Engineered CRISPR-Cas9 nucleases with altered PAM specificities

Benjamin P. Kleinstiver^{1,2,3}, Michelle S. Prew^{1,2}, Shengdar Q. Tsai^{1,2,3}, Ved V. Topkar^{1,2}, Nhu T. Nguyen^{1,2}, Zongli Zheng^{1,3,4}, Andrew P. W. Gonzales^{5,6,7}, Zhuyun Li⁵, Randall T. Peterson^{5,6,7}, Jing-Ruey Joanna Yeh^{5,8}, Martin J. Aryee^{1,3,9} & J. Keith Joung^{1,2,3}

Although CRISPR-Cas9 nucleases are widely used for genome editing^{1,2}, the range of sequences that Cas9 can recognize is constrained by the need for a specific protospacer adjacent motif (PAM)³⁻⁶. As a result, it can often be difficult to target doublestranded breaks (DSBs) with the precision that is necessary for various genome-editing applications. The ability to engineer Cas9 derivatives with purposefully altered PAM specificities would address this limitation. Here we show that the commonly used Streptococcus pyogenes Cas9 (SpCas9) can be modified to recognize alternative PAM sequences using structural information, bacterial selection-based directed evolution, and combinatorial design. These altered PAM specificity variants enable robust editing of endogenous gene sites in zebrafish and human cells not currently targetable by wild-type SpCas9, and their genome-wide specificities are comparable to wild-type SpCas9 as judged by GUIDE-seq analysis⁷. In addition, we identify and characterize another SpCas9 variant that exhibits improved specificity in human cells, possessing better discrimination against off-target sites with non-canonical NAG and NGA PAMs and/or mismatched spacers. We also find that two smaller-size Cas9 orthologues, Streptococcus thermophilus Cas9 (St1Cas9) and Staphylococcus aureus Cas9 (SaCas9), function efficiently in the bacterial selection systems and in human cells, suggesting that our engineering strategies could be extended to Cas9s from other species. Our findings provide broadly useful SpCas9 variants and, more importantly, establish the feasibility of engineering a wide range of Cas9s with altered and improved PAM specificities.

CRISPR-Cas⁹ nucleases enable efficient genome editing in a wide variety of organisms and cell types^{1,2}. Target site recognition by Cas⁹ is programmed by a chimaeric single guide RNA (sgRNA) that encodes a sequence complementary to a target protospacer⁵, but also requires recognition of a short neighbouring PAM³⁻⁶. SpCas⁹, the most robust and widely used Cas⁹ to date, primarily recognizes NGG PAMs and is consequently restricted to sites that contain this motif^{5,8}. It can therefore be challenging to implement genome editing applications that require precision, such as homology-directed repair, which is most efficient when DSBs are placed within 10–20 base pairs of a desired alteration^{9–11}; the introduction of variable-length insertion or deletion (indel) mutations into small size genetic elements such as microRNAs, splice sites, short open reading frames, or transcription factor binding sites by non-homologous end-joining; and allele-specific editing, where PAM recognition might be exploited to differentiate alleles.

One potential solution to address targeting range limitations would be to engineer Cas9 variants with novel PAM specificities. A previous attempt to alter SpCas9 PAM specificity mutated R1333 and R1335 residues that contact the guanine nucleotides at the second and third PAM positions; however, the R1333Q/R1335Q variant failed to cleave a site harbouring the expected NAA PAM *in vitro*¹². Using a human U2OS-cell-based enhanced green fluorescent protein (EGFP) reporter gene disruption assay in which nuclease-induced indels lead to loss of fluorescence^{13,14}, we confirmed that an R1333Q/R1335Q SpCas9 variant failed to efficiently cleave target sites with NAA PAMs (Fig. 1a). Additionally, we found that single R1333Q and R1335Q variants each failed to efficiently cleave target sites containing the expected NAG and NGA PAMs, respectively (Fig. 1a), suggesting that re-engineering PAM specificity might require additional mutations.

To identify such mutations, we adapted a bacterial selection system (hereafter referred to as the positive selection) previously used to study properties of homing endonucleases^{15,16}. In our adaptation of this system, survival is enabled by Cas9-mediated cleavage of a selection plasmid encoding an inducible toxic gene (Fig. 1b, Extended Data Fig. 1a). We mutagenized the PAM-interacting domains of wild-type and R1335Q SpCas9 and performed selections against an NGA PAM target site (Extended Data Fig. 1b, Methods). Sequences of surviving clones from both libraries revealed the most frequent substitutions were D1135V/Y/N/E, R1335Q, and T1337R (Extended Data Fig. 2a). After testing all combinations of these mutations using the human cell-based EGFP disruption assay, two variants were chosen for further characterization because they possessed the greatest discrimination between NGA and NGG PAMs: D1135V/R1335Q/T1337R and D1135E/R1335Q/T1337R (hereafter referred to as the VQR and EQR variants, respectively) (Fig. 1c).

To define the global PAM specificity profiles of these SpCas9 variants, we used a bacterial-based negative selection system (Fig. 1d, Extended Data Fig. 3a) similar to other methods previously used to identify PAM preferences of Cas9 (refs 8, 17). In this site-depletion assay, a library of plasmids bearing 6 randomized base pairs adjacent to a protospacer is tested for cleavage by Cas9 in Escherichia coli (Extended Data Fig. 3b). Plasmids with PAM sequences refractory to Cas9 enable cell survival due to the presence of an antibiotic resistance gene, whereas plasmids bearing targetable PAMs are depleted from the library (Fig. 1d, Extended Data Fig. 3b). Sequencing the uncleaved population of plasmids enables the calculation of a postselection PAM depletion value (PPDV), an estimate of Cas9 activity against those PAMs (post-selection frequency relative to the preselection frequency). Site-depletion data obtained with catalytically inactive Cas9 (dCas9) on two randomized PAM libraries (each with a different protospacer) enabled us to define what represents a statistically significant change in PPDV for any given PAM or group of PAMs (Extended Data Fig. 3c, d), and PPDVs observed for wild-type SpCas9 recapitulated its previously described profile of targetable PAMs⁸ (Fig. 1e).

Using the site-depletion assay, we obtained PAM specificity profiles for the VQR and EQR variants. The VQR variant strongly depleted sites bearing NGAN and NGCG PAMs, while the EQR variant seemed more specific for an NGAG PAM (Fig. 1f). Human cell

¹Molecular Pathology Unit & Center for Cancer Research, Massachusetts General Hospital, Charlestown, Massachusetts 02129, USA. ²Center for Computational and Integrative Biology, Massachusetts General Hospital, Charlestown, Massachusetts 02129, USA. ³Department of Pathology, Harvard Medical School, Boston, Massachusetts 02115, USA. ⁴Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm SE-171 77, Sweden. ⁵Cardiovascular Research Center, Massachusetts General Hospital, Charlestown, Massachusetts 02129, USA. ⁶Department of Systems Biology, Harvard Medical School, Boston, Massachusetts 02115, USA. ⁷Broad Institute, Cambridge, Massachusetts 02124, USA. ⁸Department of Medicine, Harvard Medical School, Boston, Massachusetts 02115, USA. ⁹Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, Massachusetts 02115, USA.



Figure 1 Evolution and characterization of SpCas9 variants with altered PAM specificities. a, Activity of wild-type and mutant SpCas9s assessed via U2OS human cell-based EGFP disruption. Frequencies were quantified by flow cytometry; error bars represent s.e.m., n = 3; mean level of background EGFP loss represented by the dashed red line for this and subsequent panels (c, g, h and j). b, Schematic of the positive selection assay (see also Extended Data Fig. 1). c, Combinatorial assembly and testing of mutations obtained from the positive selection for SpCas9 variants that can cleave a target site containing an NGA PAM, using the human cell EGFP disruption assay. d, Schematic of the negative selection assay, adapted to profile Cas9 PAM specificity by generating a library of plasmids that contain a randomized sequence adjacent to the 3' end of the protospacer (see also Extended Data Fig. 3a, b). e, Scatterplot of the post-selection PAM depletion values (PPDVs) of

EGFP disruption experiments paralleled these results, with the VQR variant robustly cleaving sites bearing NGAN PAMs (with relative efficiencies NGAG > NGAT = NGAA > NGAC), and also sites bearing NGNG PAMs with generally lower efficiencies (Fig. 1g). Similarly, the EQR variant preferred NGAG to the other NGAN and NGNG PAMs in human cells, again at lower activities than with the VQR variant (Fig. 1g). The activities of the VQR and EQR variants in human cells therefore recapitulated what was observed with the bacterial sitedepletion assay and suggested that PPDVs of 0.2 (fivefold depletion) provide a reasonable predictive threshold for activity in human cells (Extended Data Fig. 4).

We next sought to extend the generalizability of our engineering strategy by identifying SpCas9 variants capable of recognizing an NGC PAM. Selections using libraries bearing pre-existing R1335E/T1337R and R1335T/T1337R substitutions (Methods) yielded surviving colonies harbouring a variety of additional mutations (Extended Data Fig. 2b). Testing all possible combinations of the most common mutations using the EGFP disruption assay established that the quadruple mutant VRER variant (D1135V/G1218R/R1335E/T1337R) displayed

wild-type SpCas9 with two randomized PAM libraries (each with a different protospacer). PAMs are plotted by their second/third/fourth positions. The red dashed line indicates statistically significant depletion (obtained from a dCas9 control experiment, see Extended Data Fig. 3c), and the grey dashed line represents fivefold depletion (PPDV of 0.2). **f**, PPDV scatterplots for the VQR and EQR variants. **g**, EGFP disruption frequencies for wild-type, VQR, and EQR spCas9 on sites with NGAN and NGNG PAMs. **h**, Combinatorial assembly and testing of mutations obtained from the positive selection for SpCas9 variants that can cleave a target site containing an NGC PAM, using the human cell EGFP disruption assay. **i**, PPDV scatterplot for the VRER variant. **j**, EGFP disruption frequencies for wild-type and VRER SpCas9 on sites with NGCN and NGNG PAMs.

the highest activity on an NGCG PAM and minimal activity on an NGG PAM (Fig. 1h). Analysis of the VRER variant using the sitedepletion assay revealed it to be highly specific for NGCG PAMs (Fig. 1i). Consistent with this result, EGFP disruption assays revealed efficient cleavage of sites with NGCG PAMs, and inconsistent or low activity against NGCH and NGDG PAMs (Fig. 1j). Notably, the mutations critical for altering the specificity of SpCas9 are spatially oriented near the PAM (Extended Data Fig. 5a), and the nature and effect of the mutations imply that they are most likely gain of function (Extended Data Fig. 5b). For example, the T1337R mutation seems to confer a preference for a fourth PAM base, especially in the case of the VRER variant.

To demonstrate directly that the SpCas9 variants broaden the targeting range of SpCas9, we tested their activities against endogenous genes in zebrafish embryos and human cells. In zebrafish embryos, the VQR variant efficiently modified sites bearing NGAG PAMs (range of 20 to 43%, Fig. 2a) with the indels originating at the predicted cleavage sites (Extended Data Fig. 6). In human cells, the VQR variant robustly modified endogenous sites that harboured NGA PAMs (again, with a



Figure 2 | SpCas9 PAM variants robustly modify endogenous sites in zebrafish embryos and human cells. a, Mutagenesis frequencies in zebrafish embryos induced by wild-type or VQR SpCas9 at endogenous gene sites bearing NGAG PAMs. Mutation frequencies were determined using the T7E1 assay; ND, not detectable by T7E1; error bars represent s.e.m., n = 5 to 9 embryos. b, Endogenous human gene disruption activity of the VQR variant quantified by T7E1 assay. Error bars represent s.e.m., n = 3. c, Endogenous human gene disruption activity of wild-type SpCas9 against NGA PAM sites quantified by T7E1 assay, where VQR data are re-presented from panel b for

ease of comparison. Error bars represent s.e.m., n = 3. **d**, Mutation frequencies of wild-type, VRER, and VQR SpCas9 at endogenous human gene sites containing NGCG PAMs quantified by T7E1 assay; error bars represent s.e.m., n = 3. **e**, Representation of the number of sites in the human genome with 20-nucleotide spacers potentially targetable by wild-type, VQR, and VRER SpCas9. The 5'-G is included for expression from a U6 promoter. **f**, Number of off-target cleavage sites identified by GUIDE-seq for the VQR and VRER variants using sgRNAs from panels **b** and **d**.

preference for NGAG> NGAT = NGAA, range of 6 to 53%) (Fig. 2b, Extended Data Fig. 7a). Importantly, wild-type SpCas9 was unable to robustly alter NGA PAM sites in zebrafish and human cells (Fig. 2a, c), yet was able to efficiently modify neighbouring sites bearing NGG PAMs in human cells (Extended Data Fig. 7b). When examining VRER variant activity at endogenous human sites with NGCG PAMs, we also observed robust disruption frequencies (range of 5 to 36%) (Fig. 2d). Consistent with the site-depletion data (Fig. 1e, f), the VQR variant also altered NGCG PAM sites while wild-type SpCas9 was unable to do so (Fig. 2d). Taken together, these results demonstrate that the VQR and VRER variants enable modification of previously inaccessible sites in zebrafish embryos and human cells, and computational analysis of the reference human genome reveals that they double the targeting potential of SpCas9 (Fig. 2e). To identify target sites for the engineered variants, we have developed a web-based tool called CasBLASTR (http://www.CasBLASTR.org).

To determine the genome-wide specificity of the VQR and VRER SpCas9 nucleases, we used the recently described GUIDE-seq method⁷ to profile off-target cleavage events in human cells. The total number of detectable off-target DSBs induced by the SpCas9 variants in human cells (Fig. 2f) are comparable to (or, in the case of the VRER variant, perhaps less than) what has been previously observed with wild-type SpCas9 (ref. 7). The off-target sites observed generally possess the expected PAM sequences predicted by our site-depletion experiments (compare Figs 1f, i to Extended Data Fig. 8), and the mismatches observed in the off-target sites of the variants are similar to the profiles previously observed with wild-type SpCas9 for sgRNAs targeted to non-repetitive sequences⁷. The stringent genome-wide specificity observed with the VRER variant might result from its extension of the PAM by 1 base pair, and perhaps from the relative depletion of NGCG PAMs in the human genome (Fig. 2e)¹⁸.

Previous studies have shown that imperfect PAM recognition by SpCas9 can lead to recognition of non-canonical PAMs^{7,8,19–21}. While

engineering the VQR variant, we noticed that a D1135E mutant seemed to discriminate between NGG and NGA PAMs better than wild-type SpCas9 (Fig. 1c). Using the site-depletion assay to assess the D1135E variant, we observed a decrease in activity against noncanonical NAG, NGA, and NNGG PAMs relative to wild-type SpCas9, with this effect being more prominent for one protospacer (Fig. 3a). Improved PAM specificity was also observed in human cell EGFP disruption assays, where NAG and NGA PAM sites were less efficiently cleaved by D1135E compared to wild-type SpCas9 (Fig. 3b, mean fold decrease in activity of 1.94). Importantly, wild-type and D1135E SpCas9 had comparable activities against canonical NGG PAM sites when targeted to the EGFP reporter or endogenous human gene sites (mean fold decrease in activity of 1.04) (Fig. 3b and Extended Data Fig. 9a, respectively). It is unlikely that the enhanced specificity of the D1135E variant is the result of protein destabilization, because titration experiments revealed no substantial differences in activity compared with wild-type SpCas9 (Extended Data Fig. 9b).

To more directly assess the effect of D1135E on off-target effects, we examined the mutation rates induced by wild-type and D1135E SpCas9 on 25 previously known off-target sites of three sgRNAs^{7,14,19}. Deep-sequencing revealed that D1135E improved specificity for 19 of the 22 off-target sites with mutation frequencies above background indel rates, when compared to the relative mutation frequencies observed at the on-target sites (Fig. 3c, Extended Data Fig. 9c). Interestingly, the gains in specificity with D1135E are not restricted to sites with non-canonical PAMs. To more thoroughly assess the improvements in specificity associated with the D1135E variant, we performed GUIDE-seq using three different sgRNAs and observed a generalized improvement in genome-wide specificity relative to wild-type SpCas9 (Fig. 3d, Extended Data Fig. 9d–f). Collectively, these results show that the D1135E substitution increases the specificity of SpCas9.



Figure 3 A D1135E mutation improves the PAM recognition and spacer specificity of SpCas9. a, PPDV scatterplots for wild-type and D1135E SpCas9 for the two randomized PAM libraries. PAMs are plotted by their second/third/ fourth positions, and wild-type data are the same as shown in Fig. 1e for ease of comparison. The red dashed line indicates PAMs that are statistically significantly depleted (see Extended Data Fig. 3c), and the grey dashed line indicates fivefold depletion (PPDV of 0.2). **b**, EGFP disruption activities of wild-type and D1135E SpCas9 on sites that contain canonical and non-canonical PAMs in human cells. Disruption frequencies were quantified by

flow cytometry; mean background level of EGFP loss represented by the dashed red line; error bars represent s.e.m., n = 3; fold change in activity is shown. c, Summary of targeted deep-sequencing data demonstrating specificity gains at off-target sites when using D1135E (see also Extended Data Fig. 9c). d, Summary of GUIDE-seq detected changes in specificity between wild-type and D1135E at off-target sites (see also Extended Data Fig. 9f). Estimated fold gain in specificity at sites without read counts for D1135E are not plotted (see Extended Data Fig. 9f).

The many Cas9 orthologues from other bacteria make attractive candidates for characterizing and engineering Cas9s with novel PAM specificities^{22,23}. To explore this, we determined whether two smaller-size orthologues, *Streptococcus thermophilus* Cas9 from the CRISPR1 locus (St1Cas9)^{24,25} and *Staphylococcus aureus* (SaCas9)²³ could function in the bacterial selection assays. Although the PAM

of St1Cas9 has previously been characterized as NNAGAA^{17,22,24,25}, our attempts to bioinformatically derive the SaCas9 PAM using a previously described approach²² failed to yield a consensus sequence. Therefore, we used the site-depletion assay to determine the PAM for SaCas9 and, as a positive control, St1Cas9. For St1Cas9, we identified two novel PAMs in addition to six PAMs that had been



Figure 4 Characterization of St1Cas9 and SaCas9 in bacteria and human cells. a, b, PPDV scatterplots for St1Cas9 (a) and SaCas9 (b), with PAMs plotted by their third/fourth/fifth/sixth positions. The red dashed line indicates PAMs that are significantly depleted (Extended Data Fig. 3c), and the grey dashed line represents fivefold depletion (PPDV of 0.2); α , PAM previously predicted by a bioinformatic approach²⁵; β , PAMs previously identified under stringent experimental conditions¹⁷; *, novel PAMs discovered in this study;

 γ , PAMs previously identified under moderate experimental conditions¹⁷. **c**, Survival percentages of St1Cas9 and SaCas9 in the bacterial positive selection when challenged with selection plasmids that harbour different spacer sequences and PAMs. NS, no survival. **d**, **e**, Mutation frequencies of St1Cas9 (**d**) and SaCas9 (**e**) quantified by T7E1 assay at sites in four endogenous human genes. Error bars represent s.e.m., n = 3; ND, not detectable by T7E1; nt, nucleotide. previously described^{17,22,25} (Fig. 4a, Extended Data Fig. 10a, b). For SaCas9, only three PAMs were depleted more than fivefold in all experiments (NNGGGT, NNGAAT, NNGAGT, Fig. 4b), although additional PAMs were targetable when using the second protospacer library (Extended Data Fig. 10c, d). These results are consistent with a recent definition of SaCas9 PAM specificity²³. We also found that St1Cas9 and SaCas9 can function efficiently in the bacterial positive selection system (Fig. 4c), suggesting that their PAM specificities could potentially be modified by mutagenesis and selection.

Because not all Cas9 orthologues function efficiently outside of their native context^{17,23}, we tested whether St1Cas9 and SaCas9 can modify sites in human cells. St1Cas9 has been previously shown to function as a nuclease in human cells but only on four sites^{17,23,26}, and a recently published manuscript assessed SaCas9 activity²³. In EGFP disruption experiments, St1Cas9 displayed high activity at three of five target sites and SaCas9 efficiently targeted eight sites (Extended Data Fig. 10e). No obvious correlation between activity and length of spacer was observed (Extended Data Fig. 10e, f). When examining activity on endogenous loci, St1Cas9 efficiently targeted 7 out of 11 sites (1 to 25% disruption; Fig. 4d), SaCas9 displayed more robust activity at 16 sites (1% to 37%; Fig. 4e), and again no distinct spacer length requirement was observed (Extended Data Fig. 10g). Collectively, these results demonstrate that St1Cas9 and SaCas9 function in human cells, making them attractive candidates for engineering additional variants with novel PAM specificities.

The VQR and VRER variants engineered in this study enhance the opportunities to utilize the CRISPR-Cas9 platform to practice efficient homology-directed repair, to generate non-homologous end-joiningmediated indels in small genetic elements, and to exploit the requirement for a PAM to distinguish between different alleles in the same cell. Importantly, the VQR, VRER, and D1135E variants all have similar (or better) genome-wide specificities compared to wild-type SpCas9. These variants can be rapidly incorporated into existing and widely used SpCas9 vectors by simple site-directed mutagenesis, and we expect that the variants should also work with other previously described improvements to the SpCas9 platform (for example, truncated sgRNAs^{7,27}, SpCas9 nickases^{20,28}, or dimeric FokI-dCas9 fusions^{29,30}). Collectively, our results establish engineering PAM recognition and characterization of additional Cas9 orthologues (as previously described)17,22,23 as complementary approaches to provide researchers with an expanded repertoire of genome-editing reagents, while also demonstrating the feasibility of engineering Cas9 nucleases with useful new properties.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 1 March 2015; accepted 1 June 2015. Published online 22 June 2015.

- Sander, J. D. & Joung, J. K. CRISPR-Cas systems for editing, regulating and targeting genomes. *Nature Biotechnol.* **32**, 347–355 (2014).
- Doudna, J. A. & Charpentier, E. Genome editing. The new frontier of genome engineering with CRISPR-Cas9. Science 346, 1258096 (2014).
- Mojica, F. J., Diez-Villasenor, C., Garcia-Martinez, J. & Almendros, C. Short motif sequences determine the targets of the prokaryotic CRISPR defence system. *Microbiology* 155, 733–740 (2009).
- Shah, S. A., Erdmann, S., Mojica, F. J. & Garrett, R. A. Protospacer recognition motifs: mixed identities and functional diversity. *RNA Biol.* 10, 891–899 (2013).
- Jinek, M. et al. A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. Science 337, 816–821 (2012).
- Sternberg, S. H., Redding, S., Jinek, M., Greene, E. C. & Doudna, J. A. DNA interrogation by the CRISPR RNA-guided endonuclease Cas9. *Nature* 507, 62–67 (2014).
- Tsai, S. Q. et al. GUIDE-seq enables genome-wide profiling of off-target cleavage by CRISPR-Cas nucleases. Nature Biotechnol. 33, 187–197 (2015).

- Jiang, W., Bikard, D., Cox, D., Zhang, F. & Marraffini, L. A. RNA-guided editing of bacterial genomes using CRISPR-Cas systems. *Nature Biotechnol.* **31**, 233–239 (2013).
- Yang, L. et al. Optimization of scarless human stem cell genome editing. Nucleic Acids Res. 41, 9049–9061 (2013).
- Elliott, B., Richardson, C., Winderbaum, J., Nickoloff, J. A. & Jasin, M. Gene conversion tracts from double-strand break repair in mammalian cells. *Mol. Cell. Biol.* 18, 93–101 (1998).
- Findlay, G. M., Boyle, E. A., Hause, R. J., Klein, J. C. & Shendure, J. Saturation editing of genomic regions by multiplex homology-directed repair. *Nature* **513**, 120–123 (2014).
- Anders, C., Niewoehner, O., Duerst, A. & Jinek, M. Structural basis of PAMdependent target DNA recognition by the Cas9 endonuclease. *Nature* 513, 569–573 (2014).
- Reyon, D. et al. FLASH assembly of TALENs for high-throughput genome editing. Nature Biotechnol. 30, 460–465 (2012).
- Fu, Y. et al. High-frequency off-target mutagenesis induced by CRISPR-Cas nucleases in human cells. Nature Biotechnol. 31, 822–826 (2013).
- Chen, Z. & Zhao, H. A highly sensitive selection method for directed evolution of homing endonucleases. *Nucleic Acids Res.* 33, e154 (2005).
- Doyon, J. B., Pattanayak, V., Meyer, C. B. & Liu, D. R. Directed evolution and substrate specificity profile of homing endonuclease I-Scel. J. Am. Chem. Soc. 128, 2477–2484 (2006).
- Esvelt, K. M. et al. Orthogonal Cas9 proteins for RNA-guided gene regulation and editing. Nature Methods 10, 1116–1121 (2013).
- Lander, E. S. et al. Initial sequencing and analysis of the human genome. Nature 409, 860–921 (2001).
- Hsu, P. D. et al. DNA targeting specificity of RNA-guided Cas9 nucleases. Nature Biotechnol. 31, 827–832 (2013).
- Mali, P. et al. CAS9 transcriptional activators for target specificity screening and paired nickases for cooperative genome engineering. *Nature Biotechnol.* **31**, 833–838 (2013).
- Zhang, Y. et al. Comparison of non-canonical PAMs for CRISPR/Cas9-mediated DNA cleavage in human cells. Sci. Rep. 4, 5405 (2014).
- Fonfara, I. *et al.* Phylogeny of Cas9 determines functional exchangeability of dual-RNA and Cas9 among orthologous type II CRISPR-Cas systems. *Nucleic Acids Res.* 42, 2577–2590 (2014).
- Ran, F. A. et al. In vivo genome editing using Staphylococcus aureus Cas9. Nature, (2015).
- Deveau, H. et al. Phage response to CRISPR-encoded resistance in Streptococcus thermophilus. J. Bacteriol. 190, 1390–1400 (2008).
- Horvath, P. et al. Diversity, activity, and evolution of CRISPR loci in Streptococcus thermophilus. J. Bacteriol. 190, 1401–1412 (2008).
- Cong, L. et al. Multiplex genome engineering using CRISPR/Cas systems. Science 339, 819–823 (2013).
- Fu, Y., Sander, J. D., Reyon, D., Cascio, V. M. & Joung, J. K. Improving CRISPR-Cas nuclease specificity using truncated guide RNAs. *Nature Biotechnol.* 32, 279–284 (2014).
- Ran, F. A. et al. Double nicking by RNA-guided CRISPR Cas9 for enhanced genome editing specificity. Cell 154, 1380–1389 (2013).
- Guilinger, J. P., Thompson, D. B. & Liu, D. R. Fusion of catalytically inactive Cas9 to Fokl nuclease improves the specificity of genome modification. *Nature Biotechnol.* 32, 577–582 (2014).
- Tsai, S. Q. et al. Dimeric CRISPR RNA-guided Fokl nucleases for highly specific genome editing. Nature Biotechnol. 32, 569–576 (2014).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank D. Edgell for providing the bacterial strain and plasmids related to the bacterial selection; J. Angstman and V. Pattanayak for discussion and comments on the manuscript. This work was supported by a National Institutes of Health (NIH) Director's Pioneer Award (DP1 GM105378) and NIH R01 GM088040 to J.K.J. NIH R01 GM088040 to J.K.J. and R.T.P., The Jim and Ann Orr Research Scholar Award (to J.K.J.), and a National Sciences and Engineering Research Council of Canada Postdoctoral Fellowship (to B.P.K.).

Author Contributions B.P.K., M.S.P., S.Q.T. and N.T.N. performed all bacterial and human cell-based experiments. A.P.W.G. and Z.L. performed all zebrafish experiments. S.Q.T., V.T., Z.Z. and M.J.A. analysed the site-depletion, targeted deep-sequencing, and GUIDE-seq data. B.P.K., R.T.P., J.-R.J.Y. and J.K.J. directed the research and interpreted experiments. B.P.K. and J.K.J. wrote the manuscript with input from all the authors.

Author Information All new reagents described in this work have been deposited with the non-profit plasmid distribution service Addgene (http://www.addgene.org/ crispr-cas). A web-tool to design sgRNA sites for the engineered variants and orthogonal Cas9 nucleases described in this study can be found at http://www.CasBLASTR.org. The sequences generated in this study have been deposited in the Sequences Read Archive under accession number SRP058629. Reprints and permissions information is available at www.nature.com/reprints. Readers are welcome to comment on the online version of the paper. The authors declare competing financial interests: details are available in the online version of the paper. Correspondence and requests for materials should be addressed to J.K.J. (jjoung@mgh.harvard.edu).

METHODS

No statistical methods were used to predetermine sample size, and the investigators were not blinded to allocation during experiments and outcome assessment. **Plasmids and oligonucleotides.** DNA sequences for parent constructs used in this study can be found in Supplementary Information. Sequences of oligonucleotides used to generate the positive selection plasmids, negative selection plasmids, and site-depletion libraries are available in Supplementary Table 1. Sequences of all sgRNA targets in this study are available in Supplementary Table 2. Point mutations in Cas9 were generated by PCR. For cloning purposes, please note the low copy number origins of these plasmids. All new plasmids described in this study will be deposited with the non-profit plasmid distribution service Addgene: http://www.addgene.org/crispr-cas.

Bacterial Cas9/sgRNA expression plasmids were constructed with two T7 promoters to separately express Cas9 and the sgRNA. These plasmids encode human codon optimized versions of Cas9 for S. pyogenes (BPK764, SpCas9 sequence subcloned from JDS246; ref. 14), S. thermophilus Cas9 from CRISPR locus 1 (MSP1673, St1Cas9 sequence modified from previous published description¹⁷), and S. aureus (BPK2101, SaCas9 sequence codon optimized from Uniprot J7RUA5). Previously described sgRNA sequences were used for SpCas9 (refs 31, 32) and St1Cas9 (ref. 17), while the SaCas9 sgRNA sequence was determined by searching the European Nucleotide Archive sequence HE980450 for crRNA repeats using CRISPRfinder (http://crispr.u-psud.fr/Server/) and identifying the tracrRNA using a bioinformatic approach similar to one previously described³³. Annealed oligonucleotides to complete the spacer complementarity region of the sgRNA were ligated into BsaI-cut BPK764 and BPK2101, or BspMI-cut MSP1673 (append 5'-ATAG to the spacer to generate the top oligo and append 5'-AAAC to the reverse compliment of the spacer sequence to generate the bottom oligo). A 5'-GG dinucleotide was included on all bacterial plasmid sgRNAs for proper expression from the T7 promoter.

Residues 1097-1368 of SpCas9 were randomly mutagenized using Mutazyme II (Agilent Technologies) at a rate of ~5.2 substitutions/kilobase to generate mutagenized PAM-interacting domain libraries. For NGA PAM selections, wild-type SpCas9 and R1335Q were used as templates for mutagenesis. For NGC PAM selections, we first designed Cas9 mutants bearing amino acid substitutions of R1335 that might be expected to interact with a cytosine (D, E, S, or T) and found no activity on an NGC PAM site using the positive selection system (data not shown). We then randomly mutagenized the PAM-interacting domain of each of these singly substituted variants but still failed to obtain surviving colonies in positive selections (data not shown). Because the T1337R mutation had increased the activities of our VQR and EQR variants, we combined this mutation with R1335 substitutions of A, D, E, S, T, or V, and again randomly mutagenized their PAM-interacting domains. Selections using two of these six mutagenized libraries (bearing pre-existing R1335E/T1337R and R1335T/T1337R substitutions) yielded surviving colonies harbouring a variety of additional mutations (Extended Data Fig. 2b). The theoretical complexity of each PAM-interacting domain library was estimated to be greater than 10⁷ clones based on the number of transformants obtained. Positive and negative selection plasmids were generated by ligating annealed target site oligonucleotides into XbaI/SphI or EcoRI/SphI cut p11lacY-wtx115, respectively.

Two randomized PAM libraries (each with a different protospacer sequence) were constructed using Klenow(-exo) to fill-in the bottom strand of oligonucleotides that contained six randomized nucleotides directly adjacent to the 3' end of the protospacer (see Supplementary Table 1). The double-stranded product was cut with EcoRI to leave EcoRI/SphI ends for ligation into cut p11-lacY-wtx1. The theoretical complexity of each randomized PAM library was estimated to be greater than 10^6 based on the number of transformants obtained.

SpCas9 and variants were expressed in human cells from vectors derived from JDS246 (ref. 14). For St1Cas9 and SaCas9, the Cas9 ORFs from MSP1673 and BPK2101 were subcloned into a CAG promoter vector to generate MSP1594 and BPK2139, respectively. Plasmids for U6 expression of sgRNAs (into which desired spacer oligonucleotides can be cloned) were generated using the sgRNA sequences described above for the SpCas9 sgRNA (BPK1520), the St1Cas9 sgRNA (BPK2301), and the SaCas9 sgRNA (VVT1). Annealed oligonucleotides to complete the spacer complementarity region of the sgRNA were ligated into the BsmBI overhangs of these vectors (append 5'-CACC to the spacer to generate the poligo and append 5'-AAAC to the reverse complement of the spacer sequence to generate the bottom oligo). A 5'-G of target spacer sequences was included when designing human cell sgRNAs, for proper expression from the U6 promoter (and thus included in the calculation in Fig. 2e).

Bacterial-based positive selection assay for evolving SpCas9 variants. Competent *E. coli* BW25141(λ DE3)³⁴ containing a positive selection plasmid (with embedded target site) were transformed with Cas9/sgRNA-encoding plasmids. Following a 60 min recovery in SOB media, transformations were plated on LB plates containing either chloramphenicol (non-selective) or chloramphenicol + 10 mM arabinose (selective). Cleavage of the positive selection plasmid was estimated by calculating the survival frequency: colonies on selective plates/ colonies on non-selective plates (see also Extended Data Fig. 1).

To select for SpCas9 variants that can target novel PAMs, PAM-interactingdomain mutagenized Cas9/sgRNA plasmid libraries were electroporated into *E. coli* BW25141(λ DE3) cells containing a positive selection plasmid that encodes a target site and PAM of interest. Generally ~50,000 clones were screened to obtain between 50 and 100 survivors. The PAM-interacting domains of surviving clones were subcloned into fresh backbone plasmid and re-tested in the positive selection. Clones that had greater than 10% survival in this secondary screen for activity were sequenced. Mutations observed in the sequenced clones were chosen for further assessment based on their frequency in surviving clones, type of substitution, proximity to the PAM bases in the SpCas9–sgRNA crystal structure (PDB:4UN3)¹², and (in some cases) activities in a human cell-based EGFP disruption assay.

Bacterial-based site-depletion assay for profiling Cas9 PAM specificities. Competent *E. coli* BW25141(λ DE3) containing a Cas9/sgRNA expression plasmid were transformed with negative selection plasmids harbouring cleavable or non-cleavable target sites. Following a 60 min recovery in SOB media, transformations were plated on LB plates containing chloramphenicol + carbenicillin. Cleavage of the negative selection plasmid was estimated by calculating the colony forming units per µg of DNA transformed (see also Extended Data Fig. 3).

The negative selection was adapted to determine PAM specificity profiles of Cas9 nucleases by electroporating each randomized PAM library into *E. coli* BW25141(λ DE3) cells harbouring an appropriate Cas9/sgRNA plasmid. Between 80,000 and 100,000 colonies were plated at a low density spread on LB + chloramphenicol + carbenicillin plates. Surviving colonies containing negative selection plasmids refractory to cleavage by Cas9 were harvested and plasmid DNA isolated by maxi-prep (Qiagen). The resulting plasmid library was amplified by PCR using Phusion Hot-start Flex DNA Polymerase (New England BioLabs) followed by an Agencourt Ampure XP clean-up step (Beckman Coulter Genomics). Dual-indexed Tru-seq Illumina deep-sequencing libraries were prepared using the KAPA HTP library preparation kit (KAPA BioSystems) from ~500 ng of clean PCR product for each site-depletion experiment. The Dana-Farber Cancer Institute Molecular Biology Core performed 150-bp paired-end sequencing on an Illumina MiSeq Sequencer.

The raw FASTQ files outputted for each MiSeq run were analysed with a Python program to determine relative PAM depletion. The program (see Supplementary Information) operates as follows: first, a file dialogue is presented to the user from which all FASTQ read files for a given experiment can be selected. For these files, each FASTQ entry is scanned for the fixed spacer region on both strands. If the spacer region is found, then the six variable nucleotides flanking the spacer region are captured and added to a counter. From this set of detected variable regions, the count and frequency of each window of length 2-6 nucleotides at each possible position was tabulated (see Supplementary Table 3 for the 6-nucleotide output). The site-depletion data for both randomized PAM libraries was analysed by calculating the post-selection PAM depletion value (PPDV): the post-selection frequency of a PAM in the selected population divided by the pre-selection library frequency of that PAM. PPDV analyses were performed for each experiment across all possible 2-6 length windows in the 6-bp randomized region. The windows we used to visualize PAM preferences were: the 3-nucleotide window representing the second, third and fourth PAM positions for wild-type and variant SpCas9 experiments, and the 4-nucleotide window representing the third, fourth, fifth and sixth PAM positions for St1Cas9 and SaCas9.

Two significance thresholds for PPDVs were determined based on: (1) a statistical significance threshold based on the distribution of dCas9 versus pre-selection library log read count ratios (see Extended Data Fig. 3c, d), and (2) a biological activity threshold based on an empirical correlation between depletion values and activity in human cells. The statistical threshold was set at 3.36 s.d. from the mean PPDV for dCas9 (equivalent to a relative PPDV of 0.85), corresponding to a normal distribution two-sided P value of 0.05 after adjusting for multiple comparisons (that is, P = 0.05/64). The biological activity threshold was set at fivefold depletion (equivalent to a PPDV of 0.2) because this level of depletion serves as a reasonable predictor of activity in human cells (see also Extended Data Fig. 4). The 95% confidence intervals in Extended Data Fig. 4 were calculated by dividing the standard deviation of the mean by the square root of the sample size multiplied by 1.96. Human cell culture and transfection. U2OS cells obtained from our collaborator T. Cathomen (Freiburg) and U2OS.EGFP cells harbouring a single integrated copy of a constitutively expressed EGFP-PEST reporter gene13 were cultured in Advanced DMEM media (Life Technologies) supplemented with 10% FBS, 2 mM GlutaMAX (Life Technologies), penicillin/streptomycin, at 37 °C with 5%

CO₂. Additionally, U2OS.EGFP cells were cultured in 400 μ g ml⁻¹ of G418. The identity of U2OS and U2OS.EGFP cell lines were validated by STR profiling (ATCC) and deep sequencing, and cells were tested bi-weekly for mycoplasma contamination. Cells were co-transfected with 750 ng of Cas9 plasmid and 250 ng of sgRNA plasmid (unless otherwise noted) using the DN-100 program of a Lonza 4D-nucleofector according to the manufacturer's protocols. Cas9 plasmid transfected together with an empty U6 promoter plasmid was used as a negative control for spontaneous background EGFP loss for all human cell EGFP disruption experiments, and all endogenous gene disruption experiments (none of which showed detectable activity by T7E1). Target sites for endogenous gene experiments were selected within 200 bp of NGG sites cleavable by wild-type SpCas9 (see Extended Data Fig. 7a and Supplementary Table 2).

Zebrafish care and injections. Zebrafish care and use was approved by the Massachusetts General Hospital Subcommittee on Research Animal Care. Cas9 mRNA was transcribed with PmeI-digested JDS246 (wild-type SpCas9) or MSP469 (VQR variant) using the mMESSAGE mMACHINE T7 ULTRA Kit (Life Technologies) as previously described³². All sgRNAs in this study were prepared according to the cloning-independent sgRNA generation method³⁵. sgRNAs were transcribed by the MEGAscript SP6 Transcription Kit (Life Technologies), purified by RNA Clean & Concentrator-5 (Zymo Research), and eluted with RNase-free water.

sgRNA- and Cas9-encoding mRNA were co-injected into one-cell stage zebrafish embryos. Each embryo was injected with ${\sim}2{-}4.5$ nl of solution containing 30 ng ${\mu}l^{-1}$ sgRNA and 300 ng ${\mu}l^{-1}$ Cas9 mRNA. The next day, injected embryos were inspected under a stereoscope for normal morphological development, and genomic DNA was extracted from 5 to 9 embryos.

Human cell EGFP disruption assay. EGFP disruption experiments were performed as previously described¹⁴. Transfected cells were analysed for EGFP expression \sim 52 h post-transfection using a Fortessa flow cytometer (BD Biosciences). Background EGFP loss was gated at approximately 2.5% for all experiments (graphically represented as a dashed red line).

T7E1 assay, targeted deep-sequencing, and GUIDE-seq to quantify nucleaseinduced mutations. T7E1 assays were performed as previously described for human cells¹³ and zebrafish³². For human cells, genomic DNA was extracted from transfected cells ~72 h post-transfection using the Agencourt DNAdvance Genomic DNA Isolation Kit (Beckman Coulter Genomics). Target loci from zebrafish or human cell genomic DNA were amplified using the primers listed in Supplementary Table 1. Roughly 200 ng of purified PCR product was denatured, annealed, and digested with T7E1 (New England BioLabs). Mutagenesis frequencies were quantified using a Qiaxcel capillary electrophoresis instrument (Qiagen), as previously described for human cells¹³ and zebrafish³².

For targeted deep-sequencing, previously characterized on- and off-target sites^{7,14,27} were amplified using Phusion Hot-start Flex with the primers listed in Supplementary Table 1. Genomic loci were amplified for a control condition (empty sgRNA), wild-type, and D1135E SpCas9. An Agencourt Ampure XP clean-up step (Beckman Coulter Genomics) was performed before pooling ~500 ng of DNA from each condition for library preparation. Dual-indexed Tru-seq Illumina deep-sequencing libraries were generated using the KAPA HTP library preparation kit (KAPA BioSystems). The Dana-Farber Cancer

Institute Molecular Biology Core performed 150-bp paired-end sequencing on an Illumina MiSeq Sequencer. Mutation analysis of targeted deep-sequencing data was performed as previously described³⁰. Briefly, Illumina MiSeq paired end read data was mapped to human genome reference GRChr37 using bwa³⁶. High-quality reads (quality score \geq 30) were assessed for indel mutations that overlapped the target or off-target sites. 1-bp indel mutations were excluded from the analysis unless they occurred within 1-bp of the predicted breakpoint. Changes in activity at on- and off-target sites comparing D1135E versus wild-type SpCas9 were calculated by comparing the indel frequencies from both conditions (for rates above background control amplicon indel levels).

GUIDE-seq experiments were performed as previously described7. Briefly, 100 pmol of phosphorylated, phosphorothioate-modified double-stranded oligodeoxynucleotides (dsODNs) were transfected into U2OS cells along with Cas9 and sgRNA expression plasmids, as described above. dsODN-specific amplification, high-throughput sequencing, and mapping were performed to identify genomic intervals containing DSB activity. For wild-type versus D1135E experiments, offtarget read counts were normalized to the on-target read counts to correct for sequencing depth differences between samples. The normalized ratios for wildtype and D1135E SpCas9 were then compared to calculate the fold change in activity at off-target sites. To determine whether wild-type and D1135E samples for GUIDE-seq had similar oligo tag integration rates at the intended target site, restriction fragment length polymorphism (RFLP) assays were performed by amplifying the intended target loci with Phusion Hot-Start Flex from 100 ng of genomic DNA (isolated as described above) using primers listed in Supplementary Table 1. Roughly 500 ng of PCR product was digested with 20 U of NdeI (New England BioLabs) for 3 h at 37 °C before clean-up using the Agencourt Ampure XP kit. RFLP results were quantified using a Qiaxcel capillary electrophoresis instrument (Qiagen) to approximate oligo tag integration rates. T7E1 assays were performed for a similar purpose, as described above. For the quantitative comparison of wild-type to D1135E SpCas9, we utilized an alternative sequence consolidation algorithm that is more stringent and less likely to overestimate the number of unique molecularly-indexed GUIDE-seq reads. All sequencing data was corrected for cell-type specific single nucleotide polymorphisms.

Code availability. Custom code written to analyse PAM depletion MiSeq data is shown in the Supplementary Information.

- Mali, P. et al. RNA-guided human genome engineering via Cas9. Science 339, 823–826 (2013).
- Hwang, W. Y. et al. Efficient genome editing in zebrafish using a CRISPR-Cas system. Nature Biotechnol. 31, 227–229 (2013).
- Chylinski, K., Le Rhun, A. & Charpentier, E. The tracrRNA and Cas9 families of type II CRISPR-Cas immunity systems. RNA Biol. 10, 726–737 (2013).
- Kleinstiver, B. P., Fernandes, A. D., Gloor, G. B. & Edgell, D. R. A unified genetic, computational and experimental framework identifies functionally relevant residues of the homing endonuclease I-Bmol. *Nucleic Acids Res.* 38, 2411–2427 (2010).
- Gagnon, J. A. et al. Efficient mutagenesis by Cas9 protein-mediated oligonucleotide insertion and large-scale assessment of single-guide RNAs. PLoS ONE 9, e98186 (2014).
- Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760 (2009).





Extended Data Figure 1 | **Bacterial-based positive selection used to engineer altered PAM specificity variants of SpCas9.** a, Expanded schematic of the positive selection from Fig. 1b (left panel), and validation that SpCas9 behaves as expected in the positive selection (right panel). b, Schematic of how the positive selection was adapted to select for SpCas9 variants that have altered

PAM recognition specificities. A library of SpCas9 clones with randomized PAM-interacting (PI) domains (residues 1097–1368) is challenged by a selection plasmid that harbours an altered PAM. Variants that survive the selection by cleaving the positive selection plasmid are sequenced to determine the mutations that enable altered PAM specificity.



Extended Data Figure 2 Amino acid sequences of clones that cleave target sites bearing alternate PAMs in the bacterial-based positive selection system. a, Sequences of variants that survived >10% when re-tested in the positive selection assay against an NGA PAM site (see Methods). Variants were selected from libraries containing randomly mutagenized PAM-interacting domains (residues 1097–1368) with or without a starting R1335Q mutation. Sequence differences compared with wild-type SpCas9 are highlighted. The histogram represents the number of changes at each position (not counting the

starting R1335Q mutation). **b**, Sequences of variants that survived >10% when re-tested in the positive selection assay against a site containing an NGC PAM. Variants were selected from libraries containing randomly mutagenized PAM-interacting domains (residues 1097–1368) with starter mutation pairs of R1335E/T1337R or R1335T/T1337R. Sequence differences compared with wild-type SpCas9 (shown at the top) are highlighted. The histogram below illustrates the number of changes at each position (not counting starter mutations at R1335 or T1337).



Extended Data Figure 3 | Bacterial cell-based site-depletion assay for profiling the global PAM specificities of Cas9 nucleases. a, Expanded schematic illustrating the negative selection from Fig. 1d (left panel), and validation that wild-type SpCas9 behaves as expected in a screen of sites with functional (NGG) and non-functional (NGA) PAMs (right panel). b, Schematic of how the negative selection was used as a site-depletion assay to screen for functional PAMs by constructing negative selection plasmid libraries containing 6 randomized base pairs in place of the PAM. Selection plasmids that contain PAMs cleaved by a Cas9/sgRNA of interest are depleted while PAMs that are not cleaved (or poorly cleaved) are retained. The frequencies of the PAMs following selection are compared to their pre-selection frequencies in the starting libraries to calculate the post-selection PAM depletion value (PPDV). **c**, **d**, A cutoff for statistically significant PPDVs was established by plotting the PPDV of PAMs for catalytically inactive SpCas9 (dCas9) (grouped and plotted by their second/third/fourth positions) for the two randomized PAM libraries (**c**). A threshold of 3.36 standard deviations from the mean PPDV for the two libraries was calculated (red lines in (**d**)), establishing that any PPDV deviation below 0.85 is statistically significant compared to dCas9 treatment (red dashed line in (**c**)). The grey dashed line in (**c**) indicates a fivefold depletion in the assay (PPDV of 0.2).

LETTER RESEARCH



Extended Data Figure 4 Concordance between the site-depletion assay and EGFP disruption activity. Data points represent the average EGFP disruption of the two NGAN and NGNG PAM sites for the VQR and EQR variants (Fig. 1g) plotted against the mean PPDV observed for library 1 and 2 (Fig. 1f) for the corresponding PAM. The red dashed line indicates PAMs that are statistically significantly depleted (PPDV of 0.85, see Extended Data Fig. 3c), and the grey dashed line represents fivefold depletion (PPDV of 0.2). Mean values are plotted with the 95% confidence interval.



Extended Data Figure 5 Structural and functional roles of D1135, G1218, and T1337 in PAM recognition by SpCas9. a, Structural representations of the six residues implicated in PAM recognition. The left panel illustrates the proximity of D1135 to S1136, a residue that makes a water-mediated, minor groove contact to the third base position of the PAM¹². The right panel illustrates the proximity of G1218, E1219, and T1337 to R1335, a residue that makes a direct, base-specific major groove contact to the third base position of the PAM¹². Angstrom distances indicated by yellow dashed lines; non-target strand guanine bases dG2 and dG3 of the PAM are shown in blue; other DNA bases shown in orange; water molecules shown in red; images generated using

PyMOL from PDB:4UN3. **b**, Mutational analysis of six residues in SpCas9 that are implicated in PAM recognition. Clones containing one of three types of mutations at each position were tested for EGFP disruption with two sgRNAs targeted to sites harbouring NGG PAMs. For each position, we created an alanine substitution and two non-conservative mutations. S1136 and R1335 were previously reported to mediate contacts to the third guanine of the PAM¹², and D1135, G1218, E1219, and T1337 are reported in this study. EGFP disruption activities were quantified by flow cytometry; background control represented by the dashed red line; error bars represent s.e.m., n = 3.

th1 - Mutations in 15/17 sequences					
CGTAAGGAGCGCGAGGCGC	3CGGCCGC <mark>GGCGGCGGAGGCTGC</mark> 2	AGGACTGAGCGAGCAGATCG	IGTTTGAG	G Wild	l-type
CGTAAGAACCGCAAGGCA CGTAAGAACCGCAAGGC CGTAAGAACGCGAAGGC CGTAAGAACGCGAAGGCG CGTAAGAACGCGAAGCGC CGTAAGAACGCGAAGCGC CGTAAGAACGCGAAGGCG CGTAAGAACGCGAAGCGC CGTAAGAACCGCAAGGCG CGTAAGAACCGCAAGCGC	CGGCCGCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG	- GCAGCAGAGCAGATCG - GAGCAGACAGATCG - GAGCAGCAGATCG - GAGCAGACAGATCG - GAGCAGACCAGATCG - GCGAGCAGATCG - GCGAGCAGATCG - GCGAGCAGATCG - GCGAGCAGATCG - GCGAGCAGATCG - GCGAGCAGATCG	IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG IGTTTGAG	G -41 G -34 G -29 G -28 G -16 G -15 G -15 G -15 G -15 G -13 G -13 G -13 G -8	(-21,+5) (-16,+1) (-32,+19) (-37,+24)
CGTAAGGAGCGCGAGGCGC CGTAAGGAGCGCGAGGCGC CGTAAGGAGCGCGAGGCGC	CGGCCGCGGCGGCGGAGGCTGC CGGCCGCGGCGGCGGAGGCTG <mark>AG</mark> CCGCCGCGCGCGGCGGGGGGGGGGGGGGGGGGGGG	AG-ACTGAGCGAGCAGATCG gcgaGACTGAGCGAGCAGATCG AGcgagcagagcgagcagat	IGTTTGAG CGTGTTTG GTGTTTG	G -1 A +2 A +2	[2x] (-3,+5) (-17,+19)
tiall - TGTCGGGAACCTCTCCAGG TGTCGGGAACCTCT TGTCGGGAACCTCTCC TGTCGGGAACCTCTCC TGTCGGGAACCTCTCCAG TGTCGGGAACCTCTCCAG TGTCGGGAACCTCTCCAG TGTCGGGAACCTCTCCAG TGTCGGGAACCTCTCCAG TGTCGGGAACCTCTCCAG	Mutations in 17/27 GATGTTACGGAGGCCCTCATCC GTTACGGAGGCCCTCATCC TGTTACGGAGGCCCTCATCC TGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GATGTTACGGAGGCCCTCATCC GGATGTTACGGAGGCCCTCATCC GTTACGGATGTTACGGAGGCCCTC GTTACGGATGTTACGGAGGCCCTC GTTGGGATGTTACGGAGGCCCTC GTTGGGATGTTACGGAGGCCCTC	SEQUENCES IGCAAGTGTTCTCTCAGATC IGCAAGTGTTCTCTCAGATC IGCAAGTGTTCTCTCAGATC IGCAAGTGTTCTCTCAGATC IGCAAGTGTTCTCTCAGATC IGCAAGTGTTCTCTCAGATC IGCAAGTGTTCTCTCAGATC CTGCAAGTGTTCTCTCAGATC ATCCTGCAAGTGTTCTCTCA CATCCTGCAAGTGTTCTCTCC CATCCTGCAAGTGTTCTCTCC Laaataa	Wild-ty -12 (-1 -8 [X4 -5 -4 [X3 -1 -1 (-2 0 (-1 +1 (-4 +4 (-4 +5 (-1 +5 (-1 +20 (-1)	pe 5,+3)] ,+1) ,+1) ,+5) ,+8) ,+6) ,+6) 1,+31)	
fh – CATGGCGACCGGGGGC <mark>GG</mark> A	Mutations in 6/20 s	sequences cgagaatcggggggggggggggg	Wild-ty	be	
CATGGCGACCGGGGCGGZ CATGGCGACCGGGGGGGGG CATGGCGACCGGGGGGGGG CATGGCGACCGGGGGGGGG CATGGCGACCGGGGGGGGGZ CATGGCGACCGGGGGGGGGZ	ACTACTGC ACCAGAGG ACTACTGCTCT CAGAGG ACTACTGCTCT CCAGAGG ACTACTGCTCT CCAGAGG ACTACTGCTCC CCAGAGG ACTACTGCTCTCCagaggCCAG AGCTACTGCTCTCCtactgctct	CGAGAATCGGGGGGGGGGGGGGG CGAGAATCGGGGGGGGGG	$\begin{array}{c} -6 \\ -5 \\ -4 \\ -4 \\ (-5 \\ +3 \\ (-3 \\ +10 \\ (-2 \end{array})$,+1) ,+6) ,+12)	

Extended Data Figure 6 | Insertion or deletion mutations induced by the VQR SpCas9 variant at endogenous zebrafish sites containing NGAG PAMs. For each target locus, the wild-type sequence is shown at the top with the protospacer highlighted in yellow (highlighted in green if present on the complementary strand) and the PAM is marked as red underlined text. Deletions are shown as red dashes highlighted in grey and insertions as lower

case letters highlighted in blue. The net change in length caused by each indel mutation is shown on the right (+, insertion; -, deletion). Note that some alterations have both insertions and deletions of sequence and in these instances the alterations are enumerated in parentheses. The number of times each mutant allele was recovered (if more than once) is shown in brackets.



Extended Data Figure 7 | **Endogenous human genes targeted by wild-type and evolved variants of SpCas9.** a, Sequences targeted by wild-type, VQR, and VRER SpCas9 are shown in blue, red, and green, respectively. Sequences of sgRNAs and primers used to amplify these loci for T7E1 are provided in

Supplementary Tables 1 and 2. **b**, Mean mutagenesis frequencies detected by T7E1 for wild-type SpCas9 at eight target sites bearing NGG PAMs in the four different endogenous human genes (corresponding to the annotations in panel **a**). Error bars represent s.e.m., n = 3.

LETTER RESEARCH



Extended Data Figure 8 | Specificity profiles of the VQR and VRER SpCas9 variants determined using GUIDE-seq⁷. The intended on-target site is marked with a black square, and mismatched positions within off-target sites are highlighted. **a**, The specificity of the VQR variant was assessed in human cells by targeting endogenous sites containing NGA PAMs: *EMX1* site 4,

FANCF site 1, *FANCF* site 3, *FANCF* site 4, *RUNX1* site 1, *RUNX1* site 3, *VEGFA* site 1, and *ZNF629*. **b**, The specificity of the VRER variant was assessed in human cells by targeting endogenous sites containing NGCG PAMs: *FANCF* site 3, *FANCF* site 4, *RUNX1* site 1, *VEGFA* site 1, and *VEGFA* site 2.



Extended Data Figure 9 | **Activity differences between D1135E and wild-type SpCas9. a**, Mutagenesis frequencies detected by T7E1 for wild-type and D1135E SpCas9 at six endogenous sites in human cells. Error bars represent s.e.m., n = 3; mean fold change in activity is shown. **b**, Titration of the amount of wild-type or D1135E SpCas9-encoding plasmid transfected for EGFP disruption experiments in human cells. The amount of sgRNA plasmid used for all of these experiments was fixed at 250 ng. Two sgRNAs targeting different EGFP sites were used; error bars represent s.e.m., n = 3. **c**, Targeted deep-sequencing of on- and off-target sites for 3 sgRNAs using wild-type and D1135E SpCas9. The on-target site is shown at the top, with off-target sites listed below highlighting mismatches to the on-target sites greater than the change in activity at the on-target site are highlighted in green; control indel levels for each amplicon are reported. **d**, Mean frequency of GUIDE-seq oligo

tag integration at the on-target sites, estimated by restriction fragment length polymorphism analysis. Error bars represent s.e.m., n = 4. e, Mean mutagenesis frequencies at the on-target sites detected by T7E1 for GUIDE-seq experiments. Error bars represent s.e.m., n = 4. f, GUIDE-seq read count comparison between wild-type SpCas9 and D1135E at 3 endogenous human cell sites. The on-target site is shown at the top and off-target sites are listed below with mismatches highlighted. In the table, a ratio of off-target activity to on-target activity is compared between wild-type and D1135E to calculate the normalized fold changes in specificity (with gains in specificity highlighted in green). For sites without detectable GUIDE-seq reads, a value of 1 has been assigned to calculate an estimated change in specificity (indicated in orange). Off-target sites analysed by deep-sequencing in panel c are numbered to the left of the *EMX1* site 3 and *VEGFA* site 3 off-target sites.

LETTER RESEARCH



Extended Data Figure 10 | Additional PAMs for St1Cas9 and SaCas9 and activities based on spacer lengths in human cells. a, PPDV scatterplots for St1Cas9 comparing the sgRNA complementarity lengths of 20 and 21 nucleotides obtained with a randomized PAM library for spacers 1 and 2 (see also Fig. 4a). PAMs were grouped and plotted by their third/fourth/fifth/sixth positions. The red dashed line indicates PAMs that are statistically significantly depleted (see Extended Data Fig. 3c) and the grey dashed line represents fivefold depletion (PPDV of 0.2). b, Table of PAMs with PPDVs of less than 0.2 for St1Cas9 under each of the four conditions tested. PAM numbering shown on the left is the same as in Fig. 4a. c, PPDV scatterplots for SaCas9 comparing the sgRNA complementarity lengths of 21 and 23 nucleotides obtained with a randomized PAM library for spacers 1 and 2 (see also Fig. 4b). PAMs were

grouped and plotted by their third/fourth/fifth/sixth positions. The red and grey dashed lines are the same as in **a**. **d**, Table of PAMs with PPDVs of less than 0.2 for SaCas9 under each of the four conditions tested. PAM numbering shown on the left is the same as in Fig. 4b. **e**, Human cell EGFP disruption activities of St1Cas9 and SaCas9 at sites of various spacer lengths. Frequencies were quantified by flow cytometry; error bars represent s.e.m, n = 3 or 4; mean level of background EGFP loss represented by the dashed red line. **f**, Activity for all replicates of data shown in **e** plotted against spacer length. n = 3 or 4; bars illustrate mean and 95% confidence interval; number of sites per spacer length n = 3 or 4; bars illustrate mean and 95% confidence interval; number of sites per spacer length n = 3 or 4; bars illustrate mean and 95% confidence interval; number of sites per spacer length n = 3 or 4; bars illustrate mean and 95% confidence interval; number of sites per spacer length n = 3 or 4; bars illustrate mean and 95% confidence interval; number of sites per spacer length n = 3 or 4; bars illustrate mean and 95% confidence interval; number of sites per spacer length indicated.