zif268-FokI (ZFN) was used as a positive control. Error bars, means \pm s.e.m.

a



b



c

CYP27B1

```

ctttgggacagTGGCACCCTGTACGTGGCTgcccctgcaactcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca WT
ctttgggacagTGGCACCCTGTACGTGGCTgcccctgc-----CTGGCACAaggagggaacccggcccagcgctgca Δ6
ctttgggacagTGGCACCCTGTACGTGGCTgcccctgcctgcaactcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca +4
ctttgggacagTGGCACCCTGTACGTGG-----gcaactcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca Δ8
ctttgggacagTGGCACCCTGTACGTGGCTgc-----actcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca Δ6
ctttgggacagTGGCACCCTGTACGT-----cactcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca Δ11
ctttgggacagTGGCACCCTGTACGTGG-----ctgcaactcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca Δ6
ctttgggacagTGGCACCCTGTA-----gcctgca Δ58
ctttgggacagTGGCACCCTGTACGTGGCTgccc-----ctcGTCGAGGAGCTGCTGCGACAaggagggaacccggcccagcgctgca Δ5

```

FGFR3

```

caagcacctcgcccgaggggccccttacTGGACACGGCCCGAGCGGATggacaagaagctGCTGGCCGTGCCGGCCGCCAacaccgtccgct WT
caagcacctcgcccgaggggccccttacTGGACACGGCCCGAGC-----tGCTGGCCGTGCCGGCCGCCAacaccgtccgct Δ15
caagcacctcgcccgaggggccccttacTGGACACGGC-----tGCTGGCCGTGCCGGCCGCCAacaccgtccgct Δ21
caagcacctcgcccgaggggccccttacTGGACACGGCCGA-----GCTGGCCGTGCCGGCCGCCAacaccgtccgct Δ18
tgg-----//-----gcagg Δ241

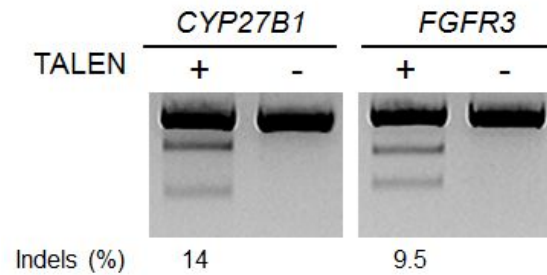
```

Figure 14. TALEN-driven mutations in drug-treated cells. (a) I treated cells with either trichostatin A, an inhibitor of histone deacetylase(HDAC), or 5-aza-2-deoxycytidine (5-aza-dC), an inhibitor of DNA methyltransferase. The two TALENs both induced mutations when cells were pretreated with the inhibitor of DNA methylation but not

with the HDAC inhibitor, as shown by the T7E1 assay. (b) Cytosine methylation at two initially unmodified sites revealed by bisulfite sequencing. The DNA sequences of two half-sites and spacers are shown in upper case and lower case, respectively. CpG dinucleotides are underlined. Closed and open circles indicate methylated and unmethylated cytosines, respectively. (c) Sequencing validation of TALEN induced mutations in 5-aza-dC treated cells.

To avoid treating cell lines with 5-aza-dC, a mutagen, First, I synthesized two new TALENs that HD RVD modules replaced with NG RVD modules in genome-inactive TALENs for binding methylated cytosine as previous study (Deng et al. 2012). The two methylated cytosine binding TALENs (mC-TALENs) were rescued mutagenesis activity (Figure 15a). Next, I tested whether the two genome-inactive TALENs could be replaced with other TALENs that target the same gene. I synthesized new TALENs specific to the *CYP27B1* and *FGFR3* genes and found that these TALENs were able to induce targeted genome modifications at the two new sites (Figure 15b).

a



b

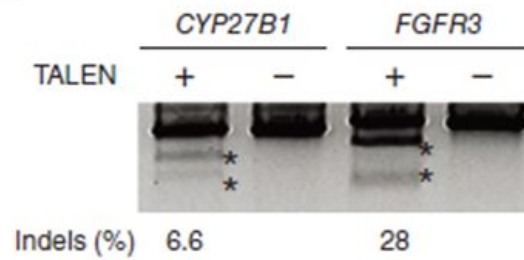


Figure 15. Targeted mutagenesis using alternative sets of TALENs. (a) HD RVDs modules in inactive TALENs were replaced with NG RVDs and transfected to HEK293 cells. The mutation bands were detected in the T7E1 assay. (b) Targeted mutagenesis using a second set of TALENs in TALEN library.

5. Undetectable off-target mutations

Both ZFNs and TALENs can induce mutations at sites other than their intended target sites (Gabriel et al. 2011; Hockemeyer et al. 2011; Mussolino et al. 2011; Pattanayak et al. 2011). As noted above, I chose all the TALEN sites carefully to minimize off-target effects. I investigated whether my designed TALENs could still induce unintended mutations at sites whose sequences are homologous to those of on-target sites. I chose ten TALENs that were highly active, showing mutation frequencies of 30% or greater (36% on average) at their on-target sites, searched for the most likely off-target sites based on the sequence homology in the genome, and tested the genome-editing activities of these TALENs at these potential off-target sites using the T7E1 assay (Table 5). None of these TALENs elicited any measurable mutations at their most likely off-target sites (assay sensitivity, ~0.5%), demonstrating exquisite specificities of these TALENs (Figure 16).

The specificities of RVDs are not absolute. Among the four RVDs, NN appears most promiscuous, as it cannot distinguish G from A efficiently (Boch et al. 2009; Moscou and Bogdanove 2009; Miller et al. 2011). To address this issue, I searched for additional sites of the ten TALENs that would be homologous to the on-target TALEN sites if the promiscuity of NN is considered. Nine TALEN sites were not associated with any potential off-target sites that carry six or fewer mismatches and thus still belong to group A. Only one TALEN site, which contained an unusually high number (17) of guanines in the

40-bp sequence, was associated with additional off-target sites; in this case, there were three potential off-target sites, each of which carries 6-base mismatches (Table 5). The T7E1 assay showed that this TALEN did not induce off-target mutations at these sites (Figure 16). Among the 17,120 group A sites, the vast majority (86%) still satisfy the group A criterion when I take the promiscuity of NN into account. As expected, the other 14% of the sites contain an unusually high number of guanines in their target sequences. These results suggest that we would still choose most of the same sites even when I consider the promiscuity of NN. Because I designed up to ten sites in each gene, potential users who concerned about NN promiscuity may avoid sites that contain too many guanines.

Table 5. Potential off-target sites of highly active TALENs in the human genome.

	Chromosome number	Gene name	Left half-site (5' to 3')	No. of mismatches	Right half-site (5' to 3')	No. of mismatches	Spacer (bp)
ON		APPL2	TTGGCAGACACAAATGGTCT		TGAGATCTTTTCTCGGAAT		12
OFF	16	CTCF	TgGGgachACAgtGGTeCa	6	TGA G cT c CTTTT C agtgAAAT	5	14
ON		SLC19A3	TATGTCCGCTACAAGCCAGT		TAATGATGAAACTGATACCT		12
OFF	7	N/A	T t T t TC C cC a ACAA G CCAGT	4	T a T t aT a AA A CT a ATgC a	6	13
ON		PPF1R1B	TCAGGAGAGGGGGCACCATCT		TGTAGGCACAGGGGTGTGGT		12
OFF	18	ZFPM1	cCAGGgacGcGGCACCATCT	5	gGTgGGC a aAGGGGTGTGG a	4	13
ON		SP7	TTTGGTGGCTCTAGCCCTCT		TGCCCTGCCTTGGCCAGAGTT		12
OFF	8	STK3	TTTGGA G cC a aaAGCCCTCT	5	T a cCTG C aT a G c aCAGAG a T	5	13
ON		B3GALT1	TCATCCAGAGTCAAAGTAAT		TTTAACTGCTCTGCTCTCTT		12
OFF	19	N/A	TCATCCAG t GTCAA A gggt	4	TgT c ACTGCTGT T CTC a CTT	5	13
ON		FKTN	TGAGCTGTGTGTGTCAAAT		TGTGCTATCAAATCCAACTT		12
OFF	13	N/A	TG G c t GT c TTTGT a t a AT	5	TGT c CTAT a AA A CCAA T T a	4	13
ON		FAAH	TCAGAGGCCTGCTCTAGCCCT		TGTGTAACTCTTGACCCAGC		12
OFF	8	RLBP1L1	TCh a AG a GA G cAGCTgCCCG	5	TGTGTggCT c C a CA C CCAG a	6	14
ON		ABCC11	TACTCTCCAAAGATGGCCCT		TGGAGCCTCAGGAATTTCTCT		12
OFF	5	N/A	TgCTCTCC c h a cTG G CCgCT	5	c a GA C CTCAGGA a T a CTCT	5	13
ON		EVC2	TGCTCTGTCTGGACAGCATT		TTTCCACAGAGTCCCACAT		12
OFF	4	N/A	TGCTCT T cC a GGAT a C a CTT	6	TTT a c a C a AGAGTCCCG a AT	3	12
ON		MMACHC	TGTAGCTGGGGCTGCTTACT		TGGGTCA G CTCTCCACATCTT		12
OFF	7	GRB10	T a agGCTGTGG t GTCTTAC a	4	T G c c TCA G t t CTC c ATCTg	6	13
ON		PPF1R1B	TCAGGAGAGGGGGCACCATCT		TGTAGGCACAGGGGTGTGGT		12
OFF	6	N/A	aCAAGAA a AAAA c CCATCT	(3)	TGTAA G CA C AA a h a CTG a CT	(3)	13
OFF	3	MFN1	TC A AA A GA A h A GA C CA a ag	(3)	TGT a c a CA C cAGGATAG a c	(3)	13
OFF	2	ERBB4	TC A AA A G c AA A h a C a CATCT	(3)	TGTAA G CA C AA a h a h a h a CT	(3)	12

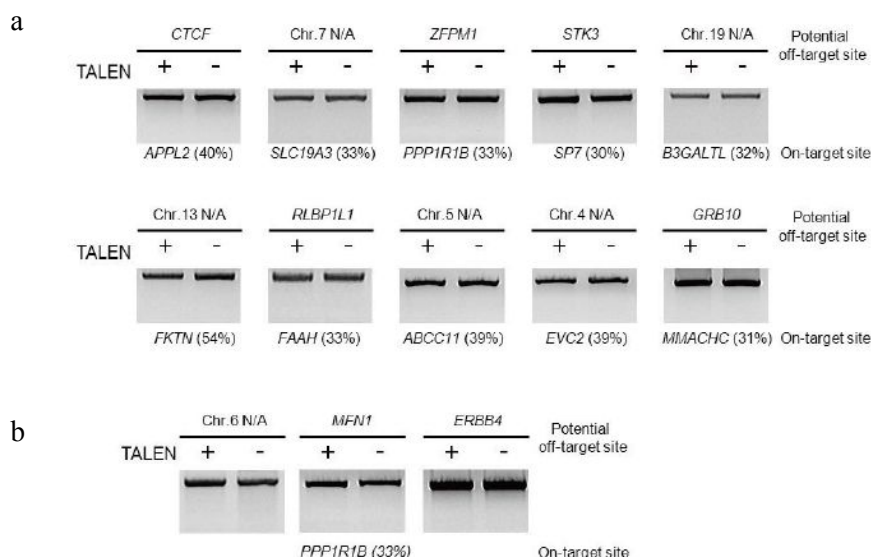


Figure 16. Undetectable off-target mutations with TALENs. (a) The genome was searched *in silico* for potential off-target sites (listed in Table 6) of 10 highly active TALENs, and the genome-editing activities of TALENs at these sites were tested by the T7E1 assay. Mutation frequencies at on-target sites are indicated at the bottom of the gels. (b) To account for the promiscuity of the NN RVD, additional sites were searched *in silico* that would be homologous to the on-target sites of the 10 highly active TALENs if all adenines were changed to guanines in the human genome. Only the *PPP1R1B*-specific TALEN was associated with any potential off-target sites (3 sites, listed in Table 6). No mutations were detected at these sites.

6. TALEN induced genome rearrangement

In many cases, homologous genes are often clustered in a chromosome. Many of these homologous genes are redundant in their functions. A single gene knockout often fails to cause any phenotypic changes owing to these redundancies. I propose that our TALENs could be used to induce large chromosomal deletions in a targeted manner to remove a cluster of homologous and functionally redundant genes. The TALEN library reported here consists of 18,740 TALEN pairs whose target sites are distributed all across the human genome. Thus, two neighboring TALEN target sites are separated by 170 kbp, on average, in a chromosome.

As a proof of principle experiment, I investigated whether and how efficiently two TALEN pairs can induce large chromosomal deletions in human cell lines. Two combinations of TALEN pairs were tested: (i) a *BRAF*-specific TALEN and a *NOBOX*-specific TALEN and (ii) a *EDNRB*-specific TALEN and a *FGF14*-specific TALEN (Figure 17a). The two target sites of the first pair are separated by 3.6 Mbp on chromosome 7, and those of the second pair are separated by 24 Mbp on chromosome 13. Targeted deletions of these large chromosomal segments were detected by PCR and confirmed by DNA sequencing (Figure 17b, d). The deletion frequencies measured by limiting dilution and PCR detection were 0.6% (*BRAF* and *NOBOX*) and 0.4% (*EDNRB* and *FGF14*) (Figure 17c), which are in line with those obtained with ZFNs (Lee et al. 2010). Recently, two TALEN pairs were used to

induce a 7-kbp chromosomal deletion within a single gene in porcine fibroblasts at a frequency of 10% (Carlson et al. 2012). Apparently, large deletions of Mbp by TALENs are at least tenfold less efficient than are those of kbp.

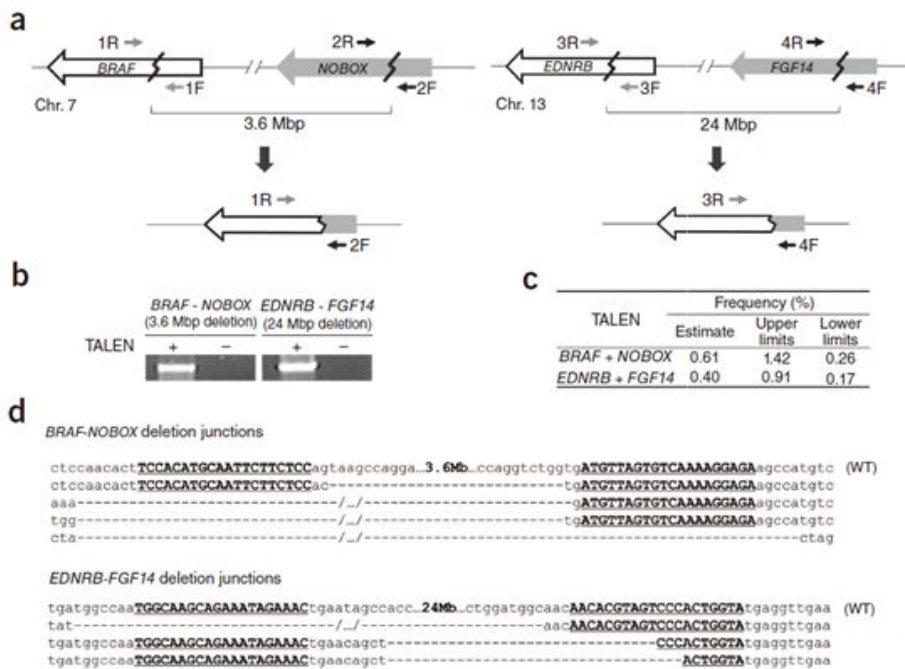


Figure 17. TALEN-mediated targeted genomic deletions. (a) Scheme of targeted chromosomal deletions induced by TALENs. Zigzag lines indicate TALEN target sites. Small arrows are PCR primers. (b) PCR products corresponding to large chromosomal deletions. (c) Deletion frequencies measured by dilution PCR. (d) DNA sequences of deletion junctions. TALEN recognition sequences are underlined and shown in boldface.

In addition, the large collection of TALENs reported here could also be used to induce other types of genome rearrangements, which include translocations, inversions, and duplications (Figure 18-20). Indeed, I was able to induce these structural variations in a targeted manner using various combinations of TALENs at frequencies ranging from 0.01% to 0.7. Structural variations in individual human genomes contribute to genetic diversity and often are associated with genetic diseases and cancer. Although thousands of structural variations have been identified, the biological consequences of these variations are largely unexplored. Previous studies have reported that the repair of two concurrent DSBs in the genome induced by ZFNs gives rise to these variations in a targeted manner in cultured human cells (Lee et al. 2012). To study their functions, one could also use TALENs to induce specific variations in cell lines. The high density of our TALENs all across the human genome could facilitate the creation of cell lines whose genomes are custom-designed.

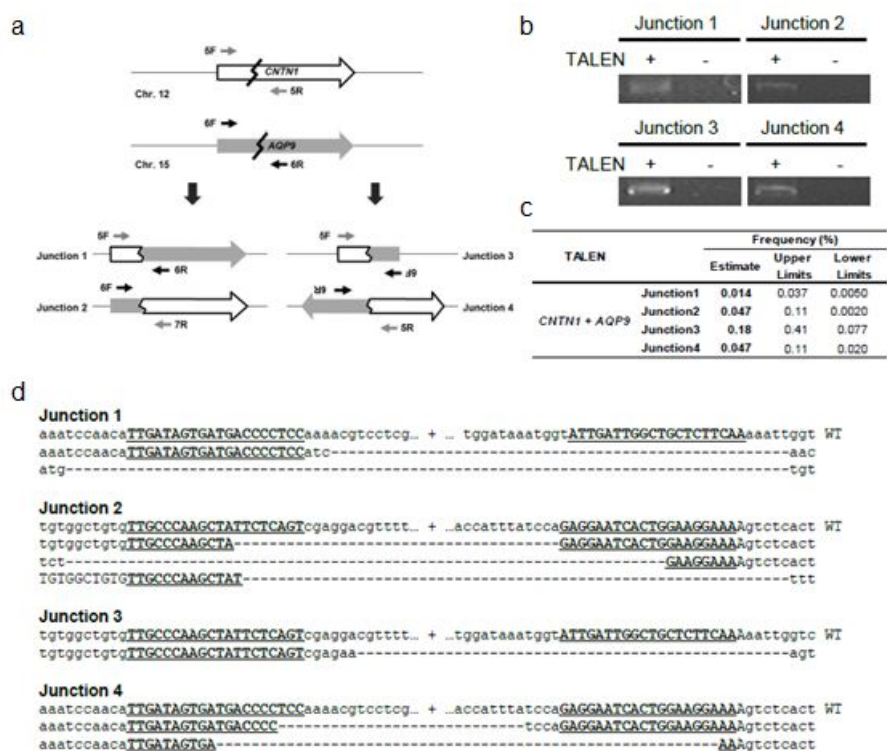


Figure 18. TALEN-mediated targeted genomic inversions. (a) Scheme of targeted chromosomal inversions induced by TALENs. Zigzag lines indicate TALEN target sites. Small arrows are PCR primers. (b) PCR products corresponding to chromosomal inversions. (c) Inversion frequencies measured by dilution PCR. (d) DNA sequences of inversion junctions. TALEN recognition sequences are underlined and shown in boldface.



Figure 19. TALEN-mediated targeted genomic duplications. (a) Scheme of targeted chromosomal duplications induced by TALENs. Zigzag lines indicate TALEN target sites. Small arrows are PCR primers. (b) PCR products corresponding to chromosomal duplications. (c) Duplication frequencies measured by dilution PCR. (d) DNA sequences of duplication junctions. TALEN recognition sequences are underlined and shown in boldface.

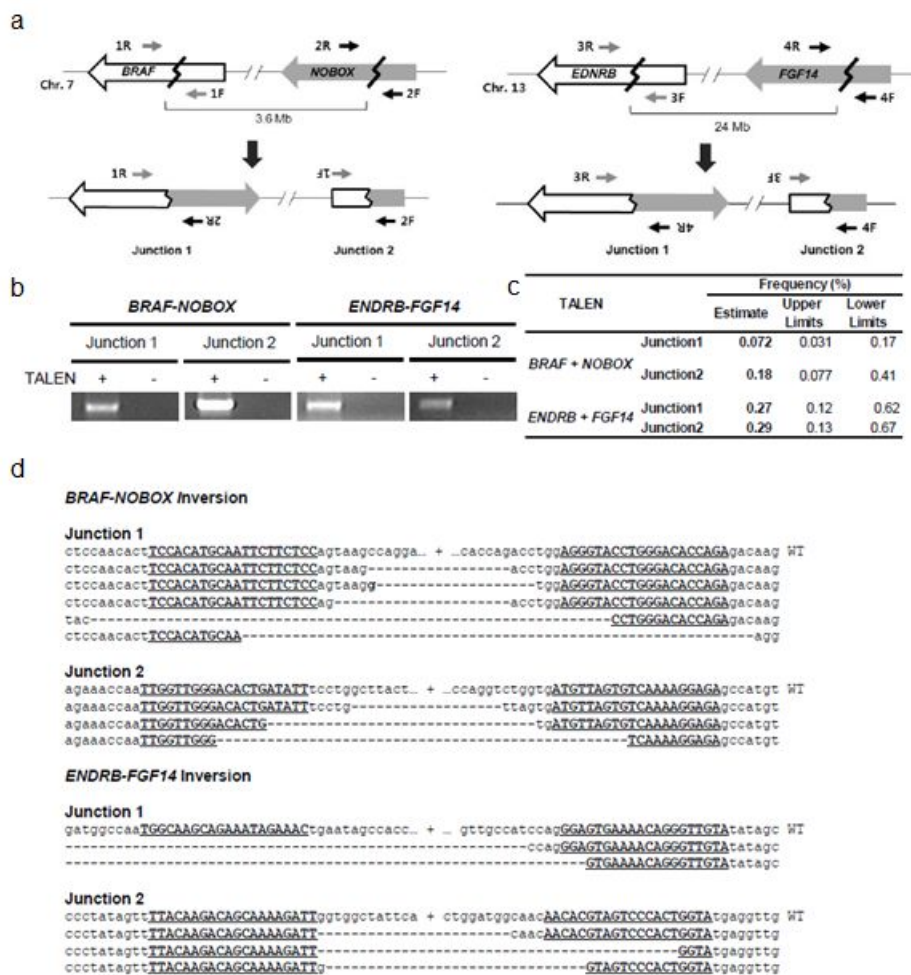


Figure 20. TALEN-mediated targeted genomic translocations. (a) Scheme of targeted chromosomal translocations induced by TALENs. Zigzag lines indicate TALEN target sites. Small arrows are PCR primers. (b) PCR products corresponding to chromosomal translocations. (c) Translocation frequencies measured by dilution PCR. (d) DNA sequences of translocation junctions. TALEN recognition sequences are underlined and shown in boldface.

C. TALEN-mediated knockout cell lines

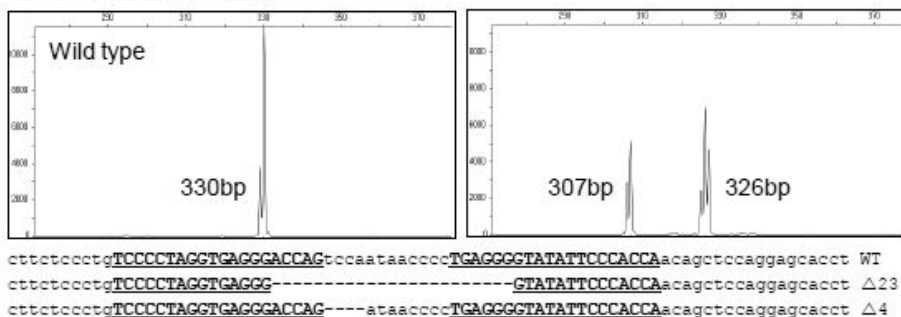
1. Establishment of knockout cell lines for NF- κ B pathway study

To investigate functional gene knockout study in human cells, the cell lines, a perpetuating strain of cells in laboratory culture, have commonly used. But most model cell lines such as HEK293 cells and HeLa cells contain more than three copies of most chromosomes (Macville et al. 1999; Bylund et al. 2004). Thus, the disruption of one or two alleles does not yield gene-knockout cell lines. In addition, some engineered nucleases are cytotoxic (Cornu et al. 2008). Thus, cells that contain nuclease-induced mutations often die out, making it difficult to isolate rare clonal populations of gene-knockout cells (Kim et al. 2009).

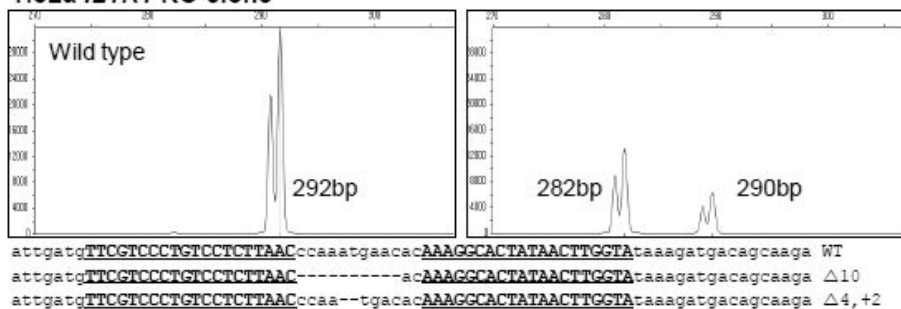
As a proof-of-principle experiment for creating gene-knockout cells, I used TALENs to disrupt a series of genes associated with NF- κ B, a transcription factor linked to inflammation, septic shock, apoptosis, oncogenesis and DNA repair (Perkins 2007; Volcic et al. 2012). I transfected HEK293 cells or HeLa cells with both TALEN plasmids and surrogate reporters, which enabled enrichment of cells in which a gene of interest is completely disrupted (Appendix 2) (Kim et al. 2011; Kim et al. 2013a). After mutant clones were isolated using T7E1 assay, I inspected clonal populations of cells using fluorescent PCR for confirming disruption of all alleles that include gene of

interests. Gene-knockout cells produced amplicon peaks corresponding to indels but not the wild-type peak. Note that all three alleles were mutated in a clonal population of cells. DNA sequencing confirmed the presence of frameshift mutations at target loci in these cells (Figure 21).

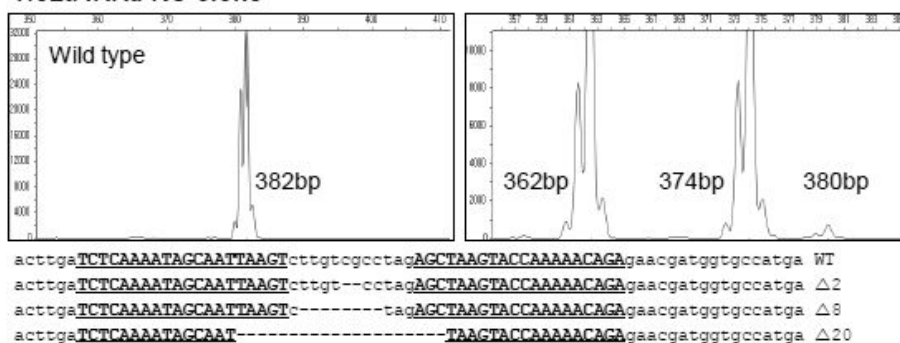
HeLa *TNFR1* KO clone



HeLa *IL1R1* KO clone

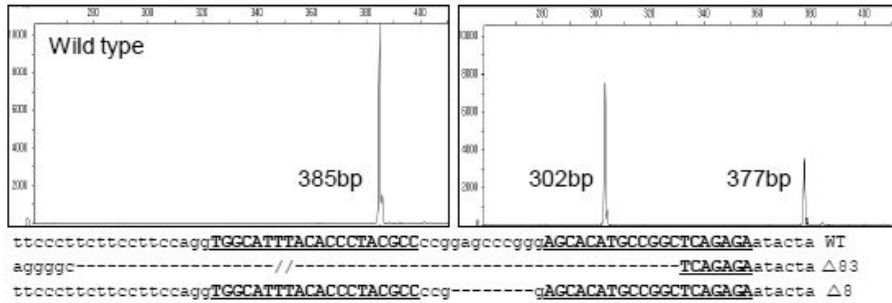


HeLa *IKK α* KO clone

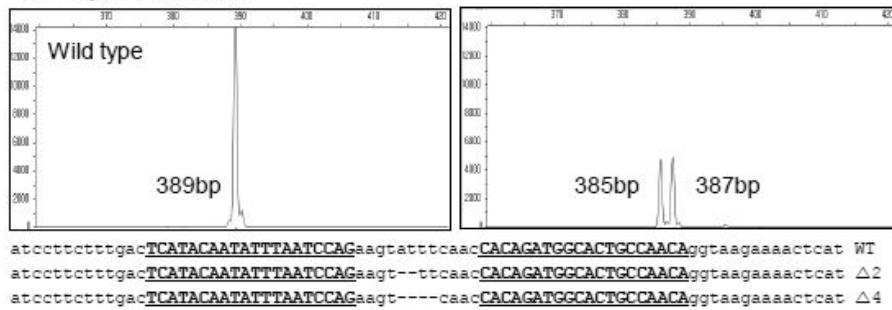


(Continued)

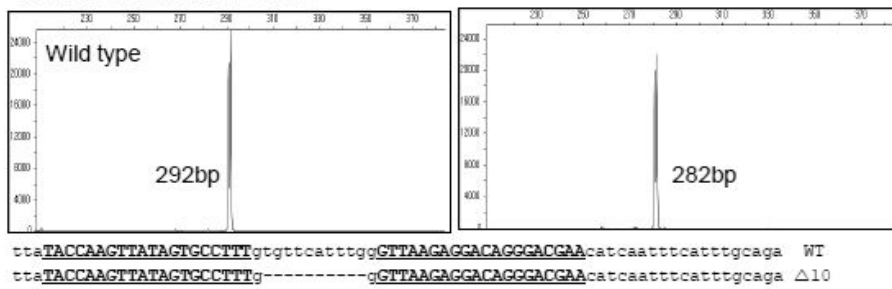
HeLa *TNFR2* KO clone



HeLa *p53* KO clone

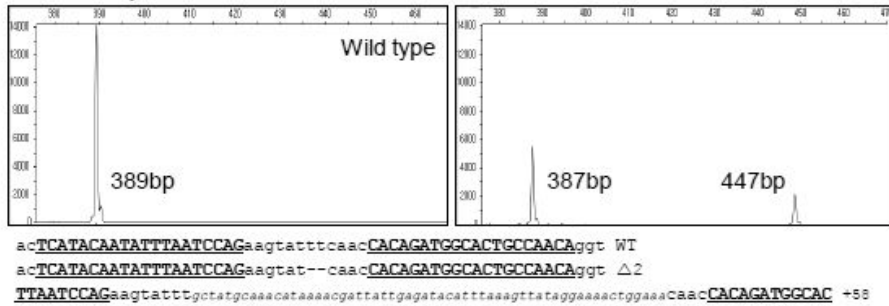


HEK 293 *IL1R1* KO clone

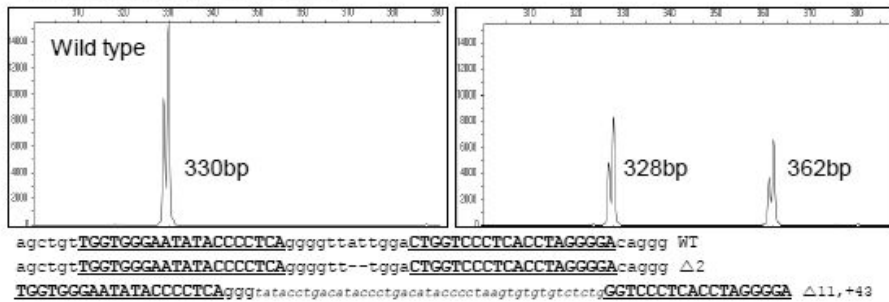


(Continued)

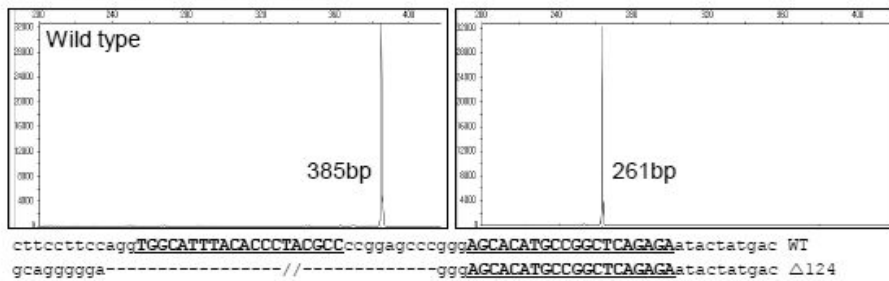
HEK 293 *p50* KO clone



HEK 293 *TNFR1* KO clone



HEK 293 *TNFR2* KO clone



(Continued)

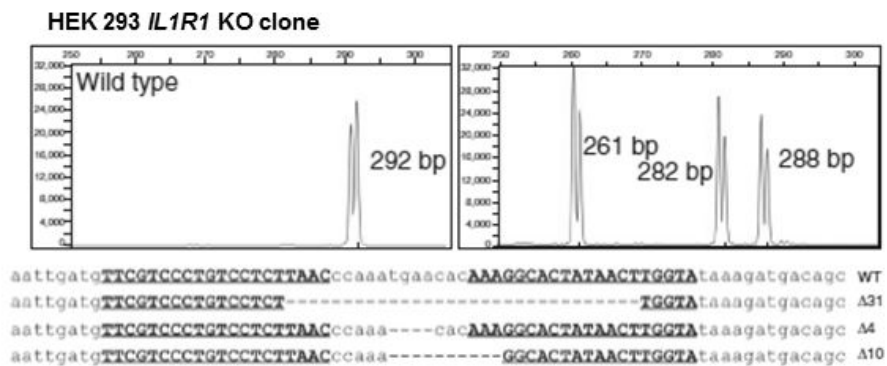


Figure 21. Validation of knockout cell lines. Genotypes of single clones were validated by F-PCR and sequencing analyses. TALEN recognition sequences are underlined and shown in boldface. The numbers of deletions and insertions are indicated. WT, wild-type.

I searched potential off-target sites that are highly homologous to the target site, and no mutations were detected using T7E1 assay (Table 6 and Figure 22). And then, to validate that target protein were not expressed in gene-knockout cells, I performed western blot analysis and RT-PCR (Figure 23a, b). This result showed that these knockout clones are functionally gene disruptions without no mutations at other potential off-target sites.

Table 6. Potential off-target sites of TALENs used for knockout cell lines

	Chromosome number	Gene name	Left half-site (5' to 3')	No. of mismatches	Right half-site (5' to 3')	No. of mismatches	Spacer (bp)
ON		<i>TNFR1</i>	TGGTGGGAATATACCCCTCA		TCCCTAGGTGAAGGACCAG		12
OFF	5	<i>FAM169A</i>	TGcTGGcAATggAACCtCaCA	6	TCCCTAGGaaAAGGctCCAG	4	14
OFF	16	N/A	cGGTaGaaAaAgAACCCTCt	6	TggCCatGGTGAAGGcCCAG	5	12
OFF	3	N/A	TGGTaTgtATtTACCCCTtA	5	atCCCaGgGaGAGGGA CtAG	6	13
OFF	5	N/A	TGGTGaGgAcATgCtCCTCA	5	TCCaCcAtGTGAAGaCaCAG	6	13
OFF	6	N/A	TcGTGGGAtcAgAgCCCcCA	6	TCCaCTAGGTGgtGccCCAAG	5	13
ON		<i>IL1R1</i>	TTCGTCCCTGTCTCTTAAC		TACCAAGTTATAGTGCCTTT		12
OFF	6	N/A	aTctTCCCTGgCCcCTTgAc	5	atCCTAGTTATAAtGgCTTT	5	13
ON	3	<i>VPBP</i>	TgCGaaCCTGcCCTgCTAAC	6	aaCaCAGTTATAccctCCTTT	5	12
OFF	2	N/A	TTgGTgCtTcTCCTCaTAAc	5	TctCAhaTTATAaTGaCgTT	6	13
ON	7	N/A	TTgcTCcCTGTCTCTgAAC	5	TcCCAhaaTaaAGTGCCTTg	5	14
OFF	4	N/A	TTtcTCtCTcTCCTCTcAAC	5	TACCcAtTTcTAGTtCtTTg	6	14

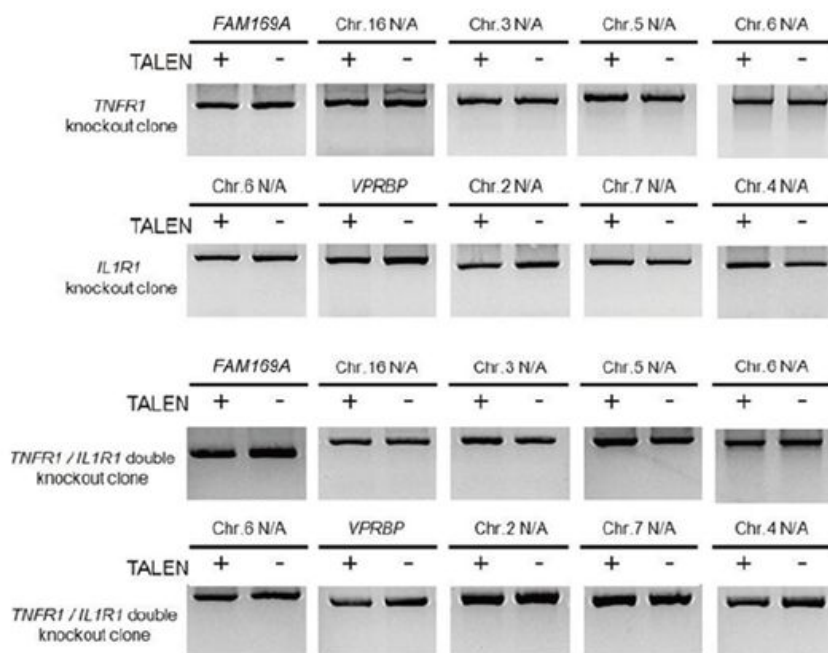
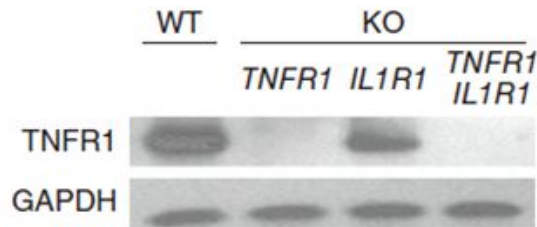


Figure 22. Undetectable off-target mutations in TALEN-mediated knockout cell lines. The top 5 potential off-target sites of *TNFR1* and *IL1R1* TALENs were searched for *in silico*. The T7E1 assay was used to investigate whether these sites were mutated in clonal populations of single- and double-gene knockout cells. No mutations were detected at these sites.

a



b

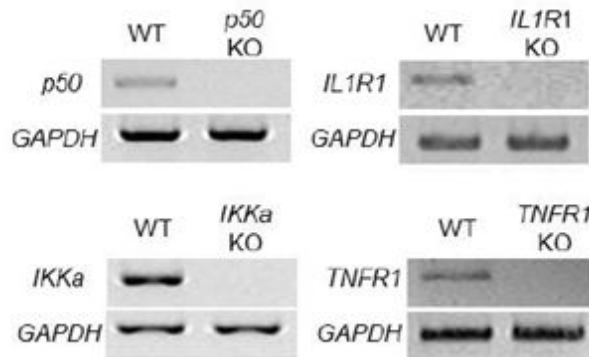


Figure 23. Undetectable gene expression level in gene knockout cells.

(a) No protein expressions in single- and double-gene-knockout cells by Western blotting. (b) Total RNA was extracted from each knockout cell line, and cDNA was synthesized using reverse transcriptase. Target mRNA levels were analyzed by PCR using appropriate primers (Kim et al. 2013b). *GAPDH* RNA levels were used as an internal control.

2. Episomal reporter assay

For the study of NF- κ B pathway study with these knockout cell lines, I constructed luciferase reporter plasmid that inserted multi-copy of p65 binding site in upstream of luciferase gene (Figure 24). Luciferase was regulated by induction of NF- κ B signaling pathway such as treatment of tumor necrosis factor alpha (TNF α) or interleukin-1 beta (IL-1 β), two well-known cytokines that activate NF- κ B signaling. First, we compared TALEN-mediated gene knockout with siRNA-mediated gene knockdown (Figure 25a). Two validated siRNAs, each specific to one of the two receptor genes, *IL1R1* or *TNFR1*, only partially suppressed the NF- κ B signaling by IL-1 β or TNF α , respectively. Thus, the two cytokines still activated the NF- κ B-dependent reporter activities in siRNA-transfected wild-type cells, compared to cells not treated with cytokines ($P < 0.01$, Student's *t*-test). In contrast, both *IL1R1* and *TNFR1* knockout cells showed complete suppression of IL-1 β - and TNF α -mediated NF- κ B signal transduction, respectively. In addition, the two cytokines did not activate the luciferase reporter in double-knockout cells in which both *IL1R1* and *TNFR1* were disrupted. Transfection of the *IL1R1* cDNA into *IL1R1* knockout cells restored IL-1 β -mediated signal transduction, demonstrating that the gene-knockout cells were not impaired in NF- κ B signaling (Figure 25b). Unexpectedly, over-expression of *TNFR1* in *TNFR1* KO cells activated NF- κ B-dependent reporters even in the absence of TNF α . Similar results have been reported by previous study

(Gaeta et al. 2000).

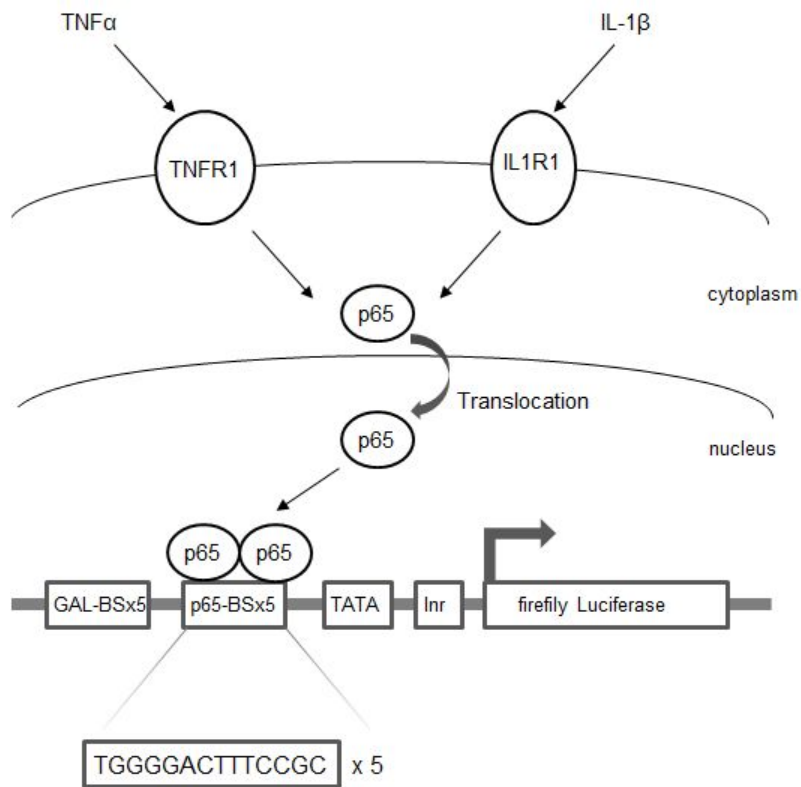
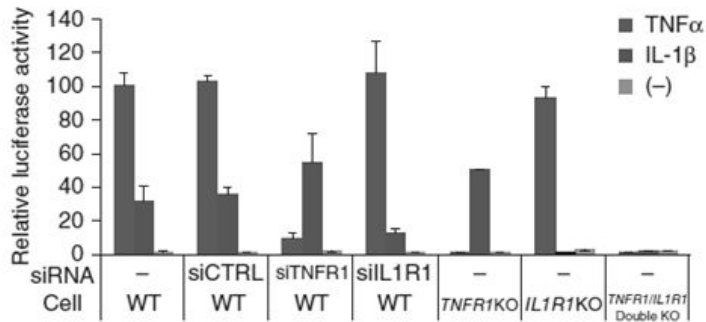


Figure 24. Schematic of reporter assay related NF- κ B signaling.

Gene-knockout or wild-type cells were co-transfected with the luciferase reporter plasmid and the Renilla luciferase plasmid. After 24 h of incubation, cells were treated with TNF α or IL-1 β and incubated for 15 h. Cells were lysed and detected luciferase activities.

a



b

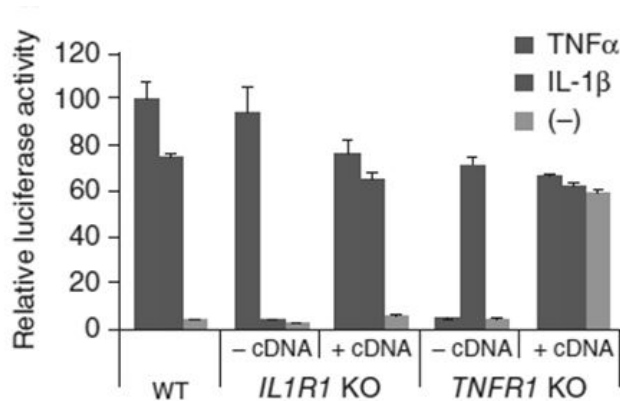


Figure 25. Functional assay in knockout cells using reporter.

(a) Comparison of TALEN-mediated gene knockout (KO) with siRNA-mediated gene knockdown. (b) Transfection of cDNA into gene-knockout cells restores cytokine-dependent NF- κ B activation. The *IL1R1* or *TNFR1* cDNA was transfected into *IL1R1* or *TNFR1* knockout cells, respectively. Note that the overexpression of the *TNFR1* cDNA in *TNFR1* knockout cells induced NF- κ B activation even in the absence of TNF α .

Next, mitoxantrone (Novantrone), an anti-cancer drug that induces nonspecific DSBs in cells, strongly activated the reporter gene in the single- and double-gene-knockout cells, demonstrating that NF- κ B activation triggered by DNA damage is independent of signaling through the two cytokines (Figure 26). These gene-knockout cells can be used for dissecting NF- κ B pathways and screening for reagents or factors that selectively modulate NF- κ B activation independent of signaling via IL-1 β or TNF α or both.

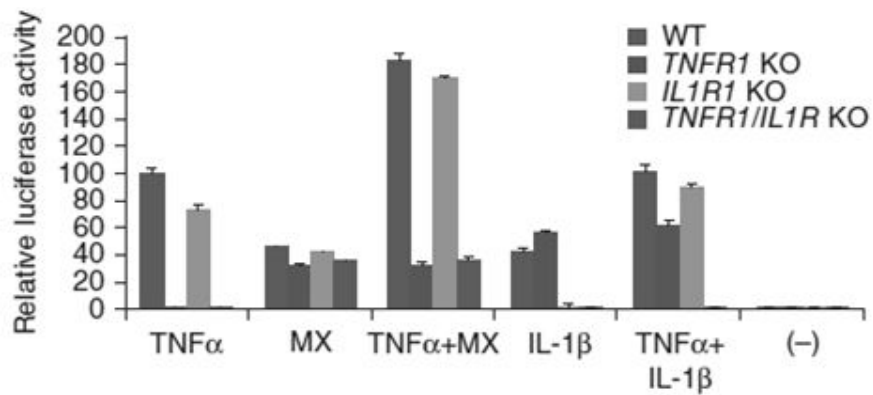
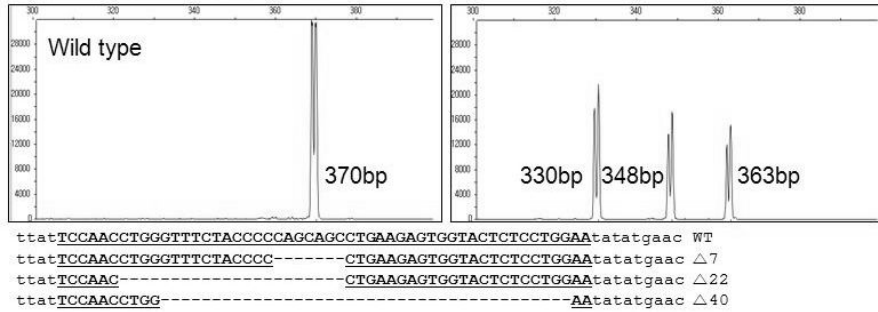


Figure 26. Cytokine-independent NF- κ B activation by mitoxantrone. Gene-knockout or wild-type cells were co-transfected with the luciferase reporter plasmid and the Renilla luciferase plasmid. After 24 h of incubation, cells were treated with $TNF\alpha$ / $IL-1\beta$ /MX and incubated for 15 h. Cells were lysed and detected luciferase activities.

Furthermore, to demonstrate the usefulness of TALEN-mediated gene knockout further, we disrupted the *NR1H4* and *SMEK1* gene in HEK293 cells (Figure 27). These genes had been identified as a potential mediator of NF- κ B activation in a recent genome-wide siRNA screen (Gewurz et al. 2012). I prepared knockout cell lines with TALENs and surrogate reporters and confirmed by fluorescent PCR and Sanger sequencing. Both IL-1 β and TNF α strongly activated the luciferase reporter gene in three or two independent knockout clones, respectively (Figure 28,29). To validate expression level of each genes, I performed RT-PCR and It appears that *NR1H4* is not expressed even in wild-type HEK 293 cells. In case with *SMEK1* gene, gene expression level was detected in wild-type cells but not detected in two independent knockout cell lines. These results showed that these genes were a false positive in siRNA screen.

a

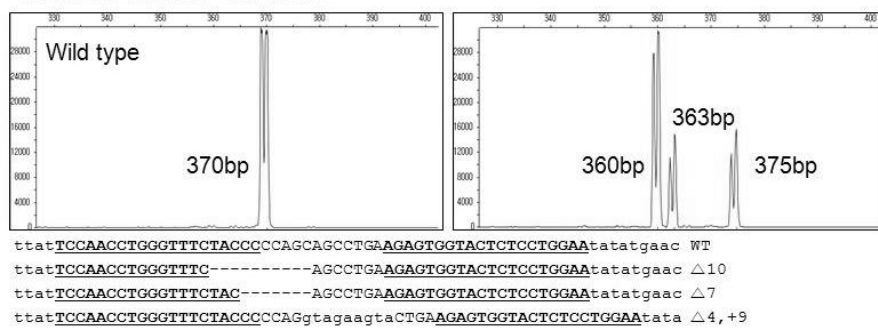
HEK 293 NR1H4 clone #1



HEK 293 NR1H4 clone #2



HEK 293 NR1H4 clone #3



(Continued)

b

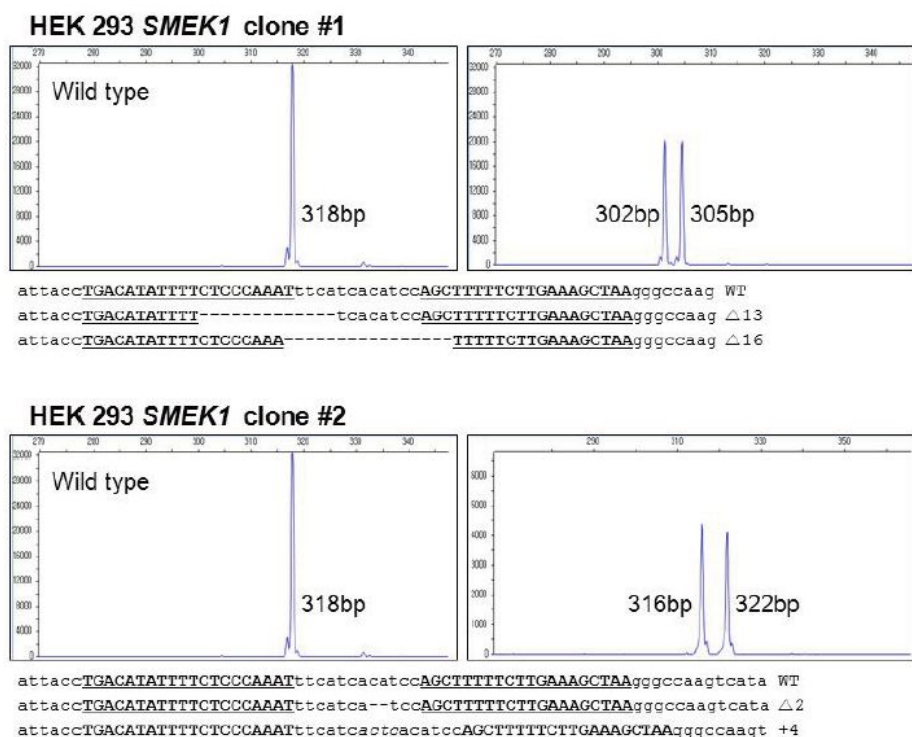
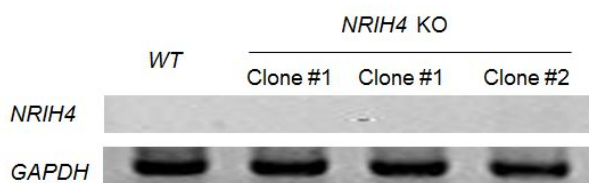


Figure 27. Validation of knockout cells related siRNA screening. Genotypes of single clones were validated by F-PCR and sequencing analyses. (a) NR1H4 gene, (b) SMEK1 gene.

a



b

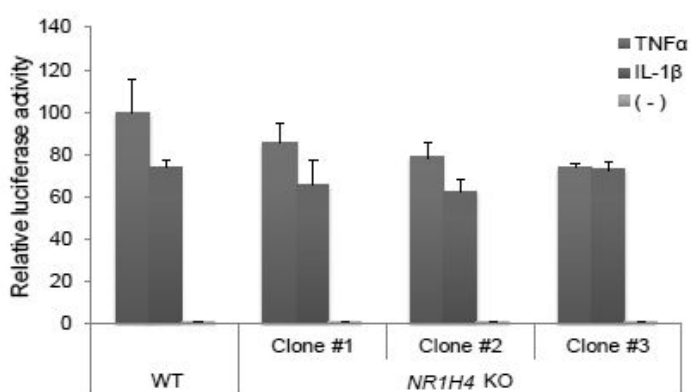
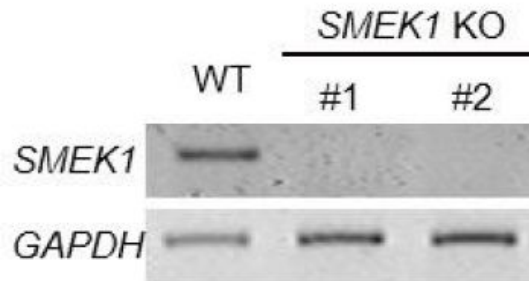


Figure 28. Invalidation of a false-positive gene identified in a genome-wide siRNA screen: *NR1H4*. (a) *NR1H4* gene was not expressed in HEK 293 cells. Total RNA was extracted from each knockout cell line, and cDNA was synthesized using reverse transcriptase. Target mRNA levels were analyzed by PCR using appropriate primers. (b) Three independent clones of *NR1H4* knockout cells were treated with cytokines, and the luciferase reporter activities were measured. Error bars indicate S.E.M. from at least three independent experiments

a



b

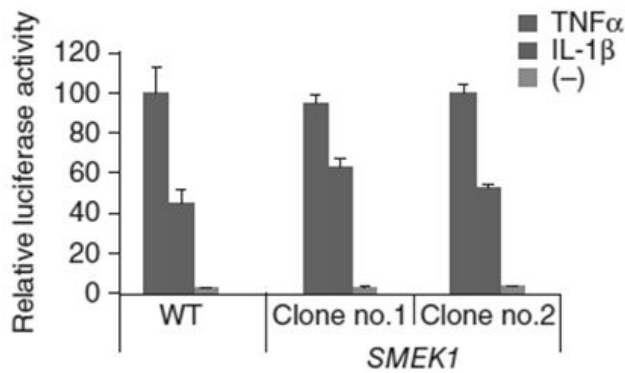


Figure 29. Invalidation of a false-positive gene identified in a genome-wide siRNA screen: *SMEK1*. (a) No expression of *SMEK1* gene in two independent *SMEK1* gene-knockout clones. (b) Two independent clones of *SMEK1* knockout cells were treated with cytokines, and the luciferase reporter activities were measured. Error bars indicate S.E.M. from at least three independent experiments.

D. Expansion of TALEN library

1. Design strategy of TALEN library

In described above, because human genome-wide TALEN library was designed for gene disruption, they do not used other functional studies such as gene correction, targeted transcript disruption, specific cancer related oncogene mutations. For the expanded use of TALEN library, I designed TALEN pairs for every exon for human protein coding genes. I obtained the DNA sequences of defined human protein coding genes from the Ensembl project in April 2013 (Birney and Ensembl 2003). Then I developed improved computational strategy to identify TALEN target sites in each exon according to the following criteria.

- 1) From the exon sequence database, I identified 40-bp target sequences with 12- or 13-bp spacers and start with the base T and end with A in each exon.
- 2) As described above, I searched for unique TALEN target sites with the minimum number of potential off-target sites using Bowtie preferred without potential methylated sequences.
- 3) Then, I confirmed that the on-target sites indeed exist in human genome sequence by searching sequences using Bowtie.
- 4) To identify additional target sites, I loosened the criteria about off-target that mentioned 2). Therefore, I selected TALEN target sites with minimum potential off-target sites preferred on-target

site alone.

As a result, I designed TALEN target sites for 99.4 % (22,567 / 22,693) of genes and 87.8% (541,985 / 617,233) of exons (included overlapped exons). I identified a total of 3,242,203 target sites (144 sites per gene, 6.0 sites per exon on average (Table 7)

To validate the TALEN mediated mutagenesis activity I selected *NRAS* gene, commonly mutated in myeloid leukemia (Brose et al. 2002), and synthesized TALENs targeted every exons in *NRAS*. *NRAS* gene consists of 7 exons, and I synthesized 6 TALEN pairs (exon 6 is too small to design target sites) (Table 8). These TALENs were transfected to HeLa cells and I performed T7E1 assay. All TALENs could induce targeted mutations with high efficiencies (Figure 30).

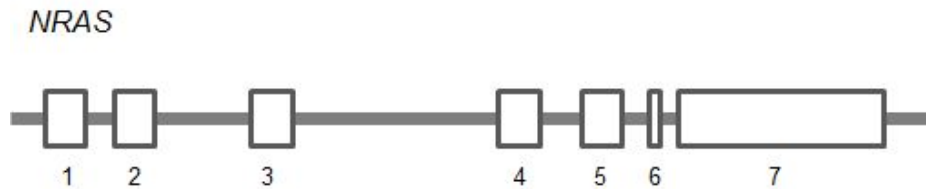
Table 7. Summary of TALEN target sites for exons in human protein-coding genes.

Human								
Criteria		Genes			Exons			Target sites
Off-target	CpG sites	number of genes	%	Avg.	number of exons	%	Avg.	
> 6 (3+3)	0	21360	94.1	101.3	340,940	55.2	6.3	2,162,853
	< 3	21644	95.4	129.7	458,616	74.3	6.1	2,807,219
	< 5	21674	95.5	134.6	483,157	78.3	6.0	2,916,278
	all	21681	95.5	136.2	491,834	79.7	6.0	2,953,457
all		22,567	99.4	143.7	541,985	87.8	6.0	3,242,203

Table 8. List of TALEN target sites for targeting *NR4S*

Exons	Target sequence	spacer (bp)
Exon1	<u>TTTTCCCGGCTGTGGTCCTAAATCTGTCCAAAGCAGAGGCAGTGGAGCTTGA</u>	12
Exon2	<u>TGGTTGGAGCAGGTGGTGTGGGAAAAGCGCACTGACAATCCAGCTAATCCA</u>	12
Exon3	<u>TGGACATACTGGATACAGCTGGACAAGAAGGTACAGTGCCATGAGAGACCAA</u>	13
Exon4	<u>TGATGTACCTATGGTGCTAGTGGGAAACAAGTGTGATTTGCCAACAAGGACA</u>	12
Exon5	<u>TTACACACTGGTAAGAGAAAATACGCCAGTACCGAATGAAAAAACTCAACAGCA</u>	13
Exon7	<u>TTCAGTCTCACAGAGAAGCTCCTGCTACTTCCCAGCTCTCAGTAGTTTAGTA</u>	13

a



b

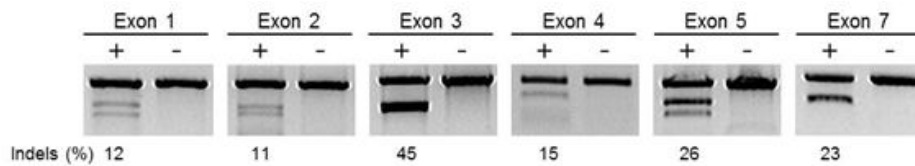
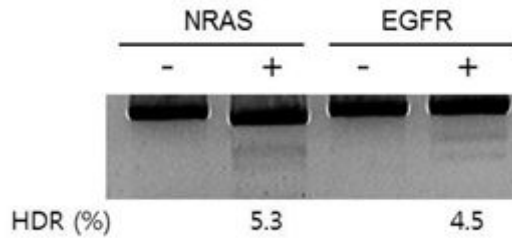


Figure 30. TALEN-mediated mutations in each exons of *NRAS*. (a) schematic of *NRAS* gene. *NRAS* is consist of 7 exons. (b) Targeted mutation in each exons of *NRAS* were detected by the T7E1 assay. all of TALENs were highly active.

As a proof-of-principle experiment for targeted gene modification, I performed TALEN-mediated oncogenic mutations at two genes, *NRAS* and *EGFR*. Q61K mutation at *NRAS* gene and L858R mutation at *EGFR* gene are well known for hyperactive oncogenic mutations in primary cancers (Brose et al. 2002;Paez et al. 2004). With designed TALEN pairs targeting *NRAS* and *EGFR* with appropriate ssODN contained PstI and HindIII respectively, I transfected human K562 cells. PCR amplicons of these chromosomal regions were partially digested by appropriate restriction enzymes, demonstrating ~5% genome-editing efficiency (Figure 31a). I determined the DNA sequences of PCR products representing the targeted genomic regions and found that substitutions were induced in targeted regions (Figure 31b).

a



b

EGFR (11/95)

ctggtgaaacacccgcagcatgtccagatccacagaTTTGGGCTGGCCAACTGCtgggtgcggagAGAAAGAATACCATGCAGAAggaggcacaagtaaggagg genomic DNA

ACACCCGACATGTCAAGATCACAGATTTTGGGCTGGCCAACTGCtgggtgcggagAGAAAGAATACCATGCAGAAggaggcacaagtaaggagg ssODN

ctggtgaaacacccgcagcatgtccagatccacagaTTTGGGCTGGCCAACTGCtgggtgcggagAGAAAGAATACCATGCAGAAggaggcacaagtaaggagg HDR x4

ctggtgaaacacccgcagcatgtccagatccacagaTTTGGGCTGGCCAACTGCtgggtgcggagAGAAAGAATACCATGCAGAAggaggcacaagtaaggagg HDR (-1 del)

ctggtgaaacacccgcagcatgtccagatccacagaTTTGGGCTGGCCAACTGCtgggtgcggagAGAAAGAATACCATGCAGAAggaggcacaagtaaggagg HDR (+1 ins)

NRAS (18/96)

ggttatagatgggtgaaacccgtttgtTGGACATACTGGATACAGCTgggcaagaaaggtACAGTCCCATGAGAGACCAAtacatgaggacagggcgaaggcttcc genomic DNA

ATAGATGGTGAACCTGTTTGTGGACATACTGGATACAGCTgggcaagaaaggtACAGTCCCATGAGAGACCAAtacatgaggacagggcgaaggcttcc ssODN

ggttatagatgggtgaaacccgtttgtTGGACATACTGGATACAGCTgggcaagaaaggtACAGTCCCATGAGAGACCAAtacatgaggacagggcgaaggcttcc HDR x5

ggttatagatgggtgaaacccgtttgtTGGACATACTGGATACAGCTgggcaagaaaggtACAGTCCCATGAGAGACCAAtacatgaggacagggcgaaggcttcc HDR (-3 del)

Figure 31. TALEN-mediated homology dependent repair (HDR). (a) HDR mediated gene-editing in *NRAS* and *EGFR*. TALENs and ssODN including each PstI and HindII restriction enzymes, respectively were transfected to K562 cells and after 3 d of incubation, genomic DNA was isolated, and the target locus was amplified with appropriate primers. TALEN with ssODN-mediated mutations were detected appropriate restriction enzyme treatment. (b) sequence analysis of targeted modifications.

2. TALEN library for several organisms

TALENs are also applicable for several organisms such as mouse, rat, zebrafish and *C. elegans*. In the previous studies, TALENs have been successfully used for making gene-knockout animals and plants. Therefore, as human TALEN library, I designed TALENs targeting protein coding genes for mouse, rat, zebrafish and *C. elegans* (Table 9).

To validate these designed TALENs, I synthesized TALEN pairs targeting several genes in zebrafish (Table 10). Then TALENs transcripts were injected to zebrafish's embryo and performed T7E1 assay to detect mutagenesis activities. The TALEN pairs were highly active from 10 to 50 % of populations as results of TALENs for human genes (Figure 32). I verified the TALEN induced mutation sequences by dideoxy sequencing (Figure 33). These result showed that TALEN libraries for several organisms are useful tools for gene editing. The TALEN library reported here provides a foundation for functional genomic studies.

Table 9. Summary of TALEN library for several organisms

Mouse								
Criteria		Genes			Exons			Target sites
Off-target	CpG sites	number of genes	%	Avg.	number of exons	%	Avg.	
> 6 (3+3)	0	20,552	90.4	60.4	214,446	58.6	5.8	1,241,791
	<3	20,790	91.5	81.9	288,750	78.9	5.9	1,703,261
	<5	20,842	91.7	84.1	299,762	82.0	5.8	1,753,450
	all	20,850	91.8	84.5	303,208	82.9	5.8	1,762,478
all		22,509	99.1	84.8	326,654	89.3	5.8	1,908,519

Zebrafish								
Criteria		Genes			Exons			Target sites
Off-target	CpG sites	number of genes	%	Avg.	number of exons	%	Avg.	
> 6 (3+3)	0	22,986	87.6	39.7	170,174	54.4	5.4	911,721
	<3	23,517	89.6	59.7	239,671	76.6	5.9	1,403,493
	<5	23,547	89.8	61.0	244,915	78.3	5.9	1,436,532
	all	23,551	89.8	61.1	245,305	78.4	5.9	1,438,631
all		26,139	99.6	62.8	274,157	87.7	6.0	1,640,393

Rat								
Criteria		Genes			Exons			Target sites
Off-target	CpG sites	number of genes	%	Avg.	number of exons	%	Avg.	
> 6 (3+3)	0	19,712	85.9	31.1	115,404	54.9	5.3	613,438
	<3	20,444	89.1	43.6	158,799	75.6	5.6	891,864
	<5	20,580	89.7	44.5	163,505	77.8	5.6	915,305
	all	20,599	89.8	44.6	164,436	78.2	5.6	917,917
all		22,397	97.6	45.2	178,838	85.1	5.7	1,011,929

Drosophila								
Criteria		Genes			Exons			Target sites
Off-target	CpG sites	number of genes	%	Avg.	number of exons	%	Avg.	
> 6 (3+3)	0	13,397	96.1	21.0	45,445	58.8	6.2	281,360
	<3	13,674	98.1	34.6	69,176	89.5	6.8	473,747
	<5	13,680	98.2	36.7	72,788	94.2	6.9	501,596
	all	13,680	98.2	36.9	73,138	94.6	6.9	504,289
all		13,931	100.0	36.7	74,054	95.8	6.9	511,880

Table 10. List of TALEN target sites for zebrafish

Genes	Target sites	spacer (bp)
BRAF	<u>TTAAACAGATGATTAAGTTAACTCAAGAGCACCTAGAAGCCCTTTTAGATAA</u>	12
Rad51	<u>TGGGCAGTGATGTTCTGGATAACGTGGCCTACGCCAGAGCCTTCAACACTGA</u>	12
ATM	<u>TCACCAGACACTGTGGAAGAACTTGATCGTACATCAGGAAGCAAAGGCTCCA</u>	12
PTENB	<u>TTTATCCTAACATAATAGCTATGGGTTTCCCTGCTGAAAGACTGGAGGGTGTA</u>	13
CASP7	<u>TTCTCTGTGTAGTACTTACAACGTTTATTTGGAATAGCAGAATACAAGACTGA</u>	13
Rock2a	<u>TGAACCATATAAGGGAGCTCCAGATGAGGCCAGAGGATTTTGACAGAGTGAA</u>	12
Fzd2	<u>TGGTGGGACACTACAACCAAGATGATGCCGGGCTTGAGGTCCATCAGTTTTA</u>	12
Stat3	<u>TCAGTTGCAGCAGTTGGAGACGCGGTATCTGGAGCAGCTGTATCACCTGTACA</u>	13

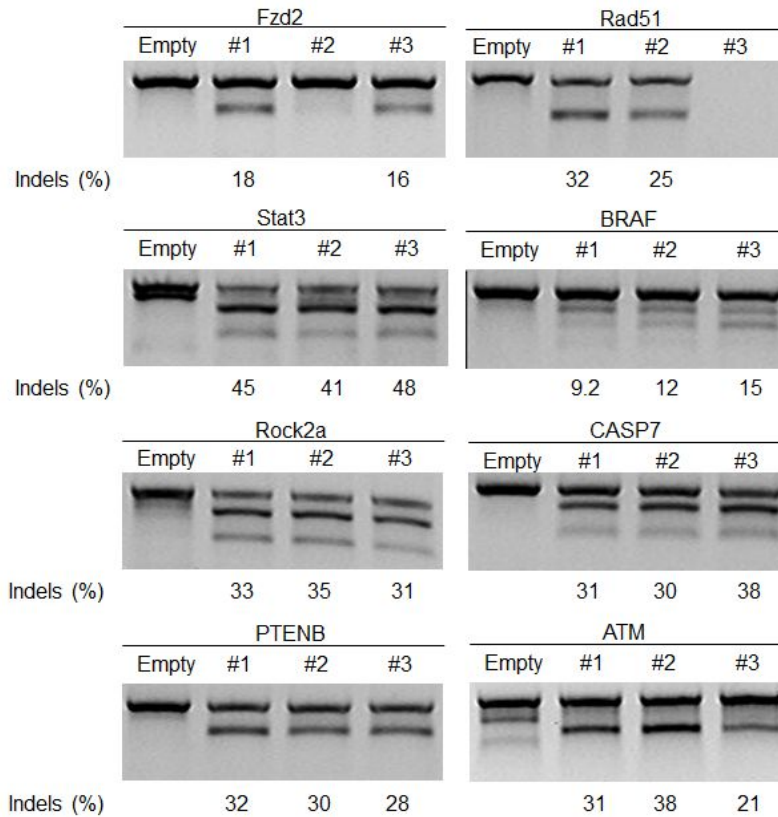


Figure 32. TALEN-mediated gene disruption in zebrafish. TALEN-encoding mRNAs were injected zebrafish embryos and after 1 d incubation, genomic DNA was extracted from pools of 8 embryos. TALEN induced mutations were detected by the T7E1 assay. each numbers of genes indicated independent pools of embryos.

PTENB-9

```
CCTTAACGCTGCTGTACAGACATTTATCCTAACATAAATAGCTATGGGTTTCCCTGCTGAAAGACTGGAGGGTGTATATAGA WT
CCTTAACGCTGCTGTACAGACATTTATCCTAACATAAATAGCTAT-----CCTGCTGAAAGACTGGAGGGTGTATATAGA (-7)
CCTTAACGCTGCTGTACAGACATTTATCCTAACATAAATAGCTAT-----GAAAGACTGGAGGGTGTATATAGA (-15)
CCTTAACGCTGCTG-----AAAGACTGGAGGGTGTATATAGA (-44)
CCTTAAC-----CCTGCTGAAAGACTGGAGGGTGTATATAGA (-44)
CCTTAACGCTGCTGTACAGACATTTATCCTAACATAAATA-----CCTGCTGAAAGACTGGAGGGTGTATATAGA (-12)
CCTTAACGCTGCTGTACAGACATTTATCCTAACATAAATAGCTAat a-----CTGAAAGACTGGAGGGTGTATATAGA (-12,+3)
CCTTAACGCTGCTGTACAGACATTTATCCTAACATAAATAGC-----CCTGCTGAAAGACTGGAGGGTGTATATAGA (-10)
```

CASP7-11

```
AATCAGCCTCAGCTTGTATTCTGCTATTCCAAATAAACGTTGTAAGTACTACACAGAGAAAAAGAACTACTTTCAAACC WT
AATCAGCCTCAGCTTGTATTCTGCTATTCCAA-----CGTTGTAAGTACTACACAGAGAAAAAGAACTACTTTCAAACC (-5)
AATCAGCCTCAGCTTGTATTCTGC-----GTTGTAAGTACTACACAGAGAAAAAGAACTACTTTCAAACC (-14)
AATCAGCCTCAGCTTGTATT-----//-----TAAA (-103)
AATCAGCCTCAGCTTGTATTCTGCTATTC-----GTTGTAAGTACTACACAGAGAAAAAGAACTACTTTCAAACC (-9)
```

Fzd2-2

```
AATGCCAAATCTGGTGGGACACTACAACCAAGATGATGCCGGGCTTGAGGTCCATCAGTTTATCCCGCTTGTCAA WT
AATGCCAAATCTGGTGGGACACTACAACCAAGATGAT-----TGAGGTCCATCAGTTTATCCCGCTTGTCAA (-8)
AATGCCAAATCTGGTGGGACACTACAACCAAGATcaa--CGGGCTTGAGGTCCATCAGTTTATCCCGCTTGTCAA (-5,+3)
```

RAD51-3

```
GTGTGTGCAGGTATGGTCTGGTGGGCAGTGATGTTCTGGATAACGTGGCTACGCCAGAGCCTTCAACACTGACCATC WT
GTGTGTGCAGGTATGGTCTGGTGGGCAGTGATGTTCTGGATA-----CGCCAGAGCCTTCAACACTGACCATC (-10)
```

ROCK2A-6

```
CTCCTCTACCAATCACCCTCACTCTGTCAAAATCCTCTGGCCTCATCTGGAGCTCCCTTATATGGTTCAATGACTTTC WT
CTCCTCTACCAATCACCCTCACTCTGTCAAAATCCTCTGG-----AGCTCCCTTATATGGTTCAATGACTTTC (-10) (X7)
CTCCTCTACCAATCACCCTCACTCTGTCAAAATCCTCTGGCgc--TCTGGAGCTCCCTTATATGGTTCAATGACTTTC (-4,+2)
```

Figure 33. Sequencing validation of TALEN-induced indels in zebrafish.

Endogenous mutations induced by 10 different TALENs were confirmed by dideoxy DNA sequencing. The numbers of inserted or deleted bases are shown on the right side of each mutant sequence. WT, wild-type sequence.

E. Comparison of TALENs and ZFNs

TALENs and ZFNs share the same FokI nuclease domain and induce site-specific DNA cleavages, whose repair via error-prone NHEJ gives rise to indels at the target sites. I investigated difference of TALENs with ZFNs in patterns of mutations (Kim et al. 2013c).

First, I compared mutation signatures of ZFNs and TALENs that have been reported in the literature. I calculated frequencies of insertions, deletions, and complex patterns that accompany both insertions and deletions at each target site in the following categories: mammalian cells, mammalian organisms, non-mammalian animals, and plants. my analysis included a total of 1,456 mutant sequences at 122 target sites reported in 43 independent studies (Figure 34). ZFN-induced mutations are much more evenly distributed between deletions and insertions in any organism or cell line than are TALEN-induced mutations. In particular, insertions are rarely obtained with TALENs. For example, ZFNs induce insertions in mammalian cells at a frequency of 43%, comparable to that of deletions (52%). In sharp contrast, TALENs induce insertions at a frequency of 1.6%, much lower than that of deletions (89%). In all other systems combined, the ZFN-induced insertion frequency is much higher than the TALEN-induced insertion frequency (24% vs. 4.1%) (Student's *t*-test, $P < 0.05$). I also found that ZFNs are associated with lower deletion frequencies (59%) as compared to TALENs (81%) ($P < 0.05$). Complex patterns are obtained at comparable frequencies (17% with ZFNs and

15% with TALENs).

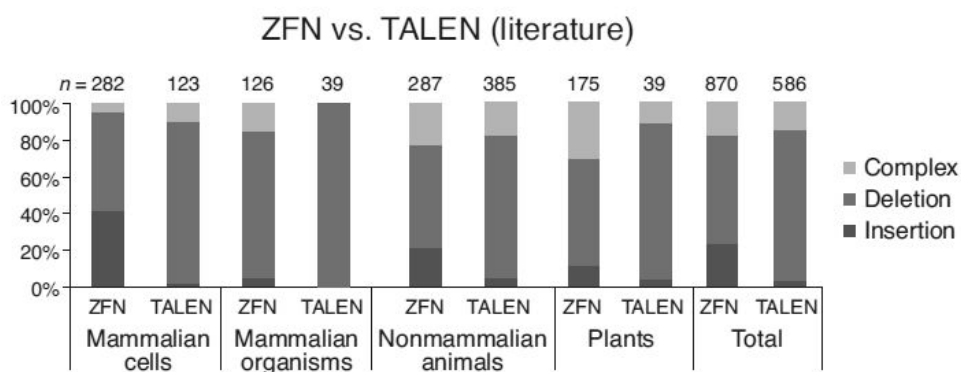


Figure 34. ZFN and TALEN mutation patterns reported in the literature. Nuclease-induced mutant sequences were classified as three different types: insertions, deletions or complex (that is, having both insertions and deletions) patterns. The number of mutant events is shown above each other.

Next, I investigated that pattern of TALENs induced mutations different with that of ZFNs by several TALENs and ZFNs whose target sites overlap with each other (Figure 35). I transfected HEK 293 cells with them, and compared mutation patterns using Sanger sequencing methods. In line with my analysis described above, TALENs, unlike ZFNs, rarely produced insertions. Thus, five TALENs induced insertions at an overall frequency of 2.6%. In sharp contrast, ZFNs induced insertions much more frequently (39%) ($P < 0.05$) (Figure 36). This results showed that the differential mutation signatures induced by TALENs and ZFNs do not arise from the differences in target loci.

F8
ZFN : AGGCAAGAATTAAGTCGGCCCCACCTTTGCccaactCAGTGGGTCTCCTTGAGAGAGGTCTGCAG
TALEN : AGGCAAGAATTAAGTCGGCCCCACCTttgcccactcaGTGGGTCTCCTTGAGAGAGGTCTGCAG

F9
ZFN : GTTACAACATATGGTTGCCAGGTACTGTGTcagggTACTAGGGGTATGGGGATAAACCCAGACTCCCT
TALEN : GTTACAACTACGGTGCCAGGTACTGtgtcagggTactAGGGGTATGGGGATAAACCAGACTCCCT

CCR5-1
ZFN : TTTGTGGGCAACATGCTGGTCATCCTCATCetgatAAACTGCAAAAGGCTGAAAGAGCATGACTGA
TALEN : TTTGTGGGCAACATTGCTGGTCATCCTcatcctgataaacTGCAAAAGGCTGAAGAGCATGACTGA

CCR5-2
ZFN : ACGCACTGCTGCATCAACCCCATCATCTATgcctttGTCGGGGAGAAGTTCAGAACTACCTCTT
TALEN : ACGCACTGCTTGCATCAACCCATCATCtatgcctttgtcGGGGAGAAGTTCAGAACTACCTCTT

CCR5-3
ZFN : CAAGATGGATTATCAAGTGTCAAGTCCAATCTATgacatcAATTATTATACATCGGAGCCCTGCC
TALEN : CAAGATGGATTATCAAGTGTCAAGTCCAatctatgacatcAATTATTATACATCGGAGCCCTGCC

Figure 35. Overlapping target sites of ZFNs and TALENs. Both the left and right half-sites are underlined. Spacers are shown in small letters.

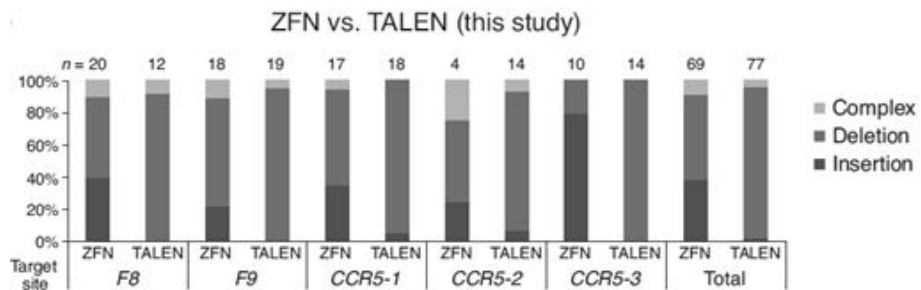


Figure 36. Comparison of ZFNs and TALENs that target overlapping sites. Nuclease-induce mutant sequences were classified as three different type: insertions, deletions or complex (that is, having both insertions and deletions) patterns. The number of mutant events is shown above each other.

Furthermore, I investigated the molecular basis of the difference in mutation patterns. I hypothesized that the larger spacers (12 to 21 bp) in the TALEN sites as compared to the ZFN sites (mostly 5 to 6 bp) might cause the difference. To test this idea, I expressed two homodimeric ZFNs in cells whose genome was modified to contain ZFN sites with 15-bp spacers (Figure 37a). Unlike ZFNs that target sites with 5- to 6-bp spacers, these ZFNs produced insertions much less frequently than deletions (1.8% vs. 86%) (Figure 37b). I speculated that the larger spacers either in ZFN sites or in TALEN sites may give rise to heterogeneous overhangs when cleaved by nucleases. In contrast, ZFNs that target sites with 5- to 6-bp spacers produce defined 4- or 5-bp overhangs, which might be efficiently filled in before ligation to yield insertions. Defined overhangs may facilitate targeted insertions of plasmid DNA at genomic sites (Maresca et al. 2013).

a

ZFN 1

Deletion

```
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA WT
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaag-----tatAAACTGCAAAAGGCTGAAGAGCATGACTGACA (3)
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTac---ttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA (2)
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttgaa-----ACTGCAAAAGGCTGAAGAGCATGACTGACA (4)
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTT-----GGTTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaa-----ACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttga-gtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGC-----aagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttgaa---AACTGCAAAAGGCTGAAGAGCATGACTGACA (2)
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTaca--cttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTac---cttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTT----cttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttg-----CAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttg---atAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagc-----AAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTT-----gaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttgaa-----CTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTT-----ttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
```

Complex

```
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA WT
ATCCTCATCCGGAATTCTAGAAACTTTTGCAG-----acttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTTacaagcttga---AACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCG---atAAACTGCAAAAGGCTGAAGAGCATGACTGACA
ATCCTCATCCGGAATTCTAGAAACTTTTGCAGTTT---cttgaagtatAAACTGCAAAAGGCTGAAGAGCATGACTGACA
```

ZFN 2

Insertion

```
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaagcttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC WT
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaagcaagcttgaagttgGATGAGGATGACTTCTAGATCTACGT
```

Deletion

```
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaagcttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC WT
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCcca--cttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaagcttga-----TGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaagc---agttgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCc---cttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC (3)
TTTGGTTTTGTGGGCAACATGCTG-----gGATGAGGATGACTTCTAGATCTACGTAAAC
AAAACCTA-----//-----gGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGG-----ATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCcc--cttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC (2)
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCcc-----aagttgGATGAGGATGACTTCTAGATCTACGTAAAC (3)
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCA-----tgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCAT-----cttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCc-----ttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCAT-----GACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCcca-----CTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATG-----aagttgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGG-----tgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCc-----TAAAC
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaa-----tgGATGAGGATGACTTCTAGATCTACGTAAAC
TTTGGTTTTGTGGGCAACATGCT-----tgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC
```

Complex

```
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCccaagcttgaagttgGATGAGGATGACTTCTAGATCTACGTAAAC WT
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCctga-----TAAAC (2)
TTTGGTTTTGTGGGCAACATGCTGGTCATCCTCATCcc----aaaagttgGATGAGGATGACTTCTAGATCTACGTAAAC
```

b

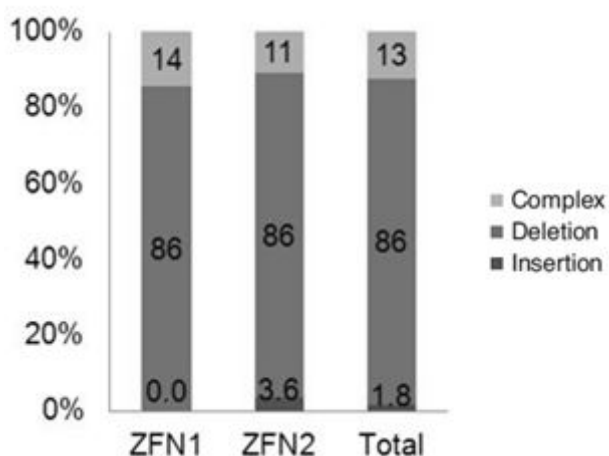


Figure 37. Mutations induced by ZFNs with 15-bp spacer. (a) I modified the *CCR5* target site of ZFN-224 by cotransfecting K562 cells with ZFN-224 plasmids and oligonucleotides as described previous study (Chen et al. 2010). The two monomers that constitute ZFN-224 were used as homodimers to cleave the modified sites with 15-bp spacers. Sequences recognized by ZFNs are underlined. Spacers are shown in small letters. (b) Mutation signatures of ZFNs that target sites with 15-bp spacers.

IV. Discussion

The TALEN technology has recently developed, which can be rationally designed *de novo* due to its highly predictable modularity. Unlike ZFNs, which have sequence-bias toward guanine-rich sites (Maeder et al. 2008), TALENs have no sequence-bias toward any base and, thus, can be designed to target any sequence.

Over the last two years, TALENs have been reported for genome editing in yeast (Li et al. 2011), roundworm (Wood et al. 2011), fruitfly (Liu et al. 2012), cricket (Watanabe et al. 2012), zebrafish (Huang et al. 2011;Sander et al. 2011;Bedell et al. 2012), frog (Lei et al. 2012), mouse (Sung et al. 2013), rat (Tesson et al. 2011), pig (Carlson et al. 2012), cow (Carlson et al. 2012), thale cress (Cermak et al. 2011), rice (Li et al. 2012), silkworm (Ma et al. 2012) and human cells (Hockemeyer et al. 2011;Miller et al. 2011;Reyon et al. 2012;Kim et al. 2013b). Although previous TALEN architectures are highly active with high success rates, for genome-scale study, they are limited some aspects. 1) The success rates of TALEN-mediated gene disruption range from 64% to 88% in mammalian cells (Cade et al. 2012;Carlson et al. 2012;Reyon et al. 2012). These success rates are certainly impressive compared to those achieved with ZFNs. On a genomic scale about 20,000 genes, however, failure rates of 12–36% correspond to 2,000–7,000 genes, if one TALEN is used per gene. 2) The previous TALEN architectures have poorly defined fusion junctions between TAL effectors and FokI cleavage domain (Christian et al.

2010;Li et al. 2011;Miller et al. 2011). Different with ZFNs that have a narrow range of spacers (5- to 6-bp) (Kim et al. 2009), these TALEN constructs induce mutagenesis in broad spacers (12- to 22-bp) that may induce indels at unwanted loci in genome.

To this end, I developed TALEN architectures with 12- to 14-bp spacers that might reduce unwanted mutations in genome. Furthermore, using obligatory heterodimeric FokI nuclease domains, the TALEN constructs were improved target specificities through elimination of TALEN monomer's homodimeric activity. And my thesis suggested that a nearly 100% success rate could be achieved with TALENs by avoiding heavily methylated sites, making it feasible to initiate genome-scale gene-knockout studies. This high success rates were important for TALEN mediated genome scale study.

I developed a scalable Golden-Gate assembly system that consists of a total of 432 plasmids (424 TAL effector array plasmids and 8 FokI plasmids) and used this system to construct a TALEN library, a collection of 18,740 TALEN pairs, to disrupt or modify every protein-coding gene in the human genome. Recently, two different solid surface-based assembly methods for the synthesis of TALEN plasmids were reported: iterative capped assembly and fast ligation-based automatable solid-phase high-throughput (FLASH) systems (Briggs et al. 2012;Reyon et al. 2012). Although these methods can be automated at least partially, multiple cycles of ligation and washing steps in addition to PCR are involved before subcloning into an expression vector. Furthermore, FLASH uses up to microgram quantities of TAL effector

array plasmids in each reaction, making it cost-inefficient and cumbersome to scale up. Compared to these methods, My Golden-Gate cloning system allowed PCR-free TALEN assembly with much less hands-on time and greatly reduced plasmid quantities. In addition, cloning efficiency was higher in general with Golden-Gate systems than with solid surface-based methods.

The one-step Golden-Gate cloning systems can use several other applications of TAL effectors. The DNA sequence based binding protein is valuable resource for targeted gene regulation by fusing VP64 activation domain or KRAB repression domain (Maeder et al. 2013b; Perez-Pinera et al. 2013). TAL effectors fused histone demethyltransferase are also useful options for targeted gene regulations (Maeder et al. 2013a). These gene regulators enable the gain of function studies. In addition, because they are not induced permanent gene modifications like engineered nucleases, they can be used for transient or inducible gene regulations. TAL effector's random or target library pools also can be used for phenotype screening as zinc finger protein library.

TAL effectors also can be used for targeted chromosomal visualizers to combine with fluorescent proteins (Miyanari et al. 2013). Unlike other chromosomal visualizers such as FISH that need fixed cells by chemicals and only processed in death cells, TALE fluorescent visualizers enable to detect target sequences in living cells. So, they can be used to detect the shape and loci of targeted chromosomes during the cell cycles.

Then, I established computational strategy for TALENs targeting almost protein coding genes in several organisms. As described results above, I identified 40-bp target sequences with 12- or 13-bp spacers in each gene, which were recognized by TALEN pairs that consist of two 18.5 RVDs with the minimum number of potential off-target sites. The ten of TALENs that were highly active at their on-target sites, did not induce any measurable mutations at their most likely off-target sites.

In the short term, TALEN-mediated gene-knockout experiments should be useful for validating or invalidating the functions of genes identified by genome-scale siRNA library screens, which have revealed many potential drug targets associated with various diseases (Sigoillot and King 2011). Unfortunately, many of the newly identified genes have turned out to be false positives that resulted from off-target effects (Lin et al. 2005; Adamson et al. 2012). To distinguish true positives from false positives, one could use TALENs to knock out candidate genes as demonstrated here, rather than using another siRNA to downregulate the genes. In the long term, TALENs could enable the creation of a series of cell lines in which a family of druggable target genes encoding proteases, kinases, phosphodiesterases and G protein-coupled receptors are disrupted. In principle, genetic screens can be performed using these cell lines instead of siRNA libraries. The TALEN library reported here provides a foundation for functional genomic studies based on gene knockout rather than knockdown in mammalian cells.

Furthermore, the TALEN pairs in advanced libraries for every exon in protein coding genes of several organisms can be used to generate not only gene disruption, but also precise genetic modifications by homology directed repair. For example, as described above, single nucleotide changes can be made in a gene of interest to study phenotypic differences associated with single nucleotide polymorphisms in an otherwise isogenic background. A reporter gene can be fused to a gene of interest to monitor the expression of the gene or to track the protein.

Although the primary use of my TALENs in drug discovery will be to identify and validate druggable target genes, I envision that some TALENs or their improved versions would be useful therapeutically themselves. An example for a targeted nuclease that is in clinical trials is Sangamo's CCR5-specific ZFN for the treatment of AIDS, which is now under clinical investigation in the United States (Holt et al, 2010). TALENs can also be used to disrupt disease-associated genes or correct genetic defects in stem and somatic cells. I propose that the one-step Golden-Gate cloning system and TALEN library presented in this thesis will be widely used for basic and biomedical research.

V Appendix

1. List of TALEN activity in reporter assay

Gene ID	sequence	Avg.	S.E.M
1544	TTGTCCAGTCTGTTCCTTCTCGGCCACAGAGCTTCTCCTGGCCTCTGCCA	28.6	1.3
5354	TCCATGCTTCCAGTATGTCACTATGGAAGTGCCTCTTCTCTTCCITTA	3.4	0.1
7007	TTAAGTTGGCCATCCAGTTTCTTCTTTGGCGTTCCCTTACCGCACTGTCTA	14.2	1.0
2875	TGGTGGCCGGCGAGGGCCACACAGCAGGGTGTGCTCATCCCATCCCCCA	32.9	2.0
11023	TCACCGGGAGCAGCTCTATCGGCTGGAGATGGAGTTCAGCGCTGCCAGTA	3.3	0.3
866	TCTCTGGAGTCTATGACCTCGGAGATGTGCTGGAGAAATGGGCAITGCAGA	20.7	3.3
146059	TCAACCCCAATCAGCTGTGTCCCAAGTTCACCAACCTCAGCCCTGGACACTA	5.3	0.9
4893	TACAAAACAAGCCACGAAGTGGCCAGAGTTACGGGATTCCATTCTTGAA	4.0	0.7
119	TGGAGAAGGGCAGCAGCTGCTTCCAGTGGACACACAGGCTTCTGTGCA	11.8	0.7
1401	TTCACTGTGTGCTCCACTTCTACAGGAAGTGTCTCGACCCGTGGGTACA	12.5	1.2
257	TCCAAAGTGCAGCTTCCCCAGCTGCCCTTGGACTGCGAGGGGGCCCCA	25.2	0.8
1675	TTCTCTGGGCGCGCACTCCCTGTGCGAGCCGAGCCCTCCAAAGCGCTGTA	6.4	0.5
8908	TCTTCAGAACTGGTGCAGCAGACATCCACAAGCACTGCCGTTTCACTA	22.1	1.7
1588	TCATCTCCACGGCAGATTCTGTGGATGGGGATCGGCACTGCTGCAACTA	22.0	1.5
958	TCTGGAGCCCTGCCAGTCGGCTTCTTCTCCAATGTGTCTCTGCTTTCGA	16.2	2.2
1759	TGGCATTGTCAACCGAGCTCCCTGTGCTTGCAGCTGGTCAATGCAACCACA	21.3	2.5
9228	TGCCCTCGACCTCCCGGGATGTTTCGAACAAGGAGTCAAGCTACCTTCGA	9.3	1.2
491459	TCTGGAGAACATGGAATTGTGGTTTCTCTTTGGGATCAATGGTCTCAGAA	38.4	1.5
5309	TGCTGTGCAGACCTGGCACTCCGACCCCAAGCTCCAGAGCACGGCTGCA	38.9	1.0
11277	TGCAGACCCCTCATCTTTTTCGACATGGAGGCCACTGGCTTGCCTTCTCCCA	14.1	1.3
2902	TTAGCCATCCACCTAGCCCAACGACCACTTCACTCCACCCCTGTCTCCTA	14.2	2.4
2321	TACTTAGAGGCCATCTCTTGTCTCAATTGTACTGTACCACTCCCTTGAA	17.0	1.0
2693	TGGCCGACCTGGACTGGGATGCTTCCCCCGGCAACGACTCGCTGGGCGACGA	17.1	2.2
54716	TCCCTACATCATCATGCTTATCGTGGAGGGAATGCCGCTCTGTACCTGGA	4.9	0.5
722	TTTCTAAGCATAAGGCTTTTCTGTACCTACAGCTGTGACCCCGCTTCTCA	12.5	1.2
5798	TACAGGACATCCCACTGGCTCCGCCCTGTGCTGCCAGCATCGGCTTCCACA	40.6	2.3
8745	TATGCTCTTCCATGGCAGTGGCACAAGTATTATCGCAGAGCCTGGCTCAA	24.0	2.7
2671	TACTACCCCGACCTGCCACCCAGAACAGCAGCAAGACATGGCCAGTTCA	17.9	2.3
4747	TGGCTACTCCCGAGACTCCAGGTCTTTGGCGATCTGCTACGGCGGTTTA	6.6	0.8
59344	TCTTCCCTCTCTCTGGATCACAGGACACGGGAGCTCGGGCCCGACAAGAA	34.3	4.2
7433	TTTATTTGCATCATCCGAATCCTGCTTCAGAAACTGCGGCCCCAGATATCA	12.3	1.0
4920	TTGCCTGTGCACGCTTCATTGGCAACCGGACCATTTATGTGGACTCGCTTCA	13.4	1.1
3816	TGGAGTGTGAGCAGCATCCAGCCCTGGCAGGCGGCTCTGTACCATTTCA	26.3	2.2
3697	TCCCAGAACTCAGCAACCATGCCTCAATACTCATCATGTGACAGATGGCGA	17.5	1.5
7441	TGTCTCGGCCCTTGGAAACCAATCCGCTCACCTGCACCCCTGAGGAACGA	2.0	0.4
50943	TGGGCCCTTCTTGGCCCTTGGCCATCCCGAGAGCCTCGCCAGCTGGA	14.6	1.2
5345	TGCTCAGAACCAACGTTGCAGAGGCTGCAACAGGTGCTGCAACGAGGCTCA	23.9	2.5
20	TGCAATCGCTGTGCCCGAGCGGCAGCGAGACGAGTTCGCTTCTGCACTA	6.7	0.9
81031	TCTGCTTTTGTGTGCTCTGTGTCTTGTGGTGGCTGACCTTTGGTTA	2.4	0.3
4602	TTCCAAAGTCTGGAAGCGTCACTTGGGAAAACAAGGTGGAACCGGGAAGA	23.4	2.7
4867	TCTTCCAGGAATGAGCATACAGGTCTCAGCAGACATGTACGCTCTGTCTA	8.3	0.5
6622	TGGAGGGAGCAGGGAGCAITGCAGCAGCACTGGCTTTGTCAAAAAGGACCA	20.9	0.7
268	TGGACCGCCCTGCGGGGCTGGCGCGGCTCGGGCTGSCCTTGAACCTGCA	26.3	3.0
2100	TTAGTGGGAACOGTTGCGCCAGCCCTGTACTGGTCCAGGTTCAAAGAGGGA	26.6	1.6
6331	TCTCAGTCCCTTCCACCCCATCGGAGAGCGGCTGTGAAGATTCTGGTTCA	20.9	2.7
6338	TACCCAGGCATGACAGAGTGTACATCCTGAGGCCACCAACATCTTTGCA	36.2	2.7
4669	TGTGTTGGGACTGGGCGCGCTGGGAGCGAGAGATAGACTGGATGGCGCTGAA	40.3	1.9
7976	TGTTCTCGGATTTCCGGCCCTTTCTTTGTGCACTCTACGCTCCTATTGTA	24.4	2.0
6531	TCTTGGTCCCTACCTGCTCTTCATGGTCAITGCTGGGATGCCACTTTCTA	16.4	3.3

Gene ID	sequence	Avg.	S.E.M
1311	TCCACTGCGGCGCCGCTTCTGCTTCCCGGGGCTGGCCTGCATCCAGACGGA	18.6	1.3
26254	TCTCCATCGATAATGATGCCTTCCGCTGCTACATGCCCTCCAGGACCTCA	14.2	1.8
1183	TTGGGAATAGCCGCTCTCGTTCTCTTTTATGTGGAATACCACACGCCCTGGTA	22.0	1.4
2239	TGGTCAGCGAAGCAGTGCATCATTTGCAAGCTGTCTTTGCTTACGTTACAA	22.0	2.5
197	TGCACCTCCGTCCTCCCTCCACTTGGCGCACTGGACTCCCTCCAGCTGGCTCA	26.5	2.5
5837	TGCTTGGCTATCGCAACAAATGTTGTCAACACCATGCGCCTCTGGTCTGCCAA	13.1	1.9
8224	TGGAAAACAGCTTGACTTCCAGGACATCACCAGCGTGGTCGCCATGGCCAA	25.7	3.2
5365	TCCGACCTTAATACCACTCTCAACCACTTGGCACTGGCACTGGCCGAGGCA	22.4	2.1
785	TCTGACTCCGATGTCTCTTTGGAAGAGGACCGGGAAGCAATTGACAGGAGA	31.5	0.7
57211	TTTGGAGGATGGACACGTCAGGATGTGTGACACACAGATTCAGATGCAA	1.6	0.2
2036	TCCCAAACACCGGTCAGCCAGAGACTGTGGAAGGTCTGCATCGAGCATCA	25.5	1.8
6017	TGGGGGAGAGGCTGGCGGTGGCGGTGGCGGAGAGGGTGCAAGAGAGGACA	25.3	4.6
4486	TGGGTGACATGGGCGACCGCTGTCAGCGGGATGTGAGTCTCTTGGGGACCA	20.1	2.5
6556	TGCTGCAGGCGGTGGGCATTGTTGGCGCCATCATGCCCCACAACATCTA	11.6	1.6
1369	TGTGCGAGGAGTTTCGGGAACAGGAACAGCGCATGTCAGCTCATCCAGGA	22.3	1.3
3786	TCTGGGCTGCTGGATGTTGCTGCGGATACAAAGGCTGGCGGGGCCGACTGAA	39.0	1.2
26047	TACAGCATCCGATTAATGCGCGCTATGTGCGCATAGTGCCTCTGGATTGSA	19.0	3.9
7137	TCAGACGCGCTCTCTCCAACTACCGGCTTATGCCAGGAGCGGACGCCAA	27.2	2.5
2645	TGGGGCTGCGACCTCGACCCACCGACTGCGACATGTCGCCGCGGCTGCGA	22.3	2.5
5459	TGGCGGCGGTGGATATCGTCTCCCAAGGCAAGAACATCGTTCAAGCCCGA	15.5	3.7
6091	TGCAGTAATGGAATGCCAACCTCCAGAGGCCATCTGAGCCCATTTTCA	14.3	2.5
11107	TGCTACCAACCCACGCGCTCTCAACTGGCTTGCCTTCGTTTCATGAGGACCA	33.3	3.1
3046	TTACCCCTGAAGTGCAGGCTGCTGGCAGAAAGCTGGTGTCTGCTGTGCGCA	11.1	1.1
63976	TCCCGTGGGGCCACCGTCCCTCTCCCAAGGAGGAGCTTCAACCCCAA	20.0	1.5
90	TTGCGAGTATGCTTTTAGCCTGCTGCTGGGAGTTGCTCTCGAAAAATTTA	5.2	0.4
81693	TGGAGGAGGACAGAGGCGGCGGCCGCTGGAGCGCCCTCGGCTTCCGCA	29.3	1.6
7054	TGCGCCGAGGGGACCTGGCGGCCCTGCTCAGTGGTGTGCGCCAGGTGTGAGA	16.6	1.0
5346	TGGCTTCCGAGGCTTGGACCACTGGAGGAAAAGATCCCGCCCTCCAGTA	2.3	0.4
10743	TGGTTTCCATGGCAACACAGAACTACGACGAGCTTCGAGGAGAACATCA	28.5	3.4
489	TGCGCCGAGTTCGACGGGCTGGTGGAGCTGGCGAACATCTGCGCCCTGTGCA	27.1	2.5
650	TTGGAAGAACTACAGAAACGAGTGGGAAACAAACCGGAGATTCTTCTTTA	8.0	1.3
3640	TGGTGGCCGACAGTAATCTCACGCTGGGACCTGGCCTGCAGCCCTGCCCCA	32.7	2.4
7224	TGTACTACAAAAAGTCAACTACTCACCGTACAGAGACCGCATCCCTCTGCA	3.0	0.1
50939	TACCAAAAAGAAACAGCCTCTGGACGCGAGAGAACTGAAAGACAGTGGTTA	1.9	0.3
5071	TGGGGGCGGTGTTATGCCCGCCCTGGCTGTGGAGCGGGGCTGCTGCCGGA	15.2	0.8
2588	TTTTCCTCTACTGGGCTGTGACGCCACGCAACCGCTCTATGCTCCAA	25.7	1.3
51206	TTTTGCCAAACCTCGCTCTCAGCCAGCCCGGCGCGGCTGTCTGTCAGGA	24.1	0.8
2516	TGCAGAACCAAGCACTACAGTGCACCGAGAGCCAGAGCTGCAAGATCGA	12.9	1.0
1789	TCGATCCTCGTCAACGGGGCTGCGAGGACAGTCTCCGACTCGCCCCAA	23.2	1.7
50508	TTTCTTCAACCTGGAACGCTACCACTGGAGCCAGTCCGAGGAGGCCAGGGA	0.8	0.1
63970	TCTTCTCTGAGGCGAGCTTCAGATCTGCTCAAGCTTCTGCGTGGGGCCA	37.3	1.1
9333	TCCAGAGCTCCAGAAATAATGTGGGCAACACAGGAGAGATCACTGTGGA	20.5	1.7
1583	TGGTGCAAGTGGCCATCTATGCTCTGGGCGGAGAGCCACCTTCTTCTTGA	22.4	1.0
2253	TCATCGTGGAGACGACACCTTTGSAAGCAGAGTTGAGTCCGAGGAGCCGA	27.1	1.5
6492	TGTCCCTCTCGACGCTGCTACCAAAAGTGGGCTGGTGGCGGTGGGCCA	36.0	3.4
6491	TCAAGCAGGATGTTACCTTTCACCTTCTCCATCAAAATGTGCACTTTGGA	19.0	0.6
2914	TCCGCCGCTCTGGAAGCTTCGAAAGCCAGGGGAGTCACTATCTTTGCCAA	13.6	0.4
11173	TCGAGTGCAGCGGGGGGCTCTCTGTCTTACGAGCTGTGGCCCGGCGCA	12.0	0.5
51196	TACAGTTGAGTCTCTACTCCATTCTCTCACCAGCTCCAGCCTCCGAGACA	23.1	1.0

Gene ID	sequence	Avg.	S.E.M
10195	TCTTCACTCCAACTTCATGGCATCTGCTTCAGCCGCTCCCTCCACTACCA	18.1	1.3
6010	TCAAGCGGAGGTCAACAACGAGTCTTTTGTCTCATCTACATGTTCTGTTGCTCA	24.8	1.7
29760	TCAAGTACAGACAGAAGCTGTGCTAGTACCAGTGTTCCTCTCTGCCAGCA	27.1	1.5
3363	TGGCTTCTCTCCGCTCCATCACCTTACCTCCACTCTTTGATGGGCTCAGA	37.1	1.5
56606	TCTCCCGACAGCCACGCTACCTGCTCTTGAGAGACACACAGGCAAGA	13.6	1.1
6530	TGCGCTGTACGTTGGCTTCTACTACAACTCATCATGCTGCTCACTCTA	7.1	0.4
6517	TTCTCTGCGGTGCTTGGCTCCCTGCAGTTTGGGTACAACATGGGGTCATCA	28.2	2.2
2532	TCTCCCTCTCACTGAGAACTCAAGTCAGTGGACTTGAAGATGTATGAA	32.1	1.3
8022	TCTCAAGGCTCTGGACGCTCAGTGGACAGCAAGTGTCTCAAGTGCAGCA	24.3	1.6
415	TCTCTCTCTTTGTTTCTCTTCTACAGTTTCACTCCCTCTTATCACTATGA	20.2	1.2
3730	TTCAAAGACAGCAACCACTCGCCCGCTGGAAGTCGGAGCTCCCTTCTATCA	13.8	1.2
1748	TCCGAGCGGGCCCTCAGGCCCCCGCCAAAAGCTCCGCAAGCGGAGGACA	29.6	1.0
3918	TGTCTGCCAGGCTTCCACATGCTCAAGGATCGGGGTGCAACCAAGACCA	22.6	0.8
367	TGCGCAATTGACTATTACTTTTCAACCCAGAGACCTGCTGATCTGTGAGA	30.2	1.1
81030	TCTACAGGTTTCTCAAAGACAAATGCTCCCGAGGGCCCTGGTCATCGCCA	30.0	1.7
4330	TGAGGCTCTCGTCCCACTGGCGGTCCCTGCTTCTGACATCTCCAACAGA	27.9	2.7
4693	TGGTCTCTCTGCGCAGGTGCGAGGGGCACTGCAAGTCCAGGCTCAGCTCCGA	34.3	1.3
84667	TCCGGTTTCCGCGAGTGGCTGCTTCGCTTGGCGGCTTCGCGCACGACCA	24.3	2.3
6528	TGTTCTACACTGACTGCGACCTCTCTCTCTGGGGCGCATCTCTGCCCAAGA	16.0	0.6
5317	TGGCTGACAATTACAACTATGGGACCAACAGCAGGAGCTACTACTCCAA	17.3	1.6
8038	TCCAGGGAGCCACCATCGGCATGGCCCAATCATGAGCATGTGCAAGGAGA	6.6	0.5
6524	TGCGCGCGCCCGCATCGCCCTTACCTGTCTGTCTCTCTCTTCTCTGTA	13.3	2.5
1592	TCTACAAGACGATCTGTTCGGCGCGCCCACTACGGGTGATGGCGCGGA	10.3	0.4
2122	TTCCCTCTCAGCATGTGCGTCCCGGGCCCATGCTGCGCGGAGTGTGGCAA	31.4	2.8
6497	TCCAGTCCCGCGCCCTTCGAAAAGGACAAGCGTCCAGCTGGCTGCGGA	31.0	2.1
3977	TCCAGGAAGGATGACAGCTTGGTGGGCCCACTGCTACAAGCTACACTTA	18.8	0.9
1180	TCTACGCGCAGATGAGGCCAGCCCTCTCTGCTGCTCTCTGCTGCGGTCA	28.9	0.5
8862	TCCGCCACCTGCTGCGAGCCAGAGGTCAGGAATGGGCCAGGGCCCTGGCA	33.9	2.7
215	TACAGCGCTCTTACAGGACCTGCGCTCGCAGATCAACTCATCTCTCTGGA	18.4	2.0
132884	TGCTCTTGCTGGAAGCATGCTGGTCTCACCATTITGGGACTCTGTGGGAA	21.7	2.5
2153	TGGCCAGCCGCCCTATAGCAITTTACCTCATGGAGTGAACCTTCTGCGCTTA	12.2	0.8
4916	TGGAGCACTGCATCGAGTTTGTGGTGGTGGCAACCCCAACCAAGCTGCA	31.6	1.7
2903	TCTGAGTGGTGGACTTCTCTGTGCCCTTTGTGGAAAGCGGAATCAGTGTCA	36.2	1.4
360	TGGCCAGGTTGTCTCAGCGGGGCAACCAAGTGGTTTCTCAACATCAA	20.0	2.1
9907	TCTTCAAGGTGCTCTCCAGCGGCCACCGCTTCAAGACGACCACTGGCTGAA	31.9	1.0
50814	TGGCTCTGGATTCTGGGGCAGCAGATGCTGGAGCAGTTGCTGGCAAGAGGA	7.1	0.6
55315	TGTCTCAGAGGACGACTTTCAGCAGCTTCAAACTCCACCTACAGAACCA	30.1	1.8
3795	TCTCCCTGCTCGAGCCCTCTGTGCTTCTGCGCTCAATGGCTCCTGGCCA	34.3	2.5
6716	TCAITCTCAGAGCACTGCTTCTGCACTGGAAATGGAGTCTTCAAGGCTA	2.9	0.4
43	TTTCTGGTTTACGGGGCCCGAGGCTTCAGCAAGACAAACGAGTCTCTCATCA	17.6	1.1
9668	TGGTGGGAGCAGGTGCACAGGAGCCAACTCATCCGCTGCTCCGGAAGA	13.3	1.6
5972	TCAAAGTCGCTTTGCACTGGTTCGTCCAATGTTTGGGTGCCCTCTCCAA	7.4	0.2
51168	TGGAATCTATGGGCGGAGCAGGTGACGAGTACAAACGAGCGGCGCTGGGA	10.5	1.5
9248	TGGGGCTGGGATCTGAGCCAGTTACAGTGCAGATGGTGGGTTCTACACA	16.9	2.1
3049	TCCCACTGAGCTGAGCCCGGCTCTCTCAAGTCAGAGCCACGGCCAGA	23.6	2.4
3954	TGGCTGCTGCACTCCCATCCACCTGTGTACACATCTCCAGAGGCGATCA	31.0	2.7
7959	TGCCCTGCAAGTGGCCCGCTCGGAGTGGCCATGGAAGTGGCGCTGGACCA	21.3	1.4
3667	TGGCGAGCCCTCGGAGAGCGATGGCTTCTCGGAGTGGCAAGGTGGGCTA	18.1	1.6
1280	TCTTGGCCCGCTGCTCTCCCGGCCCGCTGGTCCCGCTGGTCTTGGTGA	30.4	2.9

Gene ID	sequence	Avg.	S.E.M
83990	TTGCAGATGATTATAAAATTCGATTCAACAGACTTACTCCTGGACAAATCA	13.2	0.9
5653	TGGAGAACCACTCCGAGGCCCTGTGTGATGGGGTAACATCCCCTGTGGATCA	36.0	2.5
1244	TGTGGATTCCCTTGGGCTTCTATGGCTCCTGGCCCTGGCAGCTTCTCCA	29.2	2.5
6519	TGCACGACATTTCCGCGAGCTTCGCGCAGACCATGGAACCAATACAGCAGCA	3.8	0.3
10157	TACATAGGACTGGACTGGGCGGCTGGGGGTGAGCTCTCCAAGGGTCTTCA	33.3	2.3
3748	TGCCAGCCATTGACAGCCTGCCATGTCCCGGAAGACAAGAGCCCATCA	26.6	2.1
85358	TGGACTTCCGCACTCGCGATGGGCTCACTGCGCTGCATGTGCCACACGCCA	20.1	1.5
11005	TGCACAGGGAGCATAATCCTGTCCGTGGCCAGATGGCAAAATGCATGAA	16.0	1.3
3850	TGGTCAGCAGCAGCAGCACTTCCGCTCCGCGAGGTGGCTATGGAGGAGTCA	9.6	0.3
7306	TTTGTCCGGGCCCTGGATATGGCAAGCGCACAACTCACCTTTATTTGTCA	42.0	2.6
2563	TCACTTCCATGTGGCTGCGACATGGAATGGCCAAATACCCATGGACGA	26.9	2.4
6862	TGGACCCCAACGCCCATGTACTCCTTCTGCTGCACTTCGTGGCGGGGACAA	21.0	1.9
1586	TGGCGATGGCCACTTTGCCCTGTTCAAGGATGGCGATCAGAAGCTGGAGAA	5.9	0.6
6927	TGGTCAOAGGAGTGGGTGTCTACAACTGGTTTGGCAACCGGCGCAAGAGA	20.3	3.1
11081	TTATGGACCTCATCTTCGTACCTCCGTCTGGATGGAATGAAATCAACCA	12.8	1.2
6493	TGTCCCTGTACGACTCCTGCTACAGATTGTGGGGCTGGTGGCCGTGGGCA	16.4	1.7
79400	TGCTCATCCATGGCAGCCCATGGACAACTCAAAATTCCTCTCCAGGTGTA	7.7	0.7
2538	TCATTAACCTTCTTCTGTTGAGCTTCGCCATCGGATTTTATCTGCTGCTCAA	10.9	1.0
54453	TCCGCTGCGCTGTGAATTTGGGGCCCACTCAAGGAATTTGCCATAAGGA	10.8	0.9
4018	TGCGAAATCCAGATCCTGTGGCAGCCCTTGGTGTATACAACAGATCCCA	43.4	1.3
23316	TGGCCCTTAGTAGAGAAATCAGGAGGCGAGGCTGCTTTTCTGAGTGTTA	6.1	0.7
53630	TCTCTACTTGTCTCACACCATCCCGAATTCACGCAACATGCCTGATCAA	11.9	2.2
4914	TCAACGCTCTGGAGTCTCTCTCTGGAAACTGTGCAGGCGCTCTCTTACA	18.4	0.7
1103	TGCCAAACTGCCCTGCCCCGCTGCGAGCAGACCTGGCCAAGTACCTGCA	14.6	0.3
3762	TGCTCTGGGGCCACCGATTACACACAGTCTCACCTTGGAAAAGGGCTTCTA	14.3	0.9
54221	TGGCTGGACCACTTGTGCTGCTCTGCTCATGGCTCGCATCTCAAGGTACA	26.0	0.7
2675	TCAGGCCAAGAGCAACCAATGGCTGGATGCTGCCAAGGCTGCAACCTGAA	12.3	0.5
2161	TGTGAACATGCCAAGCGGCCCACTGTCTGTGTCACCAACCTCACTGGA	36.5	3.7
6506	TCAATGTTGTGGGTGACTCTTTGGGGCTGGGATAGTCTATCACTCTCCAA	0.1	0.0
89	TTGCGGGCGGCGGGGGCGGCGAGTACATGGAACAGGAGGAGGACTGGGA	26.1	0.6
2904	TCTCACCCCTTTTCGCTTTGGGACCGTGCCCAAGGCGAGCACAGAGAGAA	23.0	1.2
81607	TTTGTGCTTCTGGTGGTGGTGGTGGTCTATGTCCTCCGATACCATCGGCGCA	20.3	0.8
64218	TCAAGATTCTACCTGTGTGCCATCTCGGAGGACAAGGTGATGGAGGAAAA	21.0	1.7
9732	TCCGAGAACCGTTCCATATGGCTGTCAATTGGAATTTGAGATACAGTTCA	5.6	0.5
26960	TCTCGAGATGCCACCTGCTGCTCTGGTACTGGAGTGGGCGGCACATATCA	23.0	1.7
8817	TCCGGCTGGTACGTGGGCTTCAACGAAGGGGCGGCGGGAAGGGCCCA	44.6	2.1
51156	TTTGTGACCTTAGTGAATCTCAGCTACTGGAGAAATCTCCAAGTATCCA	28.6	2.1
27324	TCAGACCTGCCCTAGGCATGCCGATGACTGCTACCTTTCAAGCCCTCA	30.7	3.5
4284	TCCAGTTCGTCTCTGCATCTTGCACATAAGAGGAGGCGGAATGGCCA	39.0	1.1
2318	TGTCCAGTTCCCAAGGCCAAGCTCAAACCTGGTGGCCCTGTTCGATCCAA	24.0	1.7
6330	TGGCTGTCTGGGCGGGGTATCGGCTCTCATCTCATCTGCTGATCAA	26.4	1.5
6869	TACATGGCCATCATACATCCCTCCAGCCCGGCTGTGAGCCACAGCCACCA	36.4	1.0
5284	TGTCCATCAAGTGTACTACCAACCACTCTGTCAACGGGCAACCCGGA	28.3	2.5
6295	TACTTCTCCGGGTCCAGGTGTATCCTCCTGTGGGGCGGCGAGCACCCCA	25.1	1.0
6337	TGTGAGGCTGCCAGAGACTCTGCCATCCTGGAGGAGACAGCTGGGCAA	23.9	1.2
27031	TCCCCCAACACACTGATCTTCCCATTTTGTGGGAAGGCCATGTCAACCA	34.5	1.6
2138	TATGGGCAACACAGTTTACCAAGGAATGCAACAGCTACAGCCTATGCCA	22.6	2.3
640	TGGTGGCTGGCAGGTCACTGTGTCAGGAGAGAGGCTATGTGCCAGTA	11.4	1.5
92359	TCCCTCTTGGCTGCCTTGTCTGCTGCTGGGGCTGGCACTGTGGTGGCGA	26.1	1.4

Gene ID	se quence	Avg.	S.E.M
55615	TACCTGCGCCTGGAGACGCTGGTCCAGAAGGTGGTGTGCGCCATACCTGGGCA	34.3	1.2
2524	TTCCCATGGGCCACTTTCATCCTCTTTGTCTTTAAGGTTTCCACTATATTICA	34.4	1.1
5913	TGGTCCAGATCGACACGGGCCGGGAGCTGGAGGATGCCGACTTCTCTCTGGA	23.6	0.8
60529	TACCCCTGACTCTGACACTGTGGGATGACAGCAGCTACCTGAGTGTCAA	38.3	1.6
1535	TCCTCATCACCGGGGCATCGTGCCACAGCTGGGCGCTTCACCAGTGGTA	23.1	2.3
10840	TGGGGGCGAGCTCAACGTCCTGCGCTTCTGCAGCCAATTTCATCCCATGGA	42.9	3.0
1300	TATAGCAGTAAGAGGAGAGCAAGGTACTCTGGTCCACAGGCCCTGCTGGA	23.9	1.7
9723	TGGCTCGAGACCCCTTACTGTGCTGGGATGGCATACTCTGCTCCCGTATTA	35.2	1.7
1811	TGCTCAGTGAATATTGTTTCTGGTATCAGCACAGGATTGTGGCCGTACTACA	11.0	0.5
1356	TTCCGTGCTGTGGTACTTATTCAGCGCGGAAATGAGGCCGATGTACATGGA	32.9	0.7
525	TTGCCCGCGGCCAGAGATCCCATCTTCTCAGCAGCGGGCTCCCCACAA	41.6	0.9
3663	TCAGCAACCCCATGAGCTGCGGCTCTTCTACAGCCAGCTGGAGGCCACCCA	35.5	1.1
8483	TGCTGCTGATGTGCGAGCACTTCATGCTTCATGGGGCTGTCTCCCTTCCGGA	35.2	5.3
4621	TCAGTACTTTGCAACAATTGCACTACTGGGGACCTGGCCAAGAGAAGGA	2.0	0.5
79755	TGGGCCCTCTCAACCTCTCCAGAAATCAGAGATAAACCCTGGCAGCCACCCA	37.9	2.4
3982	TCTGTGCTGGGTGGGACCATCTCTGCTGGTGGCCATGGCAACAGCCA	21.1	2.5
6608	TGCCAGCCGTACCTCTGCCAGGCCACCCGAGGCCCTGTGCCATCTGGA	39.9	1.5
2917	TCAACAGCTCTCTGGAACACCCCACTCCAGGGCCGCTGCTGATTTTGGCAA	9.6	0.4
3866	TCGACAACTCCCGGTCATCTCGAGATCGACAATGCCAGGCTGGCTGCGGA	33.1	1.0
4143	TTGATGTTCCGCTATGCTACCGACGAGACAGAGAGTGCATGCCCTCAGCA	29.9	2.5
860	TGCTTCAGAACTGGGCCCTTTTCAGACCCAGGCACTTCCCAAGCATTTCA	14.6	3.0
4016	TGGTCCAGACCCCACTATGTGAAGCATCCACTTATGTGACAGAGGCCCA	34.8	3.9
7021	TGGTGTCCCGAGCCACAGCTCGCGGCTCTCCAGCTGGGCTCGGTGTGCCAA	30.6	2.8
1591	TCGGTGCACTGGGCTCGGCATGCTGCTGGAAGCGCTGTACCGCACCGAGA	21.3	1.8
3561	TCCCTCAGTGTTCCTACTCTGCCCTCCAGAGGTTGAGTGTGTTGTGTTCA	5.3	0.8
1950	TTTGCTGTTGGCAAATCTCAAGATATTCGACACATGCATTGTGATGGAACA	39.7	1.5
6564	TTGTGTTTGTCTTGGCAGTGGGATGTACAAGAAGTTCAAGCCACAGGGCAA	28.1	0.7
2623	TACCATATGCCGCTGGGCTACGSCAAGAGCGGGCTCTACCTGCTCTCAA	30.3	0.8
8013	TCTGAAGGGAGAGAGGCTGCTGCTTCCAAACCAAGAGCCCATACAA	35.8	1.7
2048	TTTCCAGTGTCATGAGAOCTCCCTCATGCTGGAGTGGACCCCTCCCGCGA	0.2	0.0
3955	TACCGCACACGCGGGACGCTACGCTCGGCAAGCCAGCCTGGACAGGCCCA	17.1	1.1
79133	TTATGTGCGGCACTTGGGCGGGGAGGTCCAGCGGAGAACTTTGGCCGTA	21.0	2.3
55613	TTGTCCCTGACCTCCGCTGCAAGAATTTCCGGGTGGCCACTTTGTTTTTA	21.8	1.0
6261	TCAAGCTCTGCTGGCCGCCGAGGGCTTCGGCAACCGCTGTGCTTCTGGA	39.2	1.0
2277	TCCCAAGATCTAATCCAGCAACCCCAAACTGCAGTTGCTTTGAGTGCAAA	1.9	0.3
3082	TCTGGTCCCTTCAATAGCATGTCAAGTGGAGTGAAGAAAGATTGGCCA	11.5	1.9
10560	TACCTCGTTATAAACCTGTGTTCTACTGCAGGGGCTCAGCCTTATTGTTA	7.2	0.6
8323	TGTTCAACCAACATTTAAACTTTCTCTGCAAGCATTTGTACCAACCTGCA	14.6	1.6
23426	TCTCATCGATGGAACAGCATGGAGTACTGTACACTTGCAGAGCAACCCA	37.4	1.0
8838	TTGTCACTGGCCCTGCAATGCCCTCAGCAGAGGCCCGTTGCCCTCTGGA	18.0	1.0
64241	TCTCTCTGAAAGTGACAAACAGCCTGTACTTCACTACAGTGGCCAGGCCAA	17.1	1.2
79444	TTCTTCTGCTATGGGGCCCTGCAGAGCTGGAAGCGGGGACGACCCCTGGA	24.9	1.6
1000	TGCTGTCTCAGGCTCCAGCAACCCCTTCAACCAACATGTTTACAATCAACAA	51.1	2.3
2735	TCTGCCGGGTGGGATGATCCACATCCTCAGTCCCGGGGACCCCTTCCAA	39.8	0.8
862	TCCAACAATGCCACTCCOCCAACTACTCAAGGAGCTCCAAGAACCATTTCA	1.0	0.2
79152	TGGTGGACTGGCGAAAGCCTCTCTGTGGCAGGTGGGCCACTTGGGAGAGAA	19.1	1.5
9719	TCTCCAGCAACCGTGTGACCTGCACTGTACACCGTGGACGGCCAGCGGCA	0.4	0.1
56000	TTAGCAGTAGTCTGAACCTGTCAATCCTGGCATGCTTCTTCATCCCATCA	34.0	0.9
84623	TTATCTCCCTGGTGACGTGGAGATGGCCAGAGCATCGTGTGCTGTCGA	47.3	3.4

Gene ID	se quence	Avg.	S.E.M
7223	TGGGACATGTGGCATOCCACTCTGGTGGCAGAGGCTTTATTGTCTATTGCAA	31.5	0.5
3984	TCATGCGATGCTGGAACACCCACACGTGCTCAAGTTTCATCGGGGTGCTCTA	19.3	1.2
1813	TGTCCTGGTATGATGATGATCTGGAGAGGAGAACTGGAGCCGGCCCTTCAA	39.1	0.7
26278	TGTGTACAAACAGCGCTTTTACCCAGAGGACTGGCAGGCAATTCAGAA	18.7	0.7
3680	TCTCCATCTTCTCCCGGGCTCCATCAACATCACAGCGCCTCAGTGTACGA	6.4	0.4
92745	TCTGCCAGCTCGTGGTCTGCTCTGGAGCAAGCCGGTCCAGTTCATGGA	12.8	1.0
1420	TCAGCGAGATCCGTTCCCTCCACGTGCTGGAGGGCTGCTGGGTCTCTACGA	24.3	1.1
288	TGTTATGGAAATTTGGCCOCCACTTACCAAAGGAGGACAGCAACTTGTTTTTA	1.9	0.2
4855	TCTGCTGTGTGAGCCTGGCTATTGGGGGCCACCTGCCACAGGACCTGGA	34.8	0.5
11283	TCTGTAACTCAACATCTTCGCAATCCATCACAAACCCCTCAGTCTGGCCAGA	44.4	1.6
5979	TGGTGGCAACGCTGCTGTCTTCGATGCAGACGTGGTACCTGCATCAGGGGA	24.3	0.7
1757	TGGAGGTGGAGCTTCTGGGCCACACTCGGCGGGTGGTGGCCGGGAGCTGGA	27.2	1.4
5393	TGCTGCAATCGTGGCCTTATGTCAATTCGGAAGACCTGATGTCTCTGTCCAA	19.9	2.6
6948	TCATGACAGCGTGGTGGACAAACCTCTGTATGCTGTGGAACTTCCACCA	26.7	0.7
3779	TGTGTGGCTCATCACTACTACATCCCTGGTACAGACTGTGCTGCCCTCTA	24.5	0.9
7068	TGTCATAGACAAAGTCAAGCGAAATCAGTGCCAGGAATGTGCTTTAAGAAA	16.3	0.8
8029	TTTCACCCACAGCCCTGCTGCAAGCTGGAGACAGACGAGTGCAGCTTCCA	29.4	0.9
2147	TCACGCGGGTCCGGCGAGCCAAACCTTCTTGGAGGAGGTGGCAAGGGCAA	45.4	1.8
2564	TGTCATATGGCCCGCAGCCCGAGCCTCTGGAATAACAGCTCTCTCTGAGGAA	3.2	0.5
7075	TGCGCTGGCGCGCAACGGTTGSCACCAGGTACGCTTCCGCGCTTCTCCAA	20.0	1.0
2925	TTGTGGGCTGTTCGCTTCTGCTGGCTCCCAATCATGTATCTACCTGTA	8.7	0.3
56159	TACTCCATACATGACATGCTTTTACAAAGGAAGGAAGAACTTCTTGCCAAAGGA	15.0	2.2
3662	TCCCATGACGTTTGGACCCCGCGGCCACCTGGCAAGGCCAGCTTGTGAA	27.8	2.2
50937	TCAGCAGAGCCCTCTCTCTGTGGTCCCTGTGGTAGCACTTATCTCTCAGGA	26.6	0.3
1181	TGCCAAAGAGGAAGCTGCTCGGATTCGCTGGAGGGCTGAACCTTGGAAA	19.9	0.9
28	TCTGGAGGGCAATTCACATCGACATCCTCAACAGCAGTTCAGGCTCCA	16.3	0.2
1594	TGCGCACCTGTGACGTGGCTGCCCTGCACTGCTGAGGAGCTGCTGGACA	24.6	1.0
2591	TGGTCCACGTTGCTTAGAAGTGTCCACAGTGTGCTCTATTCTTCACTGCAA	20.8	1.2
4481	TGCGGGAGCAATTGGCTTTCTGTGAAGTCGAGGACTCCAGGATATGCCGGA	24.4	1.1
7010	TATCTGATGCTGAAACATCTCTCACTGCAATTGCTCTGGGTGGGGCCCCA	15.0	0.3
1815	TGCGACGCCCTCATGGCCATGGACGTGATGCTGTGCACCGCTCCATCTTCA	28.5	0.6
10680	TCTATGAATTGGAAGGCGCTTTCTATACCAGGAAGTACCTGTCCCCCGCCGA	8.3	2.8
55145	TCCCCACCTCTCTTACCGCCTCCTGTTTCCAGGTTGATGCTGCTATTGGA	6.0	0.3
5799	TGCCACACTCGGGGGCTCTCTGAGGACCAAGGGTCTGCACTCTTACCTGGA	24.1	1.5
3141	TCCAGCTCCACATAGTGAAGTCCAGAGATTTTAACTTGCTCAAGTCA	21.1	2.1
3195	TCTTTCACAGCCTGCAATCTGCGAGCTGGAGAGCGCTTCCACCGCCAGA	15.9	0.2
631	TTCAATTGAACTCCCATTCCTCTGTTTCCCCAGAGCCATGGAGTCTCTCTCA	30.7	2.5
114327	TCGGACCTACATCAATCACTATCTTATGATGATACGTTGGAAATTCGA	19.0	0.7
6557	TTTTTTTCCAGCAGCTACTGGGATTTCTGCTGGTCCCAATATCTCAGGAGA	24.3	1.8
2294	TCAAGCCATGGGCTCTCTTCCATGCACTGGGCCGCGGGGCTCTACTA	18.1	0.2
6005	TGTTGTGCTTGTGGGCTCCAGTGGGGCACTATTGTACAGGGAATCCTGCA	35.8	1.3
326	TCTACAAGCAGCTGCGGCTCCTGCTCTGAGCCCGCTGCCAGGGCTGGA	19.5	1.0
3886	TCAAGAGCCAGAGGGGCAAGTGGGGCTCCACTGTGTCTCCGAATGTA	28.0	1.5
10243	TTAAAGAGGGGAATGTGTTTTGGCCAAAGGAACCATGGGCCCTCAGA	30.2	1.4
2000	TACCTGATTTACTGCATCTGACTCGGAGTGGAGTTGGAAGAGTTTACA	18.5	1.5
3815	TGTGCTGTGTGTCTGTGTCCAAAGCAAGCTATCTTCTTAGGGAAGGGGAA	21.7	2.1
5890	TCTCACCAGCATTTCTTATCTACTACCTTTCTGCTTTGGACGAAGCCCTGCA	22.4	1.7
4547	TTGATTCAAAGTACAGGGGCAATTCOCATTGTGGGGCAGGTCTTCCAGAGCCA	31.1	0.5
1326	TTGGTTCGATTTTATCTCTCGGGCGCCTTTGGAAAGGTATACCTTGGCACA	4.8	0.6

Gene ID	sequence	Avg.	S.E.M
4318	TCTCTGGGGGCTGCGCTGCTGCTTCTCCAGAAGCAACTGTCCCTGCCCGAGA	5.4	0.4
2571	TCTTCGACCCCATCTTCGTCCGCAACCTCTCGAACGCGGGAGCGGACCCCA	21.2	1.7
611	TGGGGCGGTGGGATGGGCTCAGTACCACATTGCCCTGTCTGGGCTTCTA	18.8	2.4
5626	TGTCCTCTGCGCGCTTTCCAGCTTCTGCCAGAGTCCACTGCTTGCCTCA	3.8	0.8
12	TGCCCTTTGACCCCAAGATACTCATCAGTCAAGGTTCTACTTGAGCAAGAA	16.5	1.4
3062	TGGAACTCATCTGTAGTTTCAGAGAAATGGAAAGCCCTGCAGCTGTITCA	48.1	1.2
4353	TGTCCAAGTCAAGCGGCTGCGCTACAGGAGCTGGGGTGACTTGCCCGGA	3.6	0.3
6584	TCCCAACATCTTGGATCTGCTTCGAACCTGGAATATCCGGATGGTCACCA	8.2	0.7
7512	TGTCATCAGAGCCCTCGCCAGCAGATGCAGACCCAGAACTCTCTCAGCCTA	15.1	0.8
6855	TGGTGGCTGGGGGTGAGTTCGGGTGGTCAAGGAGCCCTCGGCTTTGTGA	23.9	1.0
9118	TACCTGCTCCACCTAGAACTCTCAGTGTACAACTCCAAAGTCTCATCCA	31.7	1.7
3575	TTATCCAGCAGAAAGCTGACACTCTGTCAGAGAAAGCTCCAAACGCGCAGCA	5.2	0.9
2157	TTACTGCTTCTCTCTACTTTACCAATATGTTGGCACCTGGTCTCCTTCAA	33.0	1.7
23414	TACACTGCTGATTCCGCTGATCACTTTCACCAACACCTGTTCTCCCATCTCA	27.8	1.1
4010	TGCTTCCCTCGCGGGCAGACGGACTGCGCCAAGATGTTGGACGGCATCAAGA	18.5	1.4
2103	TCATCAAGACTGAGCGCTCCAGCCCTCTCGGGCATGATGCCCTCAGCCA	9.9	0.4
5105	TGTCCAGGGAAGCCTGGACAGCCTACCCAGGCACTGAGGGAGTTTCTCGA	14.8	1.1
4317	TGTATCCCAACTATGCTTTTCAGGGAAACAGCAACTACTCACTCCCTCAAGA	17.1	2.2
8854	TGCAAGCTTTTATGTGGATTTCAGGGGCTCATCAAACTTTCGATATTA	10.6	0.5
23093	TTTGTGTGCCAAGATCTGCCCCAGCGGGCATCAACTCGGCCAATTTATCCCA	25.4	1.8
383	TTTGGCAATTGGAAGCATCTCTGCCATGCCAGGTCACCCCTGATCTTGA	31.4	2.0
242	TGCACAAAGAGCGGTACGCCCTTCTTCCCAAGGACCCCTTGGTACTGCAACTA	29.4	1.2
6329	TCTACACCTTTGAGTCCCTCATCAAGATACTGGCCCGAGGCTTCTGTGTGA	1.8	0.2
3425	TGCCACAGCAGCGGCTGACAGTACGTCTCAGTGGGACAGCAGCTCAA	12.0	1.6
2587	TGGAATCTCTGGCTGCGCACCACATCATCCATCTCTGGGCTGAGTTTGA	21.5	2.1
6332	TTCAAGATCAATCGGCTTCCATCTGTGTACATTGTCTCCTTCAATTGTA	4.4	0.2
5967	TGCCCCAGGCGCGGATCAGCTGCCAGAGGACCAATGCCATTCGCTCCTA	14.1	0.2
84152	TCAGGAGGGGCAACCATCTCAAGTCAAGAGACCCACCCCTGTGCCTACA	26.7	1.4
5924	TGCGCTACAACGAGGGGCAAGCCCTGTACCTGCCCTTCTGGGCGCAGGA	17.9	0.6
887	TGCTGCTTCTGCTCTTGTCTTCATCCCGGGTGTGGTTATGGCGGTGGCTA	25.4	0.5
9152	TTCTACCTGTTTGCTCCTTTGTGTCTGTACTACCTGCGGCTCCTGCAACA	21.3	1.2
9499	TTACCCACCACTGTTTCAATCAAGTCCAGAGAACATGTCGATTGATGAAGGA	13.0	0.7
115352	TCCAGGAACAGAACAGGCTTACCGCTGCGGGAATCAAGGGCTGGTGCTCA	26.6	1.6
1993	TCCATGCTCGCCCAAGTTCAGCTTCTATCAGAGATGCAAAATTTATATGTCA	22.0	1.6
1124	TGGATTGACGTACACAAACAGTGTTCAGGACCTTCCCAATGACTGCCAA	22.3	0.3
6315	TCTTGCAAGCAGCCCTTTTCTGTGTGTCCAGGAAAGTCCATGCCCTGCA	21.9	0.7
23533	TCACAGGTTCTTCAAACTTCAGTTCTTCTACGTGCTGTGAAGCGAAGTCA	29.2	0.8
2824	TTCTCCGGGTGGCCCTTATTCTGCGGCTGTGGGCAITGTGCTCTCGCAGGCA	17.2	1.8
6514	TCATGGGACCTGGTTTCTACTGTATCACTGCTGTGCTGGGTTCTTCCA	15.8	0.7
50615	TCCACAGTTCGCGCCCAATGCCAGGCTGCGCCTACACCTGCCACATGGA	21.4	2.4
2166	TCAGAGGCGCTGTAGCCCTGCCCTGCTCAGTGGTGAGGAGTTACAGA	27.6	1.3
4868	TGGAGCTGCGTGTGGGGGTGAGCAGCCCTGGCAGTGGGTGCAATGGGCCAA	4.8	0.8
1308	TACAGGAGGCTCACTCACTGCTTCACTCTGCGCAACTCCCGAGGCTCAA	40.6	0.8
6928	TGCAGCAACACACATCCCCAGAGGGAGTGGTGTGATGTACCGGCTGAA	44.9	1.0
2334	TTGGCCACTCCCCAGCCCCACCTGAGTGCAGCCAGCGGGGTTCTGGCA	16.8	0.1
1013	TCTGCAAGGCGCAGCCCGAGCTCTTCAAGCATGAGAGCTCAGAGGAGA	31.3	2.9
5002	TGGAGCTGCTCTCAGCTTCACTGCTATCCCGCAGCAGCAAAAGGGGCCAA	25.4	0.7
695	TGGGAGATCTCTGCTGCTCTCAGACAGCCAAATGCTATGGGCTGCCAA	19.6	0.7
6097	TGTCAAGTTGCGCGCATGTCCAAGAGCAGAGGGACAGCTGATGCAGAA	25.3	2.2

Gene ID	sequence	Avg.	S.E.M
356	TGCTCTGGAAATGGGAAGACACCTATGGAATTGTCTGCTTCTGGAGTGAA	8.8	0.1
7345	TGGCACAATCGACTTATTACGCGAGTGGCCAATAATCAAGACAACTGGGA	10.1	1.0
11281	TCAACACAGCAATCCTCATTCCCTTCAACATGGCGGGACAGCTAGGAGGCCA	4.4	0.2
4359	TTTACACCGACAGGAGGTCCATGGTGTGTGGGCTCCCGGTGACCTGCA	33.4	1.6
5551	TTCAACGCTCCACCCAGCCGCTACCTCAGGCTTATCTCCAACCTACGCA	0.3	0.1
6850	TGCTGCACTATCGCATCGACAAAGACAGACAGGGAAGCTCTCCATCCCGA	20.0	0.6
2042	TCAGTGGTGTGATGAACATTACACACCCATCAGGACTTACAGGTGTGCAA	21.4	1.9
154881	TCGGAGGGGCTCAGTTCACTACAGGCTGTCCACACTGCGGTGTACGAAGA	15.2	3.4
2898	TCCTGGGGTGGCTGCCATCTTCGGGCTTTCACAGCTCATCAGCAAGCGA	32.8	1.1
1723	TCAGTGGTGAACACAGGTTACGGGCCAGACAGCAGAGCAGGCCAAGCTCA	18.4	1.1
4036	TCTGTGCTCTGCTTGAATGCCAGTACCAGTGCCATGAGACGCCGTATGGA	25.0	0.7
10462	TTAGCAACTTCACCTCAAACTGTGGCGGAGATCCAGGCACTGACTTCCCA	8.9	1.9
80144	TTGCAATCATGGGGAAGTCCGATGTACCCCCCAACATGCCACCGCTGTCA	2.5	0.2
2357	TCATCACTTATCTGTGATTTGAGTACCTTTGTCTCGGGTCTTGGGCAA	21.2	0.9
88	TCTGCCAAAGAGGTCTGCTGCTTTGGTGTGAGAGGAAAGTCTCTTATA	9.5	0.3
5143	TCCTGGCCAGTCTGGGACCGTTGAGGCAACGTGGCGGCCCTTGCCGCCA	16.2	2.2
5083	TCAGCATCCAGCATGGCTCCTTACCTACCCAGCCCAAGTGTGCTTACA	27.4	2.0
1258	TGGGTCCAGAGGGTGTGCTCAGCCCGAGGAGCCCTCGGAAGACCAAGA	39.5	1.4
4057	TACTCTGCCAGACAACACTCGGAAGCCAGTGGACAAGTTCAAAGACTGCCA	27.6	1.3
7187	TCTATTGTGCGAATGAAGCAGAGGTTGTGAGAGCAGTTAATGCTGGGACA	42.9	1.3
23211	TCAGCCCAAGTGAAGAGGTACCGCAAGTACAGAGATACAGCCCCCAATA	12.3	0.8
5191	TCGCTGCACTCTGGGACACTGCCAAGCTGAGGGCCACTGCAAGTCTATA	22.4	1.1
27	TGCAACCGAAGAGGTGACTGCAAGTTGTGCTGCTTACATGGCCACTCAGA	37.0	1.2
1773	TTGCCATTGTTCCCTGCAATGCGGCCCGGGGAGCAGTAGCCGAGATCGA	23.2	1.5
5452	TGCTTGTGCCAGGCCACACCTCCAGCCACCTGCTCAGTTTCTGTCTACGCA	12.1	0.3
7392	TTGCCATTATTCACAGGCTCCAGTGTGGAGATACTACGCTGTGCTGATA	25.3	1.9
54457	TGCACTGCTGATGGTGATATCCACCTTTCTCAGAAAGAACAGCTGCCTCTA	25.0	0.9
9060	TTGCCCTGGAGAGTACCTTGTCTCCATGCCATCCCTTGTACTCCCTGGA	30.5	3.2
51151	TTGAGGTGCCCTGGGTACCTTTGGGTGCTATAGACTGGGCCATCTGGA	19.6	0.7
4038	TCCAGTGTGCTTGAAGCTGCTATCTCGACATCTAACCCTGCGATGGCGA	26.6	1.4
6583	TTCTTCGTAGGCGTGTCTCTCGGCTCCTTCGTGTCCGGGCGCTGTCAGACA	21.7	0.7
5608	TCAAGCCTTCTAATGTACTCATCAATGCTCTCGGTCAAGTGAAGATGTGCGA	19.9	1.5
1814	TCTGGAGCTGAAGCGTTACTACAGCATGTGCCAGCACTGCCCTGGGTGGA	23.8	1.0
54806	TTTCAGAGTTTACCTCATCCTTCTTTGTTTACAGGCTAAATTCATCCA	1.9	0.5
229	TCTCTGTGGACAGTTCCATCAACCAGAGCATCGGGGTGTGATCCTTTTCCA	23.7	3.8
64131	TCCCAAGAGCTTTGAGAATGTGCAACAGCAACTTCGCAACCCAGGACTCA	18.4	0.8
2780	TGCTGCACTCCATCCTGGCTATCATCCGGCCATGACCACTGGGCATCGA	25.2	2.1
6683	TACTCAGGAAGTGACCTAACAGCTTTGGCAAAAGATGACAGCATGGGTCTTA	9.1	0.7
3670	TTGGTTGCGGCAATCAGATTACGATCAGTATATTCTGAGGTTTCTCCGGA	7.7	1.2
338	TGCTCCACTCACTTTACCGTCAAGACGAGGAAGGCAATGTGGCAACAGAA	22.9	1.2
5091	TCCGTGTGTTCCGGGCTGCACGAGCTGGGCATCCGACCGTAGCCATCTA	16.0	0.5
3785	TCTCTGCTGTTTCTCTGCTCTGTGCTGTGTTTCCACCAATCAAGGA	22.2	1.1
4056	TCTACCTGTTCCGGGCTCTCGCTACTTCCAGGCTACGCGCGCTCCGGCA	23.1	1.0
5630	TCCGCGCGAGTACGAGAGCATCGCGCGAAGAACCTGCAGGAGCGGAGGA	13.5	1.0
84000	TGGCAGGCGCACAGGGATCAGGTACAAGGAGCAGAGGAGAGCTGTCCCA	31.1	1.2
5654	TGGCTAGTGGGTCTGGGTTTATTGTGTGCGAAGATGGAGTGTGTCAGAAA	16.8	1.0
5618	TGTCCAGGTTCCGTGCAACCCAGACCATGGATACTGGAGTGCATGGAGTCCA	30.3	0.8
3587	TTGCCCTTGTCTGCTGCTCTCCGGAGCCCTGCTCTACTGCTGCTGGCCCTCA	29.1	1.1
5308	TTCAATGGGCTCATGCGAGCCCTACGACGACATGTACCCAGGCTATTCTTACA	22.7	1.2

Gene ID	sequence	Avg.	S.E.M
4211	TACTTGTACCCCGGAGCCGGGGGTGGCGGGGGGACGTCTGCTGTCTCA	29.5	1.3
1823	TCCAAAGCAATTTCTCCATACACCAGATACCGGTGTATCACCACACTACA	0.6	0.1
5139	TTTCTTCCACTTGGACCAACCACTCGGCCACAGGTCTACCCACCTTGGGA	5.4	0.5
4624	TGCTGCTGTGTCACCAACAATCCTACGACTACGCCCTTGTGTCTCAGGGAGA	4.2	0.3
579	TGTGCAGCGGGGCGCGGGCGGGCGGAGCGGGCGGCGAGGCTCGCGGA	3.1	0.1
84839	TGGCGCGCCAGGAGCGGCTGGAGTCAGGCTCGGGTGGCTGGCAGCTCCGA	19.0	3.8
25833	TCCACAGCAACAAAGCGGTCTCTCTCTCCACAGACTGGGCGGGACTGGCA	5.8	0.6
22806	TCAATAACGCCATCAGCTATCTTGGCGCCGAAGCCCTGGCGCCCTTGGTCCA	27.0	1.2
38	TGCCCTTGAAGCAAGGAGAATAAGGTCTTTCAGTATTTGCAATGGAGGAGA	25.3	0.5
189	TCTCTGGCTCGGGACACTGTGCCCTGGAGGCGCCCTGGTCAATGTGTGGA	7.2	0.1
247	TGCTTGGATGAAAGACAGTGGAGACTTGGAGCTCAATATCAATACTCCA	34.1	0.6
284	TGGGATAAAGTGGCACTACTTCAAGGGCCAGTACTCCTTACGTTCACCA	40.1	0.3
366	TTGCCCAAGCTATTCTCAGTCGAGGACGTTTGGAGGGGTCTCACTATCAA	20.2	0.3
393	TTGCACTGCTGGGCGGTCTGCTGCGACACAGCGGCGAGAGCGCGGAGA	30.3	1.4
411	TCAACAGTGTGCTCTTGATTTTCGAGATGGCGAAGAGTTGCAACAGGATA	8.0	0.4
487	TGGGCACTGTGGCCACCACTGGTGTGGGCAAGGATTGGGAAGATCGAGA	29.3	1.5
540	TGGTGCAAACTACAGATGGTACACCTACATCTGTGCAGGAAGTGGCTCCCA	25.1	0.5
673	TTGGTTGGGCACTGATATTTCTGGCTTACTGGAGAGAAATTGCATGTGGA	6.8	0.1
1272	TTGAAGAGCAGCCCAATCAATAOCTTTATCCAGAGCAATCACTGGAGGAAA	45.3	0.8
1287	TCTCTGGCCCTCTGGATTTCTGGAGAAAGGGTCCAGAAAGGTGATGAA	10.3	0.7
1380	TGTTGGTACCGTGATAAGGTACAGTTGTTGAGGTACCTTCCGCCCTCATGGA	39.9	1.2
1410	TGGAGAGGACAGGTTCTCTGTCAACCTGGATGTGAAGCACTTCTCCCGAGA	32.4	1.0
1497	TCTTGGCTGCCATCGGCTTCTGGTGTCTGCGTGGCTCTTCGCAITTTGTCA	21.3	1.0
1636	TATGACCGGACATCCAGGTGGTGTGGAACGAGTATGCGAGGCCCACTGGA	15.5	1.8
1730	TTACAAGTTTITAGAAATCTACTGTCAAAAAAGAAAACTCTTATTCAACA	13.6	0.1
1821	TTGTGAGCTCAGCTGCCCTACAGGGGATGTGGCCCTGGTGAACAGGAGA	52.1	1.8
1910	TTACAAGCAGCAAAAGATTGGTGGCTATTGAGTTTCTATTCTGCTTGCCA	29.2	0.5
2070	TTCTTTCTACACGACGCTCAAAACATGTCTGCTATGCAAGCCAGACTCA	34.7	1.3
2162	TATGAAACACACAGCTGATTGTGCGCAGAGGGCAGTCTTTCTATGTGAGA	17.6	0.7
2218	TGGAGCTGGTTTGTCAAAATCCAAGGAAGCCGAATTGGATTGTATAGACA	30.8	0.6
2259	TACCAAGTGGGACTACGTGTTGTTGCCATCCAGGGAGTGAACAGGGTTGTA	30.4	0.9
2261	TGGACACGGCCCGAGCGGATGGACAAGAAGCTGCTGGCCGTGCCGCGGCCA	36.4	0.8
25974	TGTAGCTGGGCTGCTTACTACTACCAACGACAAGATGTGAGGCTGACCCA	48.6	1.5
25984	TGTTATCCGAGACCCAGTCTCGGTACTCCTGCAAGCTCCAGGACATGCRAGA	22.2	0.9
27130	TGCATGCTGCTCTTTCTGGCCATGTCAGCACCGTGAAGTTATTACTGGA	5.1	0.4
27241	TATTTTGAACACAGGGAGTCAAGATTTTGCATGTTCTTTTCGGGATCTA	14.2	2.0
27286	TGGTGTATACATTAAATATCCAGGATGGAGAAGCCACATGCTACTCACCGA	33.6	0.7
51314	TTTATTGAGAGACAATGGCTTGCATACTGAAACAAATTAATGGGACCAAGA	40.9	1.3
51733	TGGAATACAGCCGTGGTATCTCAATTCGGAGCAGTCTGGGAAAGACCA	28.9	0.5
55084	TTCTGCAGCGAGAAGTGTCTTGGCGCTGCGAGAGGCTACTTCAAGAGAA	20.1	0.2
55118	TGGCTGACTTCAACCGTGTATGGCAAGTGACATCGTCTATGGCAACTGGAA	37.5	2.5
55198	TTGGCAGACCAATGGTTCTAOCATCATACAATTCCGAGAAAAAGGATCTCA	32.4	0.7
55285	TGTCTAAGAAAGAGAGCTTTGCTCTGCTACTATGTACAAGCCCTTTGGGAA	9.2	1.1
55607	TGCAATCAACATGCGATTACAGAAGCAGAGATTCAAAATTGAAGACCAA	21.2	0.7
56246	TGGACTATCTGGACCTCAITCCCGTGGAGGAGAAGAGCTGAAGCCCAACA	27.9	0.5
56649	TGGCCTTGGCAGGTGACATCCAGTACGACAAACAGCAGCTCTGTGGAGGGA	0.3	0.1
57623	TGGAGAGGCTACAGGAATAAGCCATTACGACAGACCTTATGAGCAACCA	32.2	1.9
59272	TGCCCTTTTAAAGAACAGTCCCACTTGGCCAAATGTATCCACTACAAGAA	20.6	0.5
64093	TGGGAAGCCCATCAGTGGCTCTTCTGTGCAAGATAAACTCCTGTATGTTCA	7.1	0.3

Gene ID	sequence	Avg.	S.E.M
64102	TATGACATGGAGCACACTTTCTACAGCAATGGAGAGAAGAAGATTAC	24.7	0.7
80263	TCACCTGCTTTGTAAAGGATGCCGAGGAGAAATCATGGGAGGGGAGGAGA	17.3	0.9
80704	TATGTCGCTACAGCCAGTCATCATCTTGAAGGTATCAGTTTCATCATT	11.7	0.7
85320	TACTCTCCAAGATGGCCCTGGAGTCAGCAAGAGAGAAATCCTGAGGCTCCA	40.8	1.4
91147	TCAAGTGTCTAGTTCTACCTCTCTTCCACAAATTCAGTTTAAAGGAGAA	16.2	0.6
121340	TTTGGTGGCTCTAGCCCTCTGGGGACTCAACAACCTCTGGGCAAGCAGGCA	37.2	1.0
123263	TGGAGTGGTCACAATGCCCTTCCCATCACAAAGGACTGCCAGTGAAGCA	6.6	0.2
133396	TGGCTACAAATATGGTACTATCCAGAAAGCAACACTAACCTCAGAGAAACA	35.1	1.2
135935	TCTCCTTTGACACTAACATCACCAGACCTGGAGGTACCTGGGACACCAGA	40.7	1.3
137682	TTGATATATCCACCCATCTTTACAGCAGAAGAATACATTACTTCCATTATA	3.9	0.4
139105	TGTCGGCTAGATACCTTATTCAGAACTCTTCAGAAAGATGTCTGGTCCA	38.7	1.6
139135	TTGCAGAGGTTGAGCAGTATGGACCAAGAAACGTTTACATGTTGTAGA	28.3	1.9
145173	TCATCCAGAGTCAAGATAATTTCTTTTCATGCAAGAGAGCAGAGCGATTAA	32.8	1.4
146956	TACACACAGAAAGCCAGGCTCAAATGTGAGAGCTGGAAGAGCTGGCCGA	30.4	1.8
147372	TGGCCGAGTGTCTGTACTTGTATCCGGGATACCGATATGACCGGGAGAGA	31.1	0.5
158747	TACGTTAAAGACAATGCTTTCAATATGTGAGATAAAACAGTGAAGATATA	16.8	0.4
166012	TGGGCTCCAGCTGGCTTGGTGGGAGAGAGAGCCCTGCAGAGCTCTA	19.6	0.3
33	TCTTTTCTCCAGAGCATGACATTTCCGGAAGGTGTAAAGAGTTTTCCTA	1.6	0.1
183	TGCCCCAGGCTGAGCTGCCCGCCATTCTGCACACCGAGCTGAACCTGCAAAA	26.5	0.6
240	TGGAGAACCTGTTCATCAACCGCTTCATGCACATGTTCCAGTCTCTTGGAA	19.5	0.5
287	TTTACTCCTCTAGCTGTGGCACTCCAGCAAGGACACAACCGGCGGTGGCCA	14.6	0.3
421	TGGAGACTGCAATGTGCACTGGCCGCCAGCATCTGGCCTCGGTGAAGGA	21.1	0.6
492	TGGCGGACAAACCAATGACCTGGAGAGCGCAGGCAGATCTACGGGAGAA	25.6	1.1
649	TTCCAGGAACACATTTACCCOCAGCTGCCAGAGCACCACGGGAGCCCTCA	42.8	0.3
725	TCACGTTATGTGCAATGACCACTACATCCTCAAGGGCAGCAATCGGAGCCA	12.5	0.6
889	TGTACTAATGAAAAATTTCTCTGGATGGAGAGAGATGGGACAGAGAGCA	16.8	0.5
1009	TTGATGACAAATCAGGGAACATTATGCCACCAAGACGTTGGATCGAGAGA	39.2	1.4
1056	TCTTCAAGGGCATCCCTTTCGAGCTCCCAACCAAGGCCCTGGAATACTCTCA	16.9	0.4
1260	TTGCGCTGGTGGGATCATCAGAGAATGGGCCAACAAGAAATTCGAGAGGA	39.6	1.0
1297	TGTCCCAATGCCTGTCCACAGGTGCTCAGGATATCCAGGCTACCAGGCA	44.8	0.5
1690	TTGCTATCATATGTTTACACAGGCTTGGACATCAGGAAAGAGAAAGAGA	45.5	1.0
1824	TTTTGGAGGATGGTTCAGTCTATACAACAATACTATTCTATTGTCTCGGA	11.8	0.4
2068	TTTATAGCTACCACTACCTCTGGACCCCAAGATTGCAGACCTGGTGTCCAA	29.0	0.6
2670	TGGAGGAAGAGGGGAGAGCCTCAAGGACGAGATGGCCCGCCACTTGCAGGA	36.0	0.8
2694	TTCCCTCAGCACAGGAGCCCTTGGTCAATGGAATACAAGTACTCATGGAGAA	38.9	0.7
2892	TTGCGGTGAGTTATACAACCAACCAAGACCAACCGGAGAGCCCTTCCA	26.4	0.7

2. List of TALEN activity in the T7E1 assay

	Gene name	TALEN binding site	Number of CpG sites			Mutation frequency (%)
			Left site	Right site	Total	
Genome-scale study	AGXT	TCTCTGGCTGGGACACTGTgccttgaggcgGCTTGGTCAATGTGCTGGA	1	1	2	2
	ALOX15B	TGCTGGATGAAAGACAGTggaagacttggaGCTCAATCAAACTACTCCA	0	0	0	14
	APOB	TGCTCCACTCACTTTACGGTcaagacgagggasGGCAATGTGGCACAGAAA	1	0	1	8
	AQP9	TTGCCCAAGCTATTCTCACTgaggaagttttGGAGGGGTCACTACTATCAA	0	0	0	34
	ARG1	TTTGGCAATTGGAAGCATCTctggccatgccaGGGTCCAOCTGATCTTGA	0	0	0	35
	ARHGAP4	TTGCACTGCTGGGCGGTGCTgtcgagccagGCGCAGCAGAGCCGGGAGA	1	3	4	1.5
	ARSB	TCACACGATGTCTCTTGATttttagagatggGGAAGAGTTGCAACAGGATA	1	1	2	1.6
	ATP2A1	TGGGCATGTTGGCCACCACTgggtgtgggcaGAGATTGGGAAGATCGAGA	1	2	3	4.4
	ATP7B	TGGTCCAACTACAGATGGTaccctacatctGTCCAGGAAGTGGCTCCCCA	0	0	0	2.2
	OPN1SW	TGGGGCGTGGGATGGGCTcagtacccattGCCCCCTGTCTGGGCCCTTCTA	1	0	1	46
	BRAF	TTGGTTGGGACACTGATATTtcttggttactGGAGAGAAITGCAITGGA	0	0	0	16
	CLCN4	TTGGGAATAGCGCTCtGTTctctttttatgtgGAATACAGAGCCCTGGTA	2	1	3	7
	CNTN1	TTGAAGACAGCAATCAATaccatttatccaGAGGAATCACTGGAAAGAAA	0	0	0	34
	COL4A5	TCCTCTGGGCTCTCGGATtctctggagaaGGGTGAGAAAGGTGATGAA	0	0	0	3.5
	CR2	TGTTGGTACCGTGATAAGGTacagttgttcagGTACCTTCGGCTCAITGGA	1	1	2	8.5
	CRYAB	TGGAGAAAGACAGGTTCTCTgtcaactggatGTAAACACTTCTCCCCAGA	0	0	0	3.1
	CTNS	TCCTGGGCTGCCATGGCTTcttggtgtctggGTTGGCTTTGGCAITTGTC	1	2	3	1
	CYP1A2	TTGTCCAGCTGTTGTTCCCTTctggccacagaGCTTCTCTGGGCTCTGCCA	0	0	0	20
	CYP27B1	TGGCAGCGGTGTACGTGGCTgccccgtgactGTGAGGAGCTGCTGGACA	3	3	6	<0.5
	ACE	TATGACCGGACATCCAGGTgggtgtggaagcGTATGCGAGGCAACTGGA	1	1	2	2.2
	DIAPH2	TTACAAGTTTITAGAAAATCTactgtcaaaaaGAAAAACCTCTATTCAACA	0	0	0	3.6
	SLC26A3	TGCTCAGTGATATTGTTTCTggtatcagcacaGGGATTGTGGCGTACTACA	0	1	1	9
	DRP2	TTGTCACTCACTGCCCCacagggggtgtGGCCCTGGTGCAACAGGAGA	0	0	0	16
	EDNRB	TTACAAGACAGCAAAAGATTggtggtatttcaGTTTCTATTCTGCTTGCCA	0	0	0	14
	EYA4	TTCTTTCTACACAGCAGCTcaaaacatgtctGCTATGCAGGCCAGACTCA	0	0	0	26
	F13A1	TATGAAAACCAAGCTGATgttcogcagaggGCACTCTTCTATGTGCAGA	0	0	0	17
	FAAH	TCAGAGGGCTGCTAGGCCCTgccccgtctcagCTGGTGCAAGGTACACA	1	0	1	33
	FKTN	TGGAGCTGGTTTGTCAAAATccaaaggagcgGAATTGGATTGTATAGACA	0	1	1	54
	GPC4	TGGTCAGCGAACGTGCAATcatttgcaagctGCTTTGCTTCAAGTTACAA	1	1	2	13
	FGF14	TACCACTGGGACTAGTGTgtgtgcatccagGGAGTGAACAGGGTTGTA	1	0	1	34
	FGFR3	TGGACACGGGCGAGCGGATggacaagagctGCTGGCGTGGCGCGGCCA	3	3	6	<0.5

	Gene name	TALEN binding site	Number of CpG sites			Mutation frequency (%)
			Left site	Right site	Total	
Genome-scale study	MMACHC	TGTAGCTGGGGCTGCTTACTactaccaacgacAAGATGTGGAGGCTGACCCA	0	0	0	31
	INVS	TGCATGCTGCTCTCTTTCTggccatgtcagcACCGTGAAGTTATTACTGGA	0	1	1	3.7
	SRPX2	TGGTGTATACATTAAATATccaggatggagaAGCCACATGCTACTCACCGA	0	1	1	19
	TXNDC3	TTTATTGAGAGACAATGGCTTgcaataactggaACAATTACTGGGACCAAGA	0	0	0	14
	UPB1	TGGAATACAGCCGTGGTGATctccaattccggAGCAGTCTGGGAAAGACCA	1	0	1	4
	SOBP	TTCTGCAGCGAGAAAGTGTCTTgcgccctgcccAGAGCCTACTTCAAGAGAA	1	1	2	7.7
	CRTAC1	TGGCTGACTTCAACCGGTATggcaaatggacATCGTCTATGGCAACTGGAA	1	1	2	18
	APPL2	TTGGCAGACACAATGGTTCTacctatcatacaATTCGGAGAAAAGGATCTCA	0	1	1	40
	RBM41	TGCTAAGAAAGAGAGCTTTgctcctgtactATGTACAAGCCCTTTGGGAA	0	0	0	7.2
	PPP1R9A	TGCAATCAAAATGCACTTgacagaagcagagATTCAAAAATTGAAGACCA	0	0	0	29
	NXF3	TTAGCAGTAGGTCTGAACCTgtcaatcctggcATGCATTTCTCATCCCATCA	0	0	0	18
	MRAP	TGGACTATCTGGACCTCATTccgtggcagcagAAGAAGCTGAAAGCCACAA	0	0	0	23
	TMPRSS4	TGGCCTTGGCAGGTGAGCATccagtcagcaaaACAGCAGCTCTGTGGAGGGA	0	1	1	22
	ZFAT	TGGAGAGGCTACCAAGAAATagccattcagcAGCACCTTATGAGCAACCA	1	0	1	7.6
	ACE2	TGCCTTTTAAAGGAACAGTccacacttgcacAAATGTATCCATACAGAA	0	0	0	10
	SMOC1	TGGGAAGCCCATCAGTGGCTcttctgtgcagaATAAACTCCTGTATGTTCA	0	0	0	4.1
	TNMD	TATGCATGAGAGCACACTTctacagcaatggAGAGAAGAAGATTACA	0	0	0	15
	SLC19A3	TATGTCGGCTACAAAGCCAGTcatcatcttgcaAGGTATCAGTTTCATCATT	1	0	1	33
	PPP1R1B	TCAGGAGAGGGGACCACTCTcaagtgaagagACCAACCCCTGTGCCTACA	0	0	0	33
	ABCC11	TACTCTCCAAGATGGCCCTGgaggtcagcaagAGAGAAATCCTGAGGCTCCA	0	0	0	39
	TMEM67	TCAAGTGCTCAGTTCTACTctcttctacaaaTTTCAGTTTAAAGGAGAA	0	0	0	26
	SP7	TTTGGTGGCTCTAGCCCTCTgcgggactcaacAACTCTGGSCAAAGCAGGCA	0	0	0	30
	EVC2	TGCTCTGTCTGGACAGCATTgctggtctcaccATTGGGACTCTGTGGGAAA	0	0	0	39
	IL31RA	TGGCTACAACATATGGTACTtccagaaaagcaACATAACCTCACAGAAACA	0	0	0	9.3
	NOBOX	TCCTCTTTTGACACTAACATcaccagacctggAGGGTACCTGGGACACAGA	0	0	0	23
	PASD1	TGCGAGAGGTGAGCAGATGgaccacaagaaaACGTTACATGTTGTAGA	0	1	1	6.9
	B3GALT1	TCATCCAGAGTCAAAGTAAATcttttcatgcaAAGAGAGCAGAGCAGTAAA	0	0	0	32
	CBBE1	TGGCCGAGTGTCTGTGACTTgttatccgggatACCGATATGACCGGGAGAGA	1	2	3	2.6
	CHST13	TGGGCTCCAGCTGGCTTGGTggggagaagagaAGCCCTCTGCAGAAAGCTCT	0	0	0	11
	ACADL	TCCTTTTCCAGAGCATGACAttttccggaaaAGTGAAGGAAGTTTTCCTA	0	0	0	0.6
	ACHE	TTTCTGGTTTACGGGGCCCGaggtctcagcaaaAGACAACGAGTCTCTCATCA	1	1	2	10
	ADD2	TGGAGAAGGGCAGCAGCTGCTtccagtgagACCACAGGCTCTGTCTGCA	0	0	0	22
	ALOX5	TGGAGAAGCTGTTCATCAACcgcttcatgcacATGTCCAGTCTTCTTGAA	0	0	0	33
	ANK2	TTTACTCTCTAGCTGTGGCactccagcaaggACACAACAGGCGGTGGCCA	0	1	1	7.3
	ARVCF	TGGAGAGTGCATATGTGCACtggcccgccagcATCCTGGCCTCGGTGAAGGA	0	1	1	6.3
	BMP1	TTCCAGAGAAACACTTCTTACCccagctgccagAGCAGCAACGGGAGCTCTCA	0	1	1	25
	BMP2	TTGGAAGAAGTACCAGAAAAGagtggggaaaacACCCGGAGATTCTTCTTTA	1	1	2	1.4
	C4BPB	TCACGTTTATGTGCAATGACcactacatctcAAGGCGAGCAATCGAGGCCA	1	1	2	9.1
	KRIT1	TGTACTAATGAAAAAATTTCTcttgatggagAGAAGATGGGACAGAGAAGCA	0	0	0	2.4
	CD40	TCGCGAGCCCTGCCAGTcggcttctctctccAATGTGTCTATCTGCTTTGGA	2	1	3	4.8
	CDH11	TTGATGACAAATCAGGGAACattcatgccaccAAGAGTGGATCGAGAAGA	0	2	2	17
	CEL	TCCTCAAGGGCATCCCTCTCgagctccaccAAGGCCCTGGAAAACTCTCA	1	0	1	3.2
	CHAT	TGCCCAAACTGCCCGTGCCCGcgtgcagcagACCTGGCCAAGTACCTGCA	1	1	2	16
	CH13L1	TTGACCGCTTCTCTCTGACccacatcatctacAGCTTTGCCAATATAAGCAA	1	0	1	30
	CHN2	TGGATTGAACGTACACAAACagtggttccaagcACGTTCCCAATGACTGCCAA	1	1	2	4.4
	CNGA2	TTCCGCTGGTGGGATCATCagagaatgggcccAACAAGAAATTCGGAGAGGA	1	1	2	6.2
	COL9A1	TGTCCCAATGCTGTCCACAGgtgctcaggATATCCAGGCTTACCAGGCA	0	0	0	12
	COCH	TTGCTATCATGTTTTACGagaggttgagCATCAGGAAGAGAAAGCAGA	0	0	0	14
	DNM1	TGGCAATGTCAACCGAGCTCccctggctcttgacAGTGGTCAATGCAACCA	2	0	2	5.5

	Gene name	TALEN binding site	Number of CpG sites			Mutation frequency (%)
			Left site	Right site	Total	
Genome-scale study	<i>DSC2</i>	TTTGGAGGATGGTTCAGTCTatatacaacaatACTATTCTATTGTCTC <u>GGA</u>	0	1	1	3.3
	<i>EGF</i>	TTTGCTGTTTGCCAAATCTCaagatatttcgacACATGCATTTTGATGGAAACA	0	0	0	29
	<i>ERCC2</i>	TTTATAGCTACCACTACCTCctggaccccaagATTGCAGACCTGGTGTCCAA	0	0	0	10
	<i>EYA1</i>	TATGGCAAACACAGTTTACcaccaggaatgcaACAAGCTACAGCCTATGCCA	0	0	0	4.9
	<i>FUT2</i>	TTCCCATGGCCCACTTCATCctctttgtctttAC <u>G</u> GTTCCTACTATAITTTCA	0	1	1	17
	<i>GIF</i>	TTCCCTCAGCAGAGGCCCTtggtcaatggaATACAAGTACTCATGGAGAA	0	0	0	17
	<i>GRPR</i>	TTGTGGGCTGTTCCTTCtgctggctccccAATCATGTCTACCTGTGA	1	0	1	3
	<i>TRPC6</i>	TCCATTTGGTATGAGAACTcttctggttttaacACAGCAGACAATGG <u>G</u> GTCA	0	1	1	5.7
	<i>TYRPI</i>	TTTGTC <u>G</u> GGCCCTGGATAIggcaaaagcgacAACTACCCCTTTATTTGTCA	1	0	1	27.7
	<i>USH2A</i>	TCTGGAGATCTTCTCAGATTgcatgcccaatcACATTGC <u>G</u> GTGCCCTGGCA	0	1	1	8.3
	<i>XRCC2</i>	TTATCACTAACACAGCAGATgtatacttcccaAATCAGAAGTGGCCTGGAA	1	0	1	10.6
	<i>FOXN1</i>	TGCCGGCTTCAGCTGCTC <u>G</u> TcattttgtgtccgACGGCCCTCCAGAGAGGACA	2	1	3	15.2
	<i>PRSS12</i>	TTTGAAGGCACAGTGGAGTatatgcaagtggAGTTTGGGGCACTGTCTGTA	0	0	0	21.6
	<i>TNFSF11</i>	TTAATCAGGATGGCTTTTATtacctgtatgccAACATTGCTTTGAGCATCA	0	1	1	27.2
	<i>CHRD</i>	TGGCAGGCAGGCCAGGGCTgcatcagtggaACACATTGCTGCCAGGAAGA	0	0	0	35.9
	<i>IL18RAP</i>	TTTACCAGAGCCACAGAAATcacattttctgccACAGAAATC <u>G</u> ACTCTCACCA	0	1	1	25.6
	<i>WISP3</i>	TTGTACTGGCCCTGCAAAATgcccctcagcagaAGCCCTGTTGCCCTCCTGGA	0	1	1	21.3
	<i>ALDH1A2</i>	TGCAAGCTTTTTATGTGGATttgcaggcgctcATCAAAACCTTTTGATATTA	0	1	1	10.2
	<i>SGCE</i>	TGTCTCTTTTGTTCATGTTTtggaaagagaatATTTAAGGGGAATTTCCA	0	0	0	14
	<i>PAPSS2</i>	TTGCCCTGGAGGAGTACCTTgtctcccatgccATCCCTTGTACTCCCTGGA	0	0	0	16
	<i>SLC6A5</i>	TTCTACCTGTTTGCTCCTTgtgtctgtactACCTGGGGCTCCTGCAACA	0	0	0	7.8
	<i>REPS2</i>	TACATTGCCCTGAAATTAATgtctgcagcacaATCTGGCTCCCGGTACGGA	0	2	2	25.6
	<i>KL</i>	TCAAGCTGGATGGGGTGGATgtcatcggtatACCGCATGGTCCCTCATGGA	0	1	1	8.8
Backup TALENs	<i>CYP27B1</i>	TGTGCTGGGCCCTGGGGC <u>G</u> ccctctgcgagACTGGGACCAGATGTTTGA	1	0	1	6.6
	<i>FGFR3</i>	TGGCCGTTGGCCATCCTGGCCGgagcctcctcgGAGTCCTTGGGGAAGGAGCA	3	1	4	28
NF-kB	<i>SMEK1</i>	TTAGCTTTCAAGAAAAAGCTggatgtgatgaaATTTGGGAGAAAAATATGTA	0	0	0	20
	<i>TNFR1</i>	TGGTGGGAATATACCCCTCaggggttattggaCTGGTCCCTCACCTAGGGGA	0	0	0	22
	<i>IL1R1</i>	TTGCTCCCTGTCTCTTAACccaaatgaacacAAAGGCACATAAATTGGTA	1	0	1	18
	<i>TNFR2</i>	TGGCATTTACACCTACGCGCCggagccgggAGCACATGCGCGCTCAGAGA	1	1	2	28
	<i>p50</i>	TCATACAATATTTAATCCAGaagtatttcaacCACAGATGGCACTGCCAACA	0	0	0	26
	<i>IKKA</i>	TCTCAAAATAGCAATTAAGTcttgtgcctagAGCTAAGTACCAAAAAACAGA	0	0	0	25

VI. References

- Adamson B, Smogorzewska A, Sigoillot FD, King RW, Elledge SJ. 2012. A genome-wide homologous recombination screen identifies the RNA-binding protein RBMX as a component of the DNA-damage response. *Nature cell biology* **14**(3): 318-328.
- Bedell VM, Wang Y, Campbell JM, Poshusta TL, Starker CG, Krug RG, 2nd, Tan W, Penheiter SG, Ma AC, Leung AY et al. 2012. In vivo genome editing using a high-efficiency TALEN system. *Nature* **491**(7422): 114-118.
- Bibikova M, Beumer K, Trautman JK, Carroll D. 2003. Enhancing gene targeting with designed zinc finger nucleases. *Science* **300**(5620): 764.
- Birmingham A, Anderson EM, Reynolds A, Ilsley-Tyree D, Leake D, Fedorov Y, Baskerville S, Maksimova E, Robinson K, Karpilow J et al. 2006. 3' UTR seed matches, but not overall identity, are associated with RNAi off-targets. *Nature methods* **3**(3): 199-204.
- Birney E, Ensembl T. 2003. Ensembl: a genome infrastructure. *Cold Spring Harbor symposia on quantitative biology* **68**: 213-215.
- Boch J, Scholze H, Schornack S, Landgraf A, Hahn S, Kay S, Lahaye T, Nickstadt A, Bonas U. 2009. Breaking the code of DNA binding specificity of TAL-type III effectors. *Science* **326**(5959): 1509-1512.
- Briggs AW, Rios X, Chari R, Yang L, Zhang F, Mali P, Church GM. 2012. Iterative capped assembly: rapid and scalable synthesis of

- repeat-module DNA such as TAL effectors from individual monomers. *Nucleic acids research* **40**(15): e117.
- Brose MS, Volpe P, Feldman M, Kumar M, Rishi I, Gerrero R, Einhorn E, Herlyn M, Minna J, Nicholson A et al. 2002. BRAF and RAS mutations in human lung cancer and melanoma. *Cancer research* **62**(23): 6997-7000.
- Bultmann S, Morbitzer R, Schmidt CS, Thanisch K, Spada F, Elsaesser J, Lahaye T, Leonhardt H. 2012. Targeted transcriptional activation of silent oct4 pluripotency gene by combining designer TALEs and inhibition of epigenetic modifiers. *Nucleic acids research* **40**(12): 5368-5377.
- Bylund L, Kytola S, Lui WO, Larsson C, Weber G. 2004. Analysis of the cytogenetic stability of the human embryonal kidney cell line 293 by cytogenetic and STR profiling approaches. *Cytogenetic and genome research* **106**(1): 28-32.
- Cade L, Reyon D, Hwang WY, Tsai SQ, Patel S, Khayter C, Joung JK, Sander JD, Peterson RT, Yeh JR. 2012. Highly efficient generation of heritable zebrafish gene mutations using homo- and heterodimeric TALENs. *Nucleic acids research* **40**(16): 8001-8010.
- Carlson DF, Tan W, Lillico SG, Stverakova D, Proudfoot C, Christian M, Voytas DF, Long CR, Whitelaw CB, Fahrenkrug SC. 2012. Efficient TALEN-mediated gene knockout in livestock. *Proceedings of the National Academy of Sciences of the United States of America* **109**(43): 17382-17387.

- Cermak T, Doyle EL, Christian M, Wang L, Zhang Y, Schmidt C, Baller JA, Somia NV, Bogdanove AJ, Voytas DF. 2011. Efficient design and assembly of custom TALEN and other TAL effector-based constructs for DNA targeting. *Nucleic acids research* **39**(12): e82.
- Cho SW, Kim S, Kim JM, Kim JS. 2013. Targeted genome engineering in human cells with the Cas9 RNA-guided endonuclease. *Nature biotechnology* **31**(3): 230-232.
- Christian M, Cermak T, Doyle EL, Schmidt C, Zhang F, Hummel A, Bogdanove AJ, Voytas DF. 2010. Targeting DNA double-strand breaks with TAL effector nucleases. *Genetics* **186**(2): 757-761.
- Cong L, Ran FA, Cox D, Lin S, Barretto R, Habib N, Hsu PD, Wu X, Jiang W, Marraffini LA et al. 2013. Multiplex genome engineering using CRISPR/Cas systems. *Science* **339**(6121): 819-823.
- Cornu TI, Thibodeau-Beganny S, Guhl E, Alwin S, Eichinger M, Joung JK, Cathomen T. 2008. DNA-binding specificity is a major determinant of the activity and toxicity of zinc-finger nucleases. *Molecular therapy : the journal of the American Society of Gene Therapy* **16**(2): 352-358.
- Deng C, Capecchi MR. 1992. Reexamination of gene targeting frequency as a function of the extent of homology between the targeting vector and the target locus. *Molecular and cellular biology* **12**(8): 3365-3371.
- Deng D, Yin P, Yan C, Pan X, Gong X, Qi S, Xie T, Mahfouz M,

- Zhu JK, Yan N et al. 2012. Recognition of methylated DNA by TAL effectors. *Cell research* **22**(10): 1502-1504.
- Doyon Y, McCammon JM, Miller JC, Faraji F, Ngo C, Katibah GE, Amora R, Hocking TD, Zhang L, Rebar EJ et al. 2008. Heritable targeted gene disruption in zebrafish using designed zinc-finger nucleases. *Nature biotechnology* **26**(6): 702-708.
- Duan H, Heckman CA, Boxer LM. 2005. Histone deacetylase inhibitors down-regulate bcl-2 expression and induce apoptosis in t(14;18) lymphomas. *Molecular and cellular biology* **25**(5): 1608-1619.
- Fire A, Xu S, Montgomery MK, Kostas SA, Driver SE, Mello CC. 1998. Potent and specific genetic interference by double-stranded RNA in *Caenorhabditis elegans*. *Nature* **391**(6669): 806-811.
- Gabriel R, Lombardo A, Arens A, Miller JC, Genovese P, Kaepfel C, Nowrouzi A, Bartholomae CC, Wang J, Friedman G et al. 2011. An unbiased genome-wide analysis of zinc-finger nuclease specificity. *Nature biotechnology* **29**(9): 816-823.
- Gaeta ML, Johnson DR, Kluger MS, Pober JS. 2000. The death domain of tumor necrosis factor receptor 1 is necessary but not sufficient for Golgi retention of the receptor and mediates receptor desensitization. *Laboratory investigation; a journal of technical methods and pathology* **80**(8): 1185-1194.
- Gewurz BE, Towfic F, Mar JC, Shinnars NP, Takasaki K, Zhao B, Cahir-McFarland ED, Quackenbush J, Xavier RJ, Kieff E. 2012. Genome-wide siRNA screen for mediators of NF-kappaB activation. *Proceedings of the National Academy of Sciences of*

- the United States of America* **109**(7): 2467-2472.
- Guo J, Gaj T, Barbas CF, 3rd. 2010. Directed evolution of an enhanced and highly efficient FokI cleavage domain for zinc finger nucleases. *Journal of molecular biology* **400**(1): 96-107.
- Hockemeyer D, Wang H, Kiani S, Lai CS, Gao Q, Cassady JP, Cost GJ, Zhang L, Santiago Y, Miller JC et al. 2011. Genetic engineering of human pluripotent cells using TALE nucleases. *Nature biotechnology* **29**(8): 731-734.
- Holt N, Wang J, Kim K, Friedman G, Wang X, Taupin V, Crooks GM, Kohn DB, Gregory PD, Holmes MC, Cannon PM. 2010. Human hematopoietic stem/progenitor cells modified by zinc-finger nucleases targeted to CCR5 control HIV-1 in vivo. *Nature biotechnology* **28**(8): 839-847.
- Huang P, Xiao A, Zhou M, Zhu Z, Lin S, Zhang B. 2011. Heritable gene targeting in zebrafish using customized TALENs. *Nature biotechnology* **29**(8): 699-700.
- Hwang WY, Fu Y, Reyon D, Maeder ML, Tsai SQ, Sander JD, Peterson RT, Yeh JR, Joung JK. 2013. Efficient genome editing in zebrafish using a CRISPR-Cas system. *Nature biotechnology* **31**(3): 227-229.
- Jackson AL, Bartz SR, Schelter J, Kobayashi SV, Burchard J, Mao M, Li B, Cavet G, Linsley PS. 2003. Expression profiling reveals off-target gene regulation by RNAi. *Nature biotechnology* **21**(6): 635-637.
- Khan AA, Betel D, Miller ML, Sander C, Leslie CS, Marks DS. 2009.

- Transfection of small RNAs globally perturbs gene regulation by endogenous microRNAs. *Nature biotechnology* **27**(6): 549-555.
- Kim H, Kim MS, Wee G, Lee CI, Kim H, Kim JS. 2013a. Magnetic separation and antibiotics selection enable enrichment of cells with ZFN/TALEN-induced mutations. *PloS one* **8**(2): e56476.
- Kim H, Um E, Cho SR, Jung C, Kim H, Kim JS. 2011. Surrogate reporters for enrichment of cells with nuclease-induced mutations. *Nature methods* **8**(11): 941-943.
- Kim, H.J. 2012. Targeted mutagenesis of the human *chemokine (C-C motif) receptor 5* gene using zinc-finger nucleases and TAL effector nucleases. *Ph.D. Dissertation*, Seoul National University, Seoul, Korea.
- Kim HJ, Lee HJ, Kim H, Cho SW, Kim JS. 2009. Targeted genome editing in human cells with zinc finger nucleases constructed via modular assembly. *Genome research* **19**(7): 1279-1288.
- Kim, S., Lee, M.J., Kim, H., Kang, M. and Kim, J.S. 2011. Preassembled zinc-finger arrays for rapid construction of ZFNs. *Nat Methods*, **8**: 7.
- Kim Y, Kweon J, Kim A, Chon JK, Yoo JY, Kim HJ, Kim S, Lee C, Jeong E, Chung E et al. 2013b. A library of TAL effector nucleases spanning the human genome. *Nature biotechnology* **31**(3): 251-258.
- Kim Y, Kweon J, Kim JS. 2013c. TALENs and ZFNs are associated with different mutation signatures. *Nature methods* **10**(3): 185.
- Kim YG, Cha J, Chandrasegaran S. 1996. Hybrid restriction enzymes:

- zinc finger fusions to Fok I cleavage domain. *Proceedings of the National Academy of Sciences of the United States of America* **93**(3): 1156-1160.
- Krueger U, Bergauer T, Kaufmann B, Wolter I, Pilk S, Heider-Fabian M, Kirch S, Artz-Oppitz C, Isselhorst M, Konrad J. 2007. Insights into effective RNAi gained from large-scale siRNA validation screening. *Oligonucleotides* **17**(2): 237-250.
- Lander ES Linton LM Birren B Nusbaum C Zody MC Baldwin J Devon K Dewar K Doyle M FitzHugh W et al. 2001. Initial sequencing and analysis of the human genome. *Nature* **409**(6822): 860-921.
- Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome biology* **10**(3): R25.
- Lee HJ, Kim E, Kim JS. 2010. Targeted chromosomal deletions in human cells using zinc finger nucleases. *Genome research* **20**(1): 81-89.
- Lee HJ, Kweon J, Kim E, Kim S, Kim JS. 2012. Targeted chromosomal duplications and inversions in the human genome using zinc finger nucleases. *Genome research* **22**(3): 539-548.
- Lei Y, Guo X, Liu Y, Cao Y, Deng Y, Chen X, Cheng CH, Dawid IB, Chen Y, Zhao H. 2012. Efficient targeted gene disruption in *Xenopus* embryos using engineered transcription activator-like effector nucleases (TALENs). *Proceedings of the National Academy of Sciences of the United States of America* **109**(43):

17484-17489.

- Li T, Huang S, Zhao X, Wright DA, Carpenter S, Spalding MH, Weeks DP, Yang B. 2011. Modularly assembled designer TAL effector nucleases for targeted gene knockout and gene replacement in eukaryotes. *Nucleic acids research* **39**(14): 6315-6325.
- Li T, Liu B, Spalding MH, Weeks DP, Yang B. 2012. High-efficiency TALEN-based gene editing produces disease-resistant rice. *Nature biotechnology* **30**(5): 390-392.
- Lin X, Ruan X, Anderson MG, McDowell JA, Kroeger PE, Fesik SW, Shen Y. 2005. siRNA-mediated off-target gene silencing triggered by a 7 nt complementation. *Nucleic acids research* **33**(14): 4527-4535.
- Liu J, Li C, Yu Z, Huang P, Wu H, Wei C, Zhu N, Shen Y, Chen Y, Zhang B et al. 2012. Efficient and specific modifications of the Drosophila genome by means of an easy TALEN strategy. *Journal of genetics and genomics = Yi chuan xue bao* **39**(5): 209-215.
- Ma S, Zhang S, Wang F, Liu Y, Liu Y, Xu H, Liu C, Lin Y, Zhao P, Xia Q. 2012. Highly efficient and specific genome editing in silkworm using custom TALENs. *PloS one* **7**(9): e45035.
- Macville M, Schrock E, Padilla-Nash H, Keck C, Ghadimi BM, Zimonjic D, Popescu N, Ried T. 1999. Comprehensive and definitive molecular cytogenetic characterization of HeLa cells by spectral karyotyping. *Cancer research* **59**(1): 141-150.

- Maeder ML, Angstman JF, Richardson ME, Linder SJ, Cascio VM, Tsai SQ, Ho QH, Sander JD, Reyon D, Bernstein BE et al. 2013a. Targeted DNA demethylation and activation of endogenous genes using programmable TALE-TET1 fusion proteins. *Nature biotechnology*.
- Maeder ML, Linder SJ, Reyon D, Angstman JF, Fu Y, Sander JD, Joung JK. 2013b. Robust, synergistic regulation of human gene expression using TALE activators. *Nature methods* **10**(3): 243-245.
- Maeder ML, Thibodeau-Beganny S, Osiak A, Wright DA, Anthony RM, Eichinger M, Jiang T, Foley JE, Winfrey RJ, Townsend JA et al. 2008. Rapid "open-source" engineering of customized zinc-finger nucleases for highly efficient gene modification. *Molecular cell* **31**(2): 294-301.
- Mali P, Yang L, Esvelt KM, Aach J, Guell M, DiCarlo JE, Norville JE, Church GM. 2013. RNA-guided human genome engineering via Cas9. *Science* **339**(6121): 823-826.
- Maresca M, Lin VG, Guo N, Yang Y. 2013. Obligate ligation-gated recombination (ObLiGaRe): custom-designed nuclease-mediated targeted integration through nonhomologous end joining. *Genome research* **23**(3): 539-546.
- Mendenhall EM, Williamson KE, Reyon D, Zou JY, Ram O, Joung JK, Bernstein BE. 2013. Locus-specific editing of histone modifications at endogenous enhancers. *Nature biotechnology*.
- Mercer AC, Gaj T, Fuller RP, Barbas CF, 3rd. 2012. Chimeric TALE

- recombinases with programmable DNA sequence specificity. *Nucleic acids research* **40**(21): 11163-11172.
- Miller JC, Tan S, Qiao G, Barlow KA, Wang J, Xia DF, Meng X, Paschon DE, Leung E, Hinkley SJ et al. 2011. A TALE nuclease architecture for efficient genome editing. *Nature biotechnology* **29**(2): 143-148.
- Miyanari Y, Ziegler-Birling C, Torres-Padilla ME. 2013. Live visualization of chromatin dynamics with fluorescent TALEs. *Nature structural & molecular biology* **20**(11): 1321-1324.
- Morbitzer R, Elsaesser J, Hausner J, Lahaye T. 2011. Assembly of custom TALE-type DNA binding domains by modular cloning. *Nucleic acids research* **39**(13): 5790-5799.
- Moscou MJ, Bogdanove AJ. 2009. A simple cipher governs DNA recognition by TAL effectors. *Science* **326**(5959): 1501.
- Mussolino C, Morbitzer R, Lutge F, Dannemann N, Lahaye T, Cathomen T. 2011. A novel TALE nuclease scaffold enables high genome editing activity in combination with low toxicity. *Nucleic acids research* **39**(21): 9283-9293.
- Paez JG, Janne PA, Lee JC, Tracy S, Greulich H, Gabriel S, Herman P, Kaye FJ, Lindeman N, Boggon TJ et al. 2004. EGFR mutations in lung cancer: correlation with clinical response to gefitinib therapy. *Science* **304**(5676): 1497-1500.
- Pattanayak V, Ramirez CL, Joung JK, Liu DR. 2011. Revealing off-target cleavage specificities of zinc-finger nucleases by in vitro selection. *Nature methods* **8**(9): 765-770.

- Perez-Pinera P, Ousterout DG, Brunger JM, Farin AM, Glass KA, Guilak F, Crawford GE, Hartemink AJ, Gersbach CA. 2013. Synergistic and tunable human gene activation by combinations of synthetic transcription factors. *Nature methods* **10**(3): 239-242.
- Perkins ND. 2007. Integrating cell-signalling pathways with NF-kappaB and IKK function. *Nature reviews Molecular cell biology* **8**(1): 49-62.
- Pruitt KD, Tatusova T, Brown GR, Maglott DR. 2012. NCBI Reference Sequences (RefSeq): current status, new features and genome annotation policy. *Nucleic acids research* **40**(Database issue): D130-135.
- Reyon D, Tsai SQ, Khayter C, Foden JA, Sander JD, Joung JK. 2012. FLASH assembly of TALENs for high-throughput genome editing. *Nature biotechnology* **30**(5): 460-465.
- Sander JD, Cade L, Khayter C, Reyon D, Peterson RT, Joung JK, Yeh JR. 2011. Targeted gene disruption in somatic zebrafish cells using engineered TALENs. *Nature biotechnology* **29**(8): 697-698.
- Seal RL, Gordon SM, Lush MJ, Wright MW, Bruford EA. 2011. genenames.org: the HGNC resources in 2011. *Nucleic acids research* **39**(Database issue): D514-519.
- Sigoillot FD, King RW. 2011. Vigilance and validation: Keys to success in RNAi screening. *ACS chemical biology* **6**(1): 47-60.
- Sledz CA, Holko M, de Veer MJ, Silverman RH, Williams BR. 2003. Activation of the interferon system by short-interfering RNAs. *Nature cell biology* **5**(9): 834-839.

- Smithies O, Gregg RG, Boggs SS, Koralewski MA, Kucherlapati RS. 1985. Insertion of DNA sequences into the human chromosomal beta-globin locus by homologous recombination. *Nature* **317**(6034): 230-234.
- Sung YH, Baek IJ, Kim DH, Jeon J, Lee J, Lee K, Jeong D, Kim JS, Lee HW. 2013. Knockout mice created by TALEN-mediated gene targeting. *Nature biotechnology* **31**(1): 23-24.
- Tesson L, Usal C, Menoret S, Leung E, Niles BJ, Remy S, Santiago Y, Vincent AI, Meng X, Zhang L et al. 2011. Knockout rats generated by embryo microinjection of TALENs. *Nature biotechnology* **29**(8): 695-696.
- Urnov FD, Miller JC, Lee YL, Beausejour CM, Rock JM, Augustus S, Jamieson AC, Porteus MH, Gregory PD, Holmes MC. 2005. Highly efficient endogenous human gene correction using designed zinc-finger nucleases. *Nature* **435**(7042): 646-651.
- Valton J, Dupuy A, Daboussi F, Thomas S, Marechal A, Macmaster R, Melliand K, Juillerat A, Duchateau P. 2012. Overcoming transcription activator-like effector (TALE) DNA binding domain sensitivity to cytosine methylation. *The Journal of biological chemistry* **287**(46): 38427-38432.
- Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, Smith HO, Yandell M, Evans CA, Holt RA et al. 2001. The sequence of the human genome. *Science* **291**(5507): 1304-1351.
- Volcic M, Karl S, Baumann B, Salles D, Daniel P, Fulda S, Wiesmuller L. 2012. NF-kappaB regulates DNA double-strand

- break repair in conjunction with BRCA1-CtIP complexes. *Nucleic acids research* **40**(1): 181-195.
- Watanabe T, Ochiai H, Sakuma T, Horch HW, Hamaguchi N, Nakamura T, Bando T, Ohuchi H, Yamamoto T, Noji S et al. 2012. Non-transgenic genome modifications in a hemimetabolous insect using zinc-finger and TAL effector nucleases. *Nature communications* **3**: 1017.
- Weber E, Gruetzner R, Werner S, Engler C, Marillonnet S. 2011. Assembly of designer TAL effectors by Golden Gate cloning. *PloS one* **6**(5): e19722.
- Wood AJ, Lo TW, Zeitler B, Pickle CS, Ralston EJ, Lee AH, Amora R, Miller JC, Leung E, Meng X et al. 2011. Targeted genome editing across species using ZFNs and TALENs. *Science* **333**(6040): 307.
- Zhang F, Cong L, Lodato S, Kosuri S, Church GM, Arlotta P. 2011. Efficient construction of sequence-specific TAL effectors for modulating mammalian transcription. *Nature biotechnology* **29**(2): 149-153.

국문 초록

TALEN은 최근에 새롭게 개발된 유전자 가위로 DNA 의 염기를 하나씩 특이적으로 인식하는 TALE을 이용한다. 특히 TALEN은 기존의 유전자 가위보다 쉽게 제작할 수 있고 유전자 변이 효율도 상당히 개선되어 다양한 동·식물에 응용되고 있다. 이번 연구에서는 기존의 TALEN의 구조를 개선해 게놈 상에 원치 않는 곳에 변이를 일으킬 가능성을 줄였다. 그리고 컴퓨터 프로그래밍을 통하여 인간 전체의 유전자에 대해, 게놈 상에 유사한 염기서열이 거의 없는 TALEN 적중 염기서열을 선별했다. 선별된 TALEN 적중 서열에 대해 TALEN을 쉽게 만들 수 있는 고속대량 합성법을 개발하였다. 이를 통해 100여개 이상의 유전자를 적중하는 TALEN을 만들었고 실제 변이를 유도한 결과 거의 모든 TALEN이 높은 효율로 유전자 변이를 일으켰다. 변이를 유도하지 못한 두 TALEN의 적중 서열을 분석한 결과, 시토신이 메틸화되어 있었으며 이로 인해 TALEN이 작동하지 못했음을 알 수 있었다. 또한 이 TALEN을 이용하여 유전자 변이가 유도된 세포주를 건립하였으며 실제 유전자의 기능이 없어졌음을 확인했다. 이러한 유전자 녹아웃은 기존의 siRNA를 이용하는 방법보다 매우 효율적으로 완벽히 유전자의 기능을 상실시키는 확인했다. 나아가 인간뿐만 아니라, 다양한 실험동물들에 대해 TALEN 적중 서열을 미리 선별해 놓음으로써 신약개발 등의 다양한 연구에 널리 활용될 수 있을 것으로 생각한다.

학 번: 2008-22719