

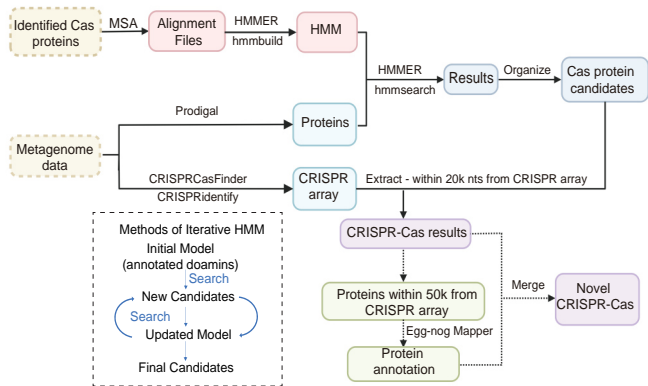
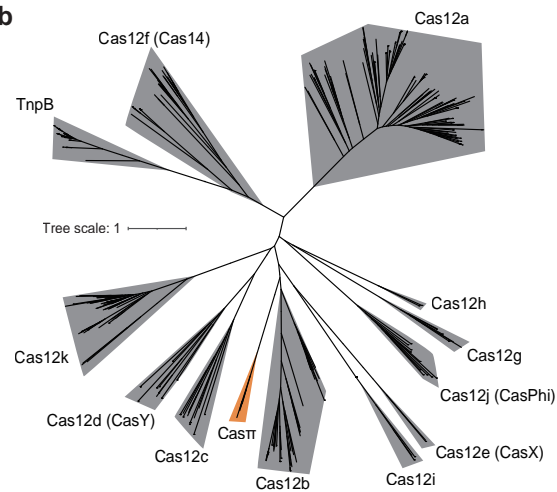
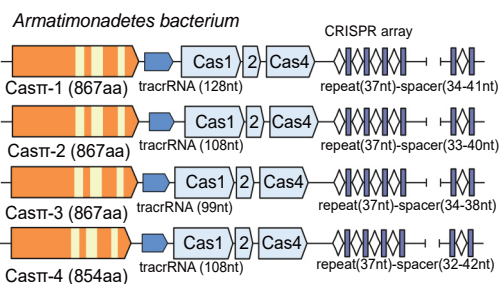
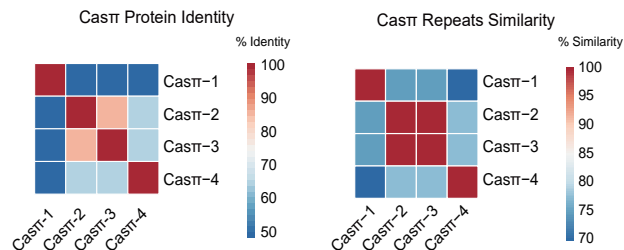
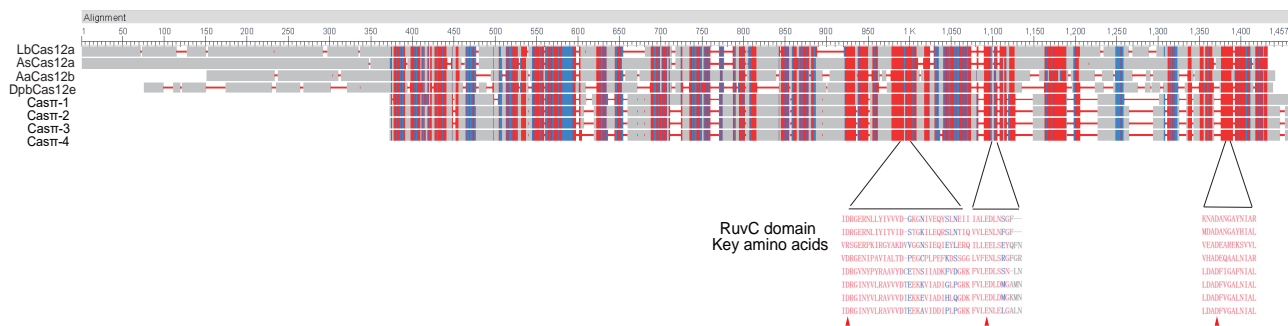
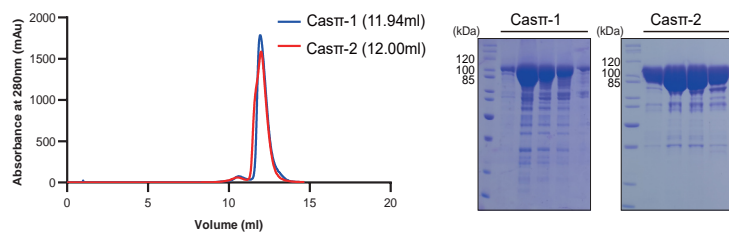
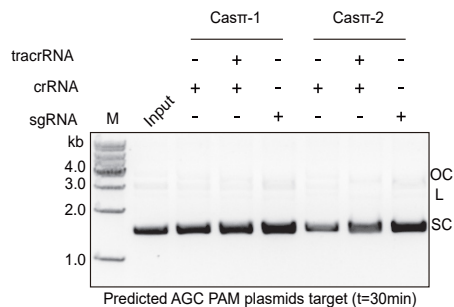
a**b****c****d****e****f****g**

Fig. S1 Cas π contains an active RuvC domain for DNA cleavage.

(a) Bioinformatics pipeline of identifying novel CRISPR-Cas systems. In general, the prediction of CRISPR-Cas systems is based on the identification of CRISPR arrays and Cas proteins. The architecture of CRISPR array is relatively conserved and can be predicted by public software, such as CRISPRCasFinder and CRISPRIdentify. Meanwhile, Cas proteins can be predicted by canonical protein homology search tools, such as HMMER and Pfam based on Hidden Markov models (HMM). However, these existed tools can only identify ‘new’ proteins which share relatively high sequence similarity to annotated Cas proteins. Therefore, we develop an unbiased workflow using iterative HMM to identify novel Cas candidates. For example, to identify novel type V Cas nucleases, RuvC domains from annotated Cas12a nucleases in the Pfam database (Pfam accession number: PF18516) were used to build an initial HMM for the first round searching of RuvC-containing proteins in all accessible databases. 1,349 protein candidates within 300 to 1600 amino acids (un-annotated in Pfam) were found, and all ‘RuvC’ domains within those candidates are used to build an improved HMM for the next round search of new candidates in all accessible databases. Using this iterative pipeline, the HMM was continuously updated with new candidates found in the current round and yielded to an updated pool of Cas-nuclease candidates from the next round search. A typical 2~10 rounds of iterative HMM and search were used for generating the final pool of novel type V Cas nucleases. Further, we filtered the candidates in the final pool by their functions annotated by eggNOG mapper, similarities to the final HMM (E-value), sizes and the confidence of putative CRISPR arrays for manual screening. Particularly in this study, Cas π showed up in the candidate pool with three rounds of iterative HMM and search.

(b) Maximum likelihood phylogenetic analysis of Cas π with type V subtypes a-k. Cas π proteins are outlined in orange, with other subtypes in grey. Bootstrap=1500, Cas π protein sequences are shown in Table S1.

- (c) Genomic architectures of all four CRISPR-Cas π systems.
- (d) Similarity matrix built and visualized using heatmap for Cas π protein sequences (left) and repeat sequences (right).
- (e) Multiple sequence alignment of Cas π with LbCas12a, AsCas12a, AaCas12b and DpbCas12e. Key amino acids of RuvC domain are marked by red arrows.
- (f) Purification of both Cas π proteins. The SEC curves were aligned referring to the injection volume. The peak fractions were analyzed by SDS-PAGE.
- (g) *In vitro* cleavage of plasmids containing predicted AGC PAM by crRNA only, tracrRNA and crRNA pair or sgRNA by Cas π effectors (SC, supercoiled plasmids; L, linearized plasmids; OC, open-circle plasmids).