

Cite as: H. Altae-Tran *et al.*, *Science*  
10.1126/science.abj6856 (2021).

# The widespread IS200/605 transposon family encodes diverse programmable RNA-guided endonucleases

Han Altae-Tran<sup>1,2,3,4,5†</sup>, Soumya Kannan<sup>1,2,3,4,5†</sup>, F. Esra Demircioglu<sup>1,2,3,4,5</sup>, Rachel Oshiro<sup>1,2,3,4,5</sup>, Suchita P. Nety<sup>1,2,3,4,5</sup>, Luke J. McKay<sup>6,7,8</sup>, Mensur Dlakic<sup>9</sup>, William P. Inskeep<sup>6,7</sup>, Kira S. Makarova<sup>10</sup>, Rhiannon K. Macrae<sup>1,2,3,4,5</sup>, Eugene V. Koonin<sup>10</sup>, Feng Zhang<sup>1,2,3,4,5\*</sup>

<sup>1</sup>Howard Hughes Medical Institute, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. <sup>2</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA. <sup>3</sup>McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. <sup>4</sup>Department of Brain and Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. <sup>5</sup>Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. <sup>6</sup>Department of Land Resources and Environmental Sciences, Montana State University, Bozeman, MT 59717, USA. <sup>7</sup>Thermal Biology Institute, Montana State University, Bozeman, MT 59717, USA. <sup>8</sup>Center for Biofilm Engineering, Montana State University, Bozeman, MT 59717, USA. <sup>9</sup>Department of Microbiology and Cell Biology, Montana State University, Bozeman, MT 59717, USA. <sup>10</sup>National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA.

†These authors contributed equally to this work.

\*Corresponding author. Email: zhang@broadinstitute.org

**IscB proteins are putative nucleases encoded in a distinct family of IS200/IS605 transposons and are likely ancestors of the RNA-guided endonuclease Cas9, but the functions of IscB and its interactions with any RNA remain uncharacterized. Using evolutionary analysis, RNA-seq, and biochemical experiments, we reconstruct the evolution of CRISPR-Cas9 systems from IS200/IS605 transposons. We show that IscB utilizes a single non-coding RNA for RNA-guided cleavage of double-stranded DNA and can be harnessed for genome editing in human cells. We also demonstrate the RNA-guided nuclease activity of TnpB, another IS200/605 transposon-encoded protein and the likely ancestor of Cas12 endonucleases. This work reveals a widespread class of transposon-encoded RNA-guided nucleases, which we name OMEGA (Obligate Mobile Element Guided Activity), with strong potential for developing as biotechnologies.**

The prokaryotic RNA-guided defense system CRISPR-Cas9 (type II CRISPR-Cas), which has been adopted for genome editing in eukaryotic cells (1, 2), is thought to have evolved from IscB proteins (3). Despite its wide distribution across prokaryotes and shared domain composition and architecture with Cas9, the function of IscB remains unknown (fig. S1). Moreover, given that IscB has not been reported to be associated with non-coding RNA (ncRNA) or CRISPR arrays, the evolutionary origins of the RNA-guided activity in Cas9 systems are unclear. IscB is encoded by a distinct subset of IS200/605 superfamily transposons that also include transposons encoding *tnpB*, a putative endonuclease distantly related to *iscB* and thought to be the ancestor of Cas12, the type V CRISPR effector (3–5). Using phylogenetic analysis, RNA-seq, and biochemical experiments, we sought to elucidate the functions of these proteins and the origin of RNA-guided activity in class 2 CRISPR systems.

## IscB is associated with an evolutionarily conserved non-coding RNA

IscB is ~400 amino acids (aa) long and contains a RuvC endonuclease domain split by the insertion of a bridge helix

(BH) and an HNH endonuclease domain, an architecture that is shared with Cas9 (Fig. 1A) (3). We performed a comprehensive search for proteins containing an HNH or a split RuvC endonuclease domain and found that Cas9 and IscB were the only proteins that contained both domains (data S1). This search also showed that IscB contains a previously unidentified N terminus that lacks clear homology to known domains and is absent in Cas9, which we denoted PLMP after its conserved sequence motifs (Fig. 1A and fig. S2). Clustering and phylogenetic analysis of the combined RuvC, BH, and HNH domains strongly suggests that all extant Cas9s descended from a single ancestral IscB (Fig. 1B and data S2 and S3). We searched for CRISPR arrays adjacent to *iscB* genes from each cluster and found 6 distinct groups of IscB, containing 16 clusters (of 603 total), that were CRISPR-associated, contrary to previous observations (3). CRISPR-associated IscBs were scattered around the IscB phylogenetic tree, suggesting they evolved independently, with one association event leading to the Cas9 lineage (Fig. 1B). In total we identified 31 unique CRISPR-associated *iscB* loci (of 2811 total).

Given their association with CRISPR arrays, we suspect-

ed that the rarely occurring CRISPR-associated IscBs may be RNA-guided nucleases. We first examined a cluster of CRISPR-associated IscBs similar to non-CRISPR associated IscBs (at ~50% aa identity). We heterologously expressed a representative locus from this clade in *E. coli* and performed small RNA-seq, which showed expression of not only the CRISPR array, but also a 329-bp intergenic region between the CRISPR array and the IscB open reading frame (ORF) (Fig. 1C). We purified the IscB protein and sequenced the co-purified RNA, demonstrating that this protein interacts with a single ncRNA component, encompassing both the CRISPR array and this intergenic region (Fig. 1C).

Given its interaction with a ncRNA that includes the CRISPR direct repeat (DR) and spacer, as well as its similar domain architecture to Cas9, we tested this IscB for RNA-guided endonuclease activity. Using a previously established protospacer adjacent motif (PAM)-discovery assay (table S1) (6), we observed depletion of specific PAM sequences (Fig. 1D and fig. S3), indicating that CRISPR-associated IscBs are reprogrammable RNA-guided nucleases. We confirmed this enzymatic activity with an in vitro cleavage assay using recombinant ribonucleoprotein (RNP) complexes (Fig. 1E).

Our finding that IscB functionally associated with CRISPR at least once, and likely on additional occasions, suggested that IscB systems more generally share a core ancestral ncRNA gene that is prone to evolving into a CRISPR array and in some cases a separate trans-acting *tracrRNA* (7). To test this hypothesis, we aligned 563 non-redundant *iscB* loci and searched for conserved nucleotide (nt) sequences either upstream or downstream of the *iscB* ORF. This analysis revealed a highly conserved intergenic region ~300 bp in length upstream of the ORF with a drop in conservation at the 5' end, which corresponds to an IS200/605 transposon end. Secondary structure predictions for individual sequences revealed the presence of multiple G:U pairs (fig. S4), suggesting that the conserved region encodes an ncRNA containing functionally important hairpins, which we named  $\omega$ RNA. Small RNA-seq on a sample of *Ktedonobacter racemifer* strain SOSP1-21, a soil bacterium that harbors 49 IscB loci in its genome (3), demonstrated expression of the predicted  $\omega$ RNA in many of these loci (Fig. 1F and figs. S5 and S6A). Moreover, we observed that the transcripts consistently extended beyond the conservation boundary at the 5' end.

An RFAM search for potential homologs of the  $\omega$ RNA showed that the conserved region of the  $\omega$ RNA partially matched the previously reported HEARO RNA, a ncRNA that was found upstream of HNH domain-containing proteins, which at the time were thought to be homing endonucleases (8, 9). However, the RFAM search did not provide any clues about the nature of the 5'-terminal non-conserved portion of these transcripts. Comparison of the consensus

CRISPR-associated IscB ncRNA and the covariance folded  $\omega$ RNA secondary structures revealed high degrees of structural and sequence similarity, particularly in shared multi-stem regions and pseudoknots (Fig. 1G, fig. S7, and supplementary text). Most importantly, we inferred that the 5'-most non-conserved sequence in the  $\omega$ RNA might function as a guide sequence, because the sequence immediately downstream was predicted to form hairpins that structurally resembled the hairpins formed by the DR/anti-repeat duplex in the CRISPR-associated IscB ncRNA (Fig. 1G).

### IscB is a reprogrammable RNA-guided DNA endonuclease

To test whether IscB was capable of cleaving DNA complementary to the putative  $\omega$ RNA guide, we performed an in vitro plasmid cleavage assay with KraIscB-1 using an in vitro transcription/translation (IVTT) expression system (Fig. 2, A and B). We found that KraIscB-1 cleaved the target in an  $\omega$ RNA-dependent manner, with an ATAAA 3' target-adjacent motif (TAM) (Fig. 2C). Retargeting of KraIscB-1 using a different guide (Fn guide) (6) also mediated cleavage of the cognate target (Fig. 2C and fig. S6B), implying that IscB is a reprogrammable RNA-guided nuclease.

Next, we biochemically characterized IscB in vitro. We identified activity in 57/86 (66%) selected phylogenetically diverse systems (table S2) as determined by the identification of a TAM (fig. S8). Of these 57 functional IscBs, 5 could be reconstituted with the respective  $\omega$ RNA in vitro to achieve efficient target cleavage, and from those, we selected AwaIscB (from *Allochromatium warmingii*) for detailed biochemical characterization (Fig. 2, D to G).

We confirmed the ability of recombinant AwaIscB to cleave multiple dsDNA targets in a programmable manner (Fig. 2E) and showed that the activity of AwaIscB is magnesium-dependent with a temperature optimum from 35–40°C (fig. S9, A and B). Appreciable activity was observed in vitro with guide lengths between 15 and 45 nt (fig. S9D). Mutation of the catalytic RuvC-II residue (E157A) abolished the nucleolytic activity on the non-target DNA strand, whereas the HNH domain catalytic mutant H212A abolished the nucleolytic activity on the target strand (Fig. 2F). Combination of the E157A and H212A mutations (dAwaIscB) abolished all dsDNA nucleolytic activity (Fig. 2F) (10, 11). Sequencing of the cleavage products showed that AwaIscB cleaves the target strand 3 nt upstream of the TAM, similar to Cas9s (12). Cleavage of the non-target strand occurred 8 or 12 nt upstream of the TAM, generating 5- or 9-nt long 5' overhangs (Fig. 2G and fig. S10). Exonuclease III mapping of a target substrate engaged by the dAwaIscB- $\omega$ RNA RNP showed that the RNP hindered exonuclease III treatment 19 nt upstream of the TAM on the target strand and 6 nt downstream of the targeted sequence on the non-target strand (fig. S11) (13). We

also found that truncation of more than 4 aa of the PLMP domain of AwaIscB abolished cleavage activity (fig. S12).

### IscB employ multiple guide-encoding mechanisms

A distinct advantage of RNA-guided systems is that they allow an effector to target many substrates by simply reprogramming the RNA guide. One way *IscB* evolved to use multiple guides is association with CRISPR arrays (Fig. 3A). However, given that *iscB* loci typically encode a single  $\omega$ RNA, it is unclear how or even whether these systems achieve such modularity in general. By searching for  $\omega$ RNAs not directly adjacent to *iscB* ORFs, we uncovered three additional potential mechanisms for guide encoding and switching:  $\omega$ RNA arrays, transposon expansion, and standalone, *trans*-acting  $\omega$ RNAs (Fig. 3A).  $\omega$ RNA arrays consist of multiple  $\omega$ RNAs, each encompassing a distinct guide, separated by up to 200 bp, and are found in 15/3356 unique *IscB*/*IsrB* loci (0.4%). Transposon expansion involves the insertion of nearly identical IS200/605 superfamily transposons in multiple locations, resulting in multiple loci per genome, each capable of expressing a nearly identical  $\omega$ RNA scaffold with a unique guide (fig. S13). By contrast, standalone  $\omega$ RNAs, which show no detectable genomic associations with *iscB*, were more common and were found in multiple copies in some genomes (table S3). *Cis*  $\omega$ RNAs from 95/3356 (2.8%) unique *IscB*/*IsrB* loci were nearly identical ( $\geq 95\%$  sequence identity) to distally encoded standalone  $\omega$ RNAs (fig. S14), implying that these standalone  $\omega$ RNAs could encode guides used by *trans*-encoded *IscBs*.

We tested this possibility by examining 10 standalone  $\omega$ RNAs in the *K. racemifer* genome, (Fig. 3B), 9 of which were found to be expressed (Fig. 3 and fig. S15). Of the 6 standalone  $\omega$ RNAs tested, we found that 5 could mediate RNA-guided DNA cleavage with a distally encoded *IscB* from the same genome (Fig. 3D), demonstrating that a single *IscB* can use multiple *trans*-encoded  $\omega$ RNAs. Guides from many  $\omega$ RNAs, both *IscB*-adjacent and *trans*-encoded, mostly target prokaryotic genomic sequences (61.5% genomic, 0.7% plasmid, 2.0% phage, 35.8% unmatched,  $N=36323$ ), suggesting a non-defense function for *IscB* systems (fig. S14 and table S3). In particular, we found that more than a third of the  $\omega$ RNAs (34.1%) targeted the same locus without the IS200/605 transposon insertion (table S3 and fig. S16).

### Evolution and diversity of *IscB* systems

We next investigated the evolutionary relationships between *IscB*, Cas9, and other homologous proteins to gain a broader insight into the evolution of RNA-guided mechanisms. In our search for proteins containing split RuvC domains, we detected another group of shorter, ~350 aa *IscB* homologs

that are also encoded in IS200/605 superfamily transposons. These proteins contain a PLMP domain and split RuvC but lack the HNH domain. We renamed these proteins *IsrB* (Insertion sequence RuvC-like OrfB) to emphasize their distinct domain architecture, replacing the previous designation, *IscB1* (3). In addition to *IscB* and *IsrB*, we identified a family of even smaller (~180 aa) proteins that only contained the PLMP domain and HNH domain but no RuvC domain, which we named *IshB* (Insertion sequence HNH-like OrfB).

To investigate the relationships between these proteins, we built a maximum likelihood (ML) tree from a multiple alignment of the split RuvC nuclease and BH domains using IQ-TREE 2 (Fig. 4A, figs. S17 and S18, data S2 and S3, and table S4) (14). The topology of the resulting tree was supported by several additional ML and Bayesian phylogenetic and robustness analyses (figs. S17 to S25 and data S2 and S3; see supplementary text for details). In the resulting tree, *IsrB*, *IscB*, and Cas9 formed distinct, strongly supported clades, suggesting that each of these nucleases originated from a unique evolutionary event (Fig. 4A, figs. S20, C and D, S21, S22, A and C, and S23, and supplementary text). We then analyzed the associations between each protein cluster and IS200/605 *tnpA* genes (3),  $\omega$ RNAs, CRISPR-Cas adaptation genes (*cas1*, *cas2*, *cas4*, and *csn2*), CRISPR arrays upstream and downstream of the respective ORF, and CRISPR anti-repeats (Fig. 4A). As discussed above, *IscB* and *isrB* were rarely associated with CRISPR arrays and were not found to be associated with CRISPR-Cas adaptation genes. The *isrBs* are associated with structurally distinct  $\omega$ RNAs. The *iscBs* are flanked by transposon ends similar to those mobilized by *TnpA* (3), but are only found near *tnpA* in 56/2811 of unique *IscB* loci (2.0%) (Fig. 4A and fig. S26D).

Additionally, we identified two distinct groups of Cas9s. The first is a new subtype, II-D, a group of relatively small *cas9s* (~700aa) that are not associated with any other known *cas* genes (15). The second is a distinct clade branching from within the II-C subtype, which includes exceptionally large *cas9s* (>1700aa) that are associated with *tnpA* (Fig. 4A and fig. S26). The *tnpA*-associated II-C loci often encompass unusually long DRs (more than 42bp in length) and in some cases encode HIRAN domain proteins between the *cas9* and other *cas* genes (Fig. 4A and fig. S27). Predicted transposon ends surround various combinations of the *tnpA*, *cas* acquisition genes, and CRISPR arrays in these loci.

These phylogenetic and association analyses confirm that IS200/605 transposon-encoded *IscBs* and *IsrBs* share a common evolutionary history with Cas9 (supplementary text). Given the deep position of the *IsrB* clade in the tree (Fig. 4A) and the lack of the HNH domain, *IsrBs* likely represent the ancestral state, probably having evolved from the compact RuvC endonuclease (16). Almost all *isrBs* are associated with

an  $\omega$ RNA, suggesting that these systems became RNA-guided at an early stage of evolution, concomitantly with the insertions in the RuvC-like domain that are likely to be involved in complex formation with  $\omega$ RNA. IsrB subsequently gained the HNH domain, possibly through insertion of another mobile element or recombination with a gene encoding an IshB-like protein, founding the IscB family (turquoise squares, Fig. 4, A and B, and supplementary text).

CRISPR arrays emerged within IscB systems on multiple, independent occasions (black circles, Fig. 4, A and B). These short arrays consist of repeats that could have evolved by duplication of segments of the ancestral  $\omega$ RNA. The resulting systems encompass a hybrid CRISPR- $\omega$ RNA that consists of a CRISPR array preceding a partial  $\omega$ RNA. These CRISPR-associated IscB proteins likely also gained REC-like insertions between the RuvC-I and RuvC-II subdomains on a number of occasions, often contemporaneously with or shortly after the CRISPR association (white squares, Fig. 4, A and B, and fig. S28). In particular, one CRISPR-associated IscB cluster (cluster 2089) apparently founded the Cas9 family (fig. S23) upon the loss of the hallmark PLMP domain (gray square, Fig. 4, A and B, and fig. S28). Moreover, the tracrRNAs of Subtype II-D, a deep branch in the Cas9 subtree (ML branch support:  $\geq 97/100$ , Bayesian posterior probability: 100%, figs. S20, B to D, and S23), shows significant similarity to IscB  $\omega$ RNAs (E-value 4.1e-8), suggesting that the Cas9 tracrRNA originally evolved from  $\omega$ RNA (fig. S29). The continued evolution of Cas9 apparently involved the gain of additional REC-like insertions between the bridge helix and the RuvC-II domains resulting in increased protein size (fig. S28). Finally, upon the association with the CRISPR adaptation machinery (*cas1*, *cas2*, and possibly *cas4*) (light blue circles, Fig. 4, A and B), a burst of Cas9 diversification and widespread dispersion among bacteria via horizontal gene transfer followed, resulting in the evolution of multiple type II CRISPR subtypes.

We also explored the evolutionary history of  $\omega$ RNAs. By iteratively building a set of  $\omega$ RNA profiles that spanned all major groups of  $\omega$ RNAs associated with *iscBs* and *isrBs*, we found that diverse  $\omega$ RNAs are associated with almost all *iscBs* and *isrBs*. Moreover, different IsrB and IscB clades are associated with distinct  $\omega$ RNA structures (Fig. 4, A and C, and figs. S18A, S24A, and S30). The transition from *isrB* to *iscB* was likely accompanied by loss of a second pseudoknot, the adaptor pseudoknot, between the transposon end region and the multi-stem loop in *isrB*-associated  $\omega$ RNAs (yellow square, Fig. 4, A to C). The inverse relationship between the complexity of the  $\omega$ RNA structure and the associated protein size is also reflected by the simplified  $\omega$ RNA structures associated with clades of large IscBs and the even smaller tracrRNAs associated with large Cas9s (Fig. 4C and fig. S30).

## IS200/IS605 elements encode diverse RNA-guided nucleases

In addition to the distinct succession of evolutionary events that yielded the abundant and diverse type II CRISPR systems, our phylogenetic analysis revealed several other events in the evolution of IscB and related proteins that led to the extant diversity, which we sought to experimentally explore.

First, we searched for IscB homologs in eukaryotic genomes and identified multiple *iscB* loci in the chloroplast genome of *Ignatius tetrasporus* UTEX B 2012, a terrestrial green alga (Fig. 5, A and B, and fig. S31). Although the ORF is disrupted by multiple stop codons in most of these loci, one locus encodes an intact IscB (~50% aa identity to related prokaryotic IscBs) and a transcriptionally active  $\omega$ RNA (Fig. 5C). This eukaryotic IscB cleaves DNA with a minimal NNG TAM (Fig. 5D), which differs from other characterized IscB TAMs (fig. S8).

Second, we investigated the clade of large IscBs, which contain a BH domain that is split in two by REC domain-like insertions (white squares, Figs. 4A and 5A). We hypothesized that these insertions might enhance DNA unwinding, similarly to the REC lobe of Cas9 (*I7*) and would therefore facilitate genome editing in the complex landscape of eukaryotic chromatin structure. We screened 6 large IscB proteins, using a pool of 12 guides each, for their ability to generate insertions/deletions (indels) in HEK293FT cells (see methods and table S5); one (OgeuIscB) produced appreciable indels (Fig. 5, E and F, and fig. S32A). To further examine OgeuIscB activity, we tested a range of guide lengths targeting 3 loci in the human genome and found that OgeuIscB achieved the maximum indel rate with a 16 nt guide (fig. S32B). On a panel of 46 sites in the human genome, we found that OgeuIscB induced indels at 28 of these sites with varying efficiency up to 4.4% (Fig. 5G, fig. S32C, and table S5). Thus, OgeuIscB seems a promising candidate for further development of IscB-based genome editing tools.

Third, we experimentally characterized the putative nuclease activity of IsrB, the apparent ancestor of IscB (Fig. 5A). *K. racemifer* contains 5 *isrBs* associated with  $\omega$ RNAs that are natively expressed (Fig. 5H and fig. S33). We found that the IsrB- $\omega$ RNA RNP nicks the non-target strand of a dsDNA substrate in a guide- and TAM-specific manner (Fig. 5, I and J, and fig. S34), which is analogous to the activity of IscB upon inactivation of the HNH domain (Fig. 2F).

Finally, we sought to determine if IS200/605 transposons in general harbor RNA-guided nucleases. In addition to the distinct IscB and IsrB families, most IS200/IS605 transposons encode RuvC-like endonucleases of another family, TnpB, which is thought to be the ancestor of Cas12s, the type V CRISPR effectors (Fig. 5A) (5). Additionally, TnpB is

the likely ancestor of larger proteins, Fanzors, encoded in diverse eukaryotic transposons (Fig. 5A) (18). The TnpB family, including Fanzor, is an order of magnitude more diverse than the IscB family; an HMMER search identified more than a million *tnpB* loci in publicly available prokaryotic genomes.

We identified conserved non-coding regions immediately downstream of the CDS of many *tnpBs*, suggesting the presence of associated ncRNAs that could function as RNA guides (fig. S35). Previous work has identified ncRNAs overlapping the 3'-end of *tnpB* genes in archaea and bacteria (19, 20), but the function of these ncRNAs has not been characterized. Small RNA-seq of *K. racemifer* revealed native expression of a ncRNA overlapping the 3' end of the associated *tnpB* ORF (Fig. 5K), which we classified as a distinct group of  $\omega$ RNAs. The reverse complement of the KraTnpB  $\omega$ RNA 3' end is nearly identical to the 5' of the  $\omega$ RNA associated with some KraIscBs, a region that corresponds to the predicted transposon end in each locus (Fig. 5L).

Analysis of non-redundant loci containing *tnpB* genes that clustered with KraTnpB showed a drop of sequence conservation at the 3' end of the loci (fig. S35), corresponding to the IS200/605 transposon end. Comparison to the small RNA-seq trace revealed expression beyond the conservation drop, indicating possible presence of a guide sequence in the transcript (Fig. 5M). In vitro plasmid cleavage assays for multiple TnpB proteins from this cluster using a reprogrammed guide demonstrated RNA-guided cleavage with a 5' TAM (Fig. 5N and fig. S36). We recombinantly purified a TnpB from *Alicyclobacillus macrosporangioides* (AmaTnpB) and confirmed its reprogrammable RNA-guided dsDNA endonuclease activity (Fig. 5O and fig. S36). We also observed that AmaTnpB robustly cleaved target-containing ssDNA substrates (Fig. 5P) and non-specifically cleaved a collateral substrate upon recognition of dsDNA or ssDNA substrates (Fig. 5Q).

## Discussion

Naturally programmable biological systems offer an efficient solution for diverse organisms to achieve scalable complexity via modularity of their components. RNA-guided defense and regulatory systems, which are widespread in prokaryotes and eukaryotes, are a prominent case in point, and have served as the basis of numerous biotechnology applications thanks to the ease with which they can be engineered and reprogrammed (21–23).

Here, through the exploration of Cas9 evolution, we discovered the programmable RNA-guided mechanism of 3 highly abundant but previously uncharacterized transposon-encoded nucleases: IscB, IsrB, and TnpB, which we collectively refer to as  $\Omega$  (OMEGA: Obligate Mobile Element

Guided Activity) (Fig. 6) because the mobile element localization and movement likely determines the identity of their guides. Although the biological functions of  $\Omega$  systems remain unknown, several hypotheses are compatible with the available evidence, including roles in facilitating TnpA-catalyzed, RNA-guided transposition, or acting as a toxin, with the transposon acting as the antitoxin, securing maintenance of IS200/605 insertions (supplementary text).

The broad distribution of the  $\Omega$  systems characterized here indicates that RNA-guided mechanisms are more widespread in prokaryotes than previously suspected and suggests that RNA-guided activities are likely ancient and evolved on multiple, independent occasions, of which only the most common ones have likely been identified so far. The TnpB family is far more abundant and diverse than the IscB family; indeed, we identified more than a million putative *tnpB* loci in bacterial and archaeal genomes, making it one of the most common prokaryotic genes altogether. These TnpBs might represent an untapped wealth of diverse RNA-guided mechanisms present not only in prokaryotes, but also in eukaryotes. Combined with our identification of a chloroplast-encoded IscB, these findings suggest that the expansion of RNA-guided systems into eukaryotic genomes could be a general phenomenon, and more broadly, that RNA-guided systems are functionally diverse and permeate all domains of life.

## REFERENCES AND NOTES

1. F. Zhang, Development of CRISPR-Cas systems for genome editing and beyond. *Q. Rev. Biophys.* **52**, e6 (2019). [doi:10.1017/S0033583519000052](https://doi.org/10.1017/S0033583519000052)
2. F. Hille, H. Richter, S. P. Wong, M. Bratovič, S. Ressel, E. Charpentier, The biology of CRISPR-Cas: Backward and forward. *Cell* **172**, 1239–1259 (2018). [doi:10.1016/j.cell.2017.11.032](https://doi.org/10.1016/j.cell.2017.11.032) [Medline](#)
3. V. V. Kapitonov, K. S. Makarova, E. V. Koonin, ISC, a novel group of bacterial and Archaeal DNA transposons that encode Cas9 homologs. *J. Bacteriol.* **198**, 797–807 (2015). [doi:10.1128/JB.00783-15](https://doi.org/10.1128/JB.00783-15) [Medline](#)
4. P. Siguier, E. Gourbeyre, M. Chandler, Bacterial insertion sequences: Their genomic impact and diversity. *FEMS Microbiol. Rev.* **38**, 865–891 (2014). [doi:10.1111/1574-6976.12067](https://doi.org/10.1111/1574-6976.12067) [Medline](#)
5. S. Shmakov, A. Smargon, D. Scott, D. Cox, N. Pyzocha, W. Yan, O. O. Abudayyeh, J. S. Gootenberg, K. S. Makarova, Y. I. Wolf, K. Severinov, F. Zhang, E. V. Koonin, Diversity and evolution of class 2 CRISPR-Cas systems. *Nat. Rev. Microbiol.* **15**, 169–182 (2017). [doi:10.1038/nrmicro.2016.184](https://doi.org/10.1038/nrmicro.2016.184) [Medline](#)
6. B. Zetsche, J. S. Gootenberg, O. O. Abudayyeh, I. M. Slaymaker, K. S. Makarova, P. Essletzbichler, S. E. Volz, J. Joung, J. van der Oost, A. Regev, E. V. Koonin, F. Zhang, Cpf1 is a single RNA-guided endonuclease of a class 2 CRISPR-Cas system. *Cell* **163**, 759–771 (2015). [doi:10.1016/j.cell.2015.09.038](https://doi.org/10.1016/j.cell.2015.09.038) [Medline](#)
7. E. Deltcheva, K. Chylinski, C. M. Sharma, K. Gonzales, Y. Chao, Z. A. Pirzada, M. R. Eckert, J. Vogel, E. Charpentier, CRISPR RNA maturation by trans-encoded small RNA and host factor RNase III. *Nature* **471**, 602–607 (2011). [doi:10.1038/nature09886](https://doi.org/10.1038/nature09886) [Medline](#)
8. I. Kalvari, E. P. Nawrocki, N. Ontiveros-Palacios, J. Argasinska, K. Lamkiewicz, M. Marz, S. Griffiths-Jones, C. Toffano-Nioche, D. Gautheret, Z. Weinberg, E. Rivas, S. R. Eddy, R. D. Finn, A. Bateman, A. I. Petrov, Rfam 14: Expanded coverage of metagenomic, viral and microRNA families. *Nucleic Acids Res.* **49**, D192–D200

- (2021). [doi:10.1093/nar/gkaa1047](https://doi.org/10.1093/nar/gkaa1047) [Medline](#)
9. Z. Weinberg, J. Perreault, M. M. Meyer, R. R. Breaker, Exceptional structured noncoding RNAs revealed by bacterial metagenome analysis. *Nature* **462**, 656–659 (2009). [doi:10.1038/nature08586](https://doi.org/10.1038/nature08586) [Medline](#)
  10. M. Jinek, K. Chylinski, I. Fonfara, M. Hauer, J. A. Doudna, E. Charpentier, A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. *Science* **337**, 816–821 (2012). [doi:10.1126/science.1225829](https://doi.org/10.1126/science.1225829) [Medline](#)
  11. G. Gasiunas, R. Barrangou, P. Horvath, V. Siksnys, Cas9-crRNA ribonucleoprotein complex mediates specific DNA cleavage for adaptive immunity in bacteria. *Proc. Natl. Acad. Sci. U.S.A.* **109**, E2579–E2586 (2012). [doi:10.1073/pnas.1208507109](https://doi.org/10.1073/pnas.1208507109) [Medline](#)
  12. G. Gasiunas, J. K. Young, T. Karvelis, D. Kazlauskas, T. Urbaitis, M. Jasnauskaitė, M. M. Grusyte, S. Paulraj, P.-H. Wang, Z. Hou, S. K. Dooley, M. Cigan, C. Alarcon, N. D. Chilcoat, G. Bigelyte, J. L. Curcuru, M. Mabuchi, Z. Sun, R. T. Fuchs, E. Schildkraut, P. R. Weigle, W. E. Jack, G. B. Robb, Č. Venclovas, V. Siksnys, A catalogue of biochemically diverse CRISPR-Cas9 orthologs. *Nat. Commun.* **11**, 5512 (2020). [doi:10.1038/s41467-020-19344-1](https://doi.org/10.1038/s41467-020-19344-1) [Medline](#)
  13. M. Jinek, F. Jiang, D. W. Taylor, S. H. Sternberg, E. Kaya, E. Ma, C. Anders, M. Hauer, K. Zhou, S. Lin, M. Kaplan, A. T. Iavarone, E. Charpentier, E. Nogales, J. A. Doudna, Structures of Cas9 endonucleases reveal RNA-mediated conformational activation. *Science* **343**, 1247997 (2014). [doi:10.1126/science.1247997](https://doi.org/10.1126/science.1247997) [Medline](#)
  14. B. Q. Minh, H. A. Schmidt, O. Chernomor, D. Schrempf, M. D. Woodhams, A. von Haeseler, R. Lanfear, IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020). [doi:10.1093/molbev/msaa015](https://doi.org/10.1093/molbev/msaa015) [Medline](#)
  15. K. S. Makarova, Y. I. Wolf, J. Iranzo, S. A. Shmakov, O. S. Alkhnbashi, S. J. J. Brouns, E. Charpentier, D. Cheng, D. H. Haft, P. Horvath, S. Moineau, F. J. M. Mojica, D. Scott, S. A. Shah, V. Siksnys, M. P. Terns, Č. Venclovas, M. F. White, A. F. Yakunin, W. Yan, F. Zhang, R. A. Garrett, R. Backofen, J. van der Oost, R. Barrangou, E. V. Koonin, Evolutionary classification of CRISPR-Cas systems: A burst of class 2 and derived variants. *Nat. Rev. Microbiol.* **18**, 67–83 (2020). [doi:10.1038/s41579-019-0299-x](https://doi.org/10.1038/s41579-019-0299-x) [Medline](#)
  16. K. A. Majorek, S. Dunin-Horkawicz, K. Steczkiewicz, A. Muszewska, M. Nowotny, K. Ginalski, J. M. Bujnicki, The RNase H-like superfamily: New members, comparative structural analysis and evolutionary classification. *Nucleic Acids Res.* **42**, 4160–4179 (2014). [doi:10.1093/nar/gkt1414](https://doi.org/10.1093/nar/gkt1414) [Medline](#)
  17. H. Nishimasu, F. A. Ran, P. D. Hsu, S. Konermann, S. I. Shehata, N. Dohmae, R. Ishitani, F. Zhang, O. Nureki, Crystal structure of Cas9 in complex with guide RNA and target DNA. *Cell* **156**, 935–949 (2014). [doi:10.1016/j.cell.2014.02.001](https://doi.org/10.1016/j.cell.2014.02.001) [Medline](#)
  18. W. Bao, J. Jurka, Homologues of bacterial TnpB<sub>IS605</sub> are widespread in diverse eukaryotic transposable elements. *Mob. DNA* **4**, 12 (2013). [doi:10.1186/1759-8753-4-12](https://doi.org/10.1186/1759-8753-4-12) [Medline](#)
  19. J. V. Gomes-Filho, L. S. Zaramela, V. C. S. Italiani, N. S. Baliga, R. Z. N. Vêncio, T. Koide, Sense overlapping transcripts in IS1341-type transposase genes are functional non-coding RNAs in archaea. *RNA Biol.* **12**, 490–500 (2015). [doi:10.1080/15476286.2015.1019998](https://doi.org/10.1080/15476286.2015.1019998) [Medline](#)
  20. Z. Weinberg, C. E. Lünse, K. A. Corbino, T. D. Ames, J. W. Nelson, A. Roth, K. R. Perkins, M. E. Sherlock, R. R. Breaker, Detection of 224 candidate structured RNAs by comparative analysis of specific subsets of intergenic regions. *Nucleic Acids Res.* **45**, 10811–10823 (2017). [doi:10.1093/nar/gkx699](https://doi.org/10.1093/nar/gkx699) [Medline](#)
  21. A. Hüttenhofer, P. Schattner, The principles of guiding by RNA: Chimeric RNA-protein enzymes. *Nat. Rev. Genet.* **7**, 475–482 (2006). [doi:10.1038/nrg1855](https://doi.org/10.1038/nrg1855) [Medline](#)
  22. A. Schneider, A short history of guide RNAs. *EMBO Rep.* **21**, e51918 (2020). [doi:10.15252/embr.202051918](https://doi.org/10.15252/embr.202051918) [Medline](#)
  23. E. V. Koonin, Evolution of RNA- and DNA-guided antiviral defense systems in prokaryotes and eukaryotes: Common ancestry vs convergence. *Biol. Direct* **12**, 5 (2017). [doi:10.1186/s13062-017-0172-2](https://doi.org/10.1186/s13062-017-0172-2) [Medline](#)
  24. H. Altae-Tran, S. Kannan, F. Zhang, Code and processed data for: The widespread IS200/605 transposon family encodes diverse programmable RNA-guided endonucleases (Version 1.0). Zenodo. <https://doi.org/10.5281/zenodo.5168777>
  25. C. Camacho, G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, T. L. Madden, BLAST+: Architecture and applications. *BMC Bioinformatics* **10**, 421 (2009). [doi:10.1186/1471-2105-10-421](https://doi.org/10.1186/1471-2105-10-421) [Medline](#)
  26. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7: Improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013). [doi:10.1093/molbev/mst010](https://doi.org/10.1093/molbev/mst010) [Medline](#)
  27. M. Steinegger, J. Söding, MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017). [doi:10.1038/nbt.3988](https://doi.org/10.1038/nbt.3988) [Medline](#)
  28. M. Steinegger, M. Meier, M. Mirdita, H. Vöhringer, S. J. Haunsberger, J. Söding, HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinformatics* **20**, 473 (2019). [doi:10.1186/s12859-019-3019-7](https://doi.org/10.1186/s12859-019-3019-7) [Medline](#)
  29. S. R. Eddy, Accelerated Profile HMM Searches. *PLOS Comput. Biol.* **7**, e1002195 (2011). [doi:10.1371/journal.pcbi.1002195](https://doi.org/10.1371/journal.pcbi.1002195) [Medline](#)
  30. S. A. Shmakov, K. S. Makarova, Y. I. Wolf, K. V. Severinov, E. V. Koonin, Systematic prediction of genes functionally linked to CRISPR-Cas systems by gene neighborhood analysis. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E5307–E5316 (2018). [doi:10.1073/pnas.1803440115](https://doi.org/10.1073/pnas.1803440115) [Medline](#)
  31. A. B. Crawley, J. R. Henriksen, R. Barrangou, CRISPRdisco: An Automated Pipeline for the Discovery and Analysis of CRISPR-Cas Systems. *CRISPR J.* **1**, 171–181 (2018). [doi:10.1089/crispr.2017.0022](https://doi.org/10.1089/crispr.2017.0022) [Medline](#)
  32. D. H. Haft, B. J. Loftus, D. L. Richardson, F. Yang, J. A. Eisen, I. T. Paulsen, O. White, TIGRFAMs: A protein family resource for the functional identification of proteins. *Nucleic Acids Res.* **29**, 41–43 (2001). [doi:10.1093/nar/29.1.41](https://doi.org/10.1093/nar/29.1.41) [Medline](#)
  33. A. L. Delcher, D. Harmon, S. Kasif, O. White, S. L. Salzberg, Improved microbial gene identification with GLIMMER. *Nucleic Acids Res.* **27**, 4636–4641 (1999). [doi:10.1093/nar/27.23.4636](https://doi.org/10.1093/nar/27.23.4636) [Medline](#)
  34. S. Naser-Khdour, B. Q. Minh, W. Zhang, E. A. Stone, R. Lanfear, The Prevalence and Impact of Model Violations in Phylogenetic Analysis. *Genome Biol. Evol.* **11**, 3341–3352 (2019). [doi:10.1093/gbe/evz193](https://doi.org/10.1093/gbe/evz193) [Medline](#)
  35. S. Kalyaanamoorthy, B. Q. Minh, T. K. F. Wong, A. von Haeseler, L. S. Jermiin, ModelFinder: Fast model selection for accurate phylogenetic estimates. *Nat. Methods* **14**, 587–589 (2017). [doi:10.1038/nmeth.4285](https://doi.org/10.1038/nmeth.4285) [Medline](#)
  36. M. N. Price, P. S. Dehal, A. P. Arkin, FastTree 2—Approximately maximum-likelihood trees for large alignments. *PLOS ONE* **5**, e9490 (2010). [doi:10.1371/journal.pone.0009490](https://doi.org/10.1371/journal.pone.0009490) [Medline](#)
  37. G. Altekar, S. Dwarkadas, J. P. Huelsenbeck, F. Ronquist, Parallel Metropolis coupled Markov chain Monte Carlo for Bayesian phylogenetic inference. *Bioinformatics* **20**, 407–415 (2004). [doi:10.1093/bioinformatics/btg427](https://doi.org/10.1093/bioinformatics/btg427) [Medline](#)
  38. A. Stamatakis, RAxML version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014). [doi:10.1093/bioinformatics/btu033](https://doi.org/10.1093/bioinformatics/btu033) [Medline](#)
  39. R. Lorenz, S. H. Bernhart, C. Höner Zu Siederdissen, H. Tafer, C. Flamm, P. F. Stadler, I. L. Hofacker, ViennaRNA Package 2.0. *Algorithms Mol. Biol.* **6**, 26 (2011). [doi:10.1186/1748-7188-6-26](https://doi.org/10.1186/1748-7188-6-26) [Medline](#)
  40. E. Rivas, J. Clements, S. R. Eddy, A statistical test for conserved RNA structure shows lack of evidence for structure in lncRNAs. *Nat. Methods* **14**, 45–48 (2017). [doi:10.1038/nmeth.4066](https://doi.org/10.1038/nmeth.4066) [Medline](#)
  41. E. P. Nawrocki, S. R. Eddy, Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics* **29**, 2933–2935 (2013). [doi:10.1093/bioinformatics/btt509](https://doi.org/10.1093/bioinformatics/btt509) [Medline](#)
  42. Z. Weinberg, R. R. Breaker, R2R—Software to speed the depiction of aesthetic consensus RNA secondary structures. *BMC Bioinformatics* **12**, 3 (2011). [doi:10.1186/1471-2105-12-3](https://doi.org/10.1186/1471-2105-12-3) [Medline](#)

43. J. Mistry, S. Chuguransky, L. Williams, M. Qureshi, G. A. Salazar, E. L. L. Sonnhammer, S. C. E. Tosatto, L. Paladin, S. Raj, L. J. Richardson, R. D. Finn, A. Bateman, Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–D419 (2021). [doi:10.1093/nar/gkaa913](https://doi.org/10.1093/nar/gkaa913) [Medline](#)
44. C. Bland, T. L. Ramsey, F. Sabree, M. Lowe, K. Brown, N. C. Kyrpides, P. Hugenholtz, CRISPR recognition tool (CRT): A tool for automatic detection of clustered regularly interspaced palindromic repeats. *BMC Bioinformatics* **8**, 209 (2007). [doi:10.1186/1471-2105-8-209](https://doi.org/10.1186/1471-2105-8-209) [Medline](#)
45. F. Asnicar, G. Weingart, T. L. Tickle, C. Huttenhower, N. Segata, Compact graphical representation of phylogenetic data and metadata with GraPhlAn. *PeerJ* **3**, e1029 (2015). [doi:10.7717/peerj.1029](https://doi.org/10.7717/peerj.1029) [Medline](#)
46. L. Maaten, G. Hinton, Visualizing Data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
47. P. G. Poličar, M. Stražar, B. Zupan, Embedding to reference t-SNE space addresses batch effects in single-cell classification. *bioRxiv* [671404](https://doi.org/10.1101/067140) [preprint]. 14 June 2019.
48. P. J. Rousseeuw, Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* **20**, 53–65 (1987). [doi:10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7)
49. R. J. G. B. Campello, D. Moulavi, J. Sander, in *Advances in Knowledge Discovery and Data Mining* (Springer, 2013), pp. 160–172.
50. D. H. Parks, M. Imelfort, C. T. Skenner, P. Hugenholtz, G. W. Tyson, CheckM: Assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015). [doi:10.1101/gr.186072.114](https://doi.org/10.1101/gr.186072.114) [Medline](#)
51. P.-A. Chaumeil, A. J. Mussig, P. Hugenholtz, D. H. Parks, GTDB-Tk: A toolkit to classify genomes with the Genome Taxonomy Database. *Bioinformatics* **btz848** (2019). [doi:10.1093/bioinformatics/btz848](https://doi.org/10.1093/bioinformatics/btz848) [Medline](#)
52. S. Frey, D. Görllich, A new set of highly efficient, tag-cleaving proteases for purifying recombinant proteins. *J. Chromatogr. A* **1337**, 95–105 (2014). [doi:10.1016/j.chroma.2014.02.029](https://doi.org/10.1016/j.chroma.2014.02.029) [Medline](#)
53. M. Martin, Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet. J.* **17**, 10–12 (2011). [doi:10.14806/ej.17.1.200](https://doi.org/10.14806/ej.17.1.200)
54. B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012). [doi:10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923) [Medline](#)
55. A. M. Sriramachandran, G. Petrosino, M. Méndez-Lago, A. J. Schäfer, L. S. Batista-Nascimento, N. Zilio, H. D. Ulrich, Genome-wide Nucleotide-Resolution Mapping of DNA Replication Patterns, Single-Strand Breaks, and Lesions by GLOE-Seq. *Mol. Cell* **78**, 975–985.e7 (2020). [doi:10.1016/j.molcel.2020.03.027](https://doi.org/10.1016/j.molcel.2020.03.027) [Medline](#)
56. K. Clement, H. Rees, M. C. Canver, J. M. Gehrke, R. Farouni, J. Y. Hsu, M. A. Cole, D. R. Liu, J. K. Joung, D. E. Bauer, L. Pinello, CRISPResso2 provides accurate and rapid genome editing sequence analysis. *Nat. Biotechnol.* **37**, 224–226 (2019). [doi:10.1038/s41587-019-0032-3](https://doi.org/10.1038/s41587-019-0032-3) [Medline](#)
57. A. Bankevich, S. Nurk, D. Antipov, A. A. Gurevich, M. Dvorkin, A. S. Kulikov, V. M. Lesin, S. I. Nikolenko, S. Pham, A. D. Pribelski, A. V. Pyshkin, A. V. Sirotkin, N. Vyahhi, G. Tesler, M. A. Alekseyev, P. A. Pevzner, SPAdes: A new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* **19**, 455–477 (2012). [doi:10.1089/cmb.2012.0021](https://doi.org/10.1089/cmb.2012.0021) [Medline](#)
58. M. Turmel, C. Otis, C. Lemieux, Divergent copies of the large inverted repeat in the chloroplast genomes of ulvophycean green algae. *Sci. Rep.* **7**, 994 (2017). [doi:10.1038/s41598-017-01144-1](https://doi.org/10.1038/s41598-017-01144-1) [Medline](#)
59. D. T. Hoang, O. Chernomor, A. von Haeseler, B. Q. Minh, L. S. Vinh, UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Mol. Biol. Evol.* **35**, 518–522 (2018). [doi:10.1093/molbev/msx281](https://doi.org/10.1093/molbev/msx281) [Medline](#)
60. P. J. Lockhart, A. W. Larkum, M. Steel, P. J. Waddell, D. Penny, Evolution of chlorophyll and bacteriochlorophyll: The problem of invariant sites in sequence analysis. *Proc. Natl. Acad. Sci. U.S.A.* **93**, 1930–1934 (1996). [doi:10.1073/pnas.93.5.1930](https://doi.org/10.1073/pnas.93.5.1930) [Medline](#)
61. A. Criscuolo, S. Gribaldo, BMGE (Block Mapping and Gathering with Entropy): A new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol. Biol.* **10**, 210 (2010). [doi:10.1186/1471-2148-10-210](https://doi.org/10.1186/1471-2148-10-210) [Medline](#)
62. J. Strecker, A. Ladha, Z. Gardner, J. L. Schmid-Burgk, K. S. Makarova, E. V. Koonin, F. Zhang, RNA-guided DNA insertion with CRISPR-associated transposases. *Science* **365**, 48–53 (2019). [doi:10.1126/science.aax9181](https://doi.org/10.1126/science.aax9181) [Medline](#)
63. S. E. Klompe, P. L. H. Vo, T. S. Halpin-Healy, S. H. Sternberg, Transposon-encoded CRISPR-Cas systems direct RNA-guided DNA integration. *Nature* **571**, 219–225 (2019). [doi:10.1038/s41586-019-1323-z](https://doi.org/10.1038/s41586-019-1323-z) [Medline](#)
64. M. Krupovic, K. S. Makarova, P. Forterre, D. Prangishvili, E. V. Koonin, Casposons: A new superfamily of self-synthesizing DNA transposons at the origin of prokaryotic CRISPR-Cas immunity. *BMC Biol.* **12**, 36 (2014). [doi:10.1186/1741-7007-12-36](https://doi.org/10.1186/1741-7007-12-36) [Medline](#)
65. S. M. Crotty, B. Q. Minh, N. G. Bean, B. R. Holland, J. Tuke, L. S. Jermini, A. V. Haeseler, GHOST: Recovering Historical Signal from Heterotachously Evolved Sequence Alignments. *Syst. Biol.* **69**, 249–264 (2020). [Medline](#)

#### ACKNOWLEDGMENTS

We thank J. Strecker, S. Hirano, and D. Strebing for advice regarding biochemistry experiments, G. Faure for advice regarding computational analyses, and all members of the Zhang lab for helpful discussions. We are grateful to the following individuals for generously providing access to their metagenomic data (IMG accessions provided in parenthesis): B. Campbell (IMG3300025818 and IMG3300007960), A. Buchan (IMG3300017968), and E. Edwards (IMG3300020812 and IMG3300023203). We appreciate assistance in DNA extraction and troubleshooting from M. Forbes for the Yellowstone Lake metagenomes. Yellowstone Lake samples were collected with support from the National Park Service–Yellowstone National Park (Research Permit YELL-2016/17-SCI-7018). **Funding:** Supported by NSF Integrated Earth Systems grant subaward A101357 (L.M. and W.I.); NSF Division of Environmental Biology grant 1950770 (M.D. and W.I.); Department of Energy–Joint Genome Institute grant CSP 1675 (W.I.); the National Library of Medicine (K.M.S. and E.V.K.); and NIH grants 1R01-HG009761 and 1DP1-HL141201, the Howard Hughes Medical Institute, the Open Philanthropy Project, the Harold G. and Leila Mathers Foundation, the Edward Mallinckrodt Jr. Foundation, the Poitras Center for Psychiatric Disorders Research at MIT, the Hock E. Tan and K. Lisa Yang Center for Autism Research at MIT, the Yang-Tan Center for Molecular Therapeutics at MIT, and the Phillips family, R. Metcalfe, and J. and P. Poitras (F.Z.). **Author contributions:** H.A.-T., S.K., and F.Z. conceived of the project. H.A.-T., S.K., E.D., R.O., S.P.N., K.S.M., E.V.K., and F.Z. designed and performed experiments. L.M., M.D., and W.I. collected metagenomic data. F.Z. supervised the research and experimental design with support from R.M. H.A.-T., S.K., R.M., E.V.K., and F.Z. wrote the manuscript with input from all authors. **Competing interests:** H.A.-T., S.K., E.D., S.P.N., and F.Z. are co-inventors on U.S. provisional patent applications filed by the Broad Institute related to this work. F.Z. is a cofounder of Editas Medicine, Beam Therapeutics, Pairwise Plants, Arbor Biotechnologies, and Sherlock Biosciences. **Data and materials availability:** Sequences of genes used in the experimental studies are available via online sequence repositories and expression plasmids listed in table S1 are available from Addgene under a material transfer agreement with the Broad Institute. Raw reads from microbial small RNA-seq are available on SRA under BioProject PRJNA744508. Scripts for data analysis and visualization are available at Zenodo (24). Additional information available via the Zhang Lab website (<https://zhanglab.bio>).

#### SUPPLEMENTARY MATERIALS

<https://science.org/doi/10.1126/science.abj6856>

Materials and Methods

Supplementary Text

Figs. S1 to S36

Tables S1 to S6

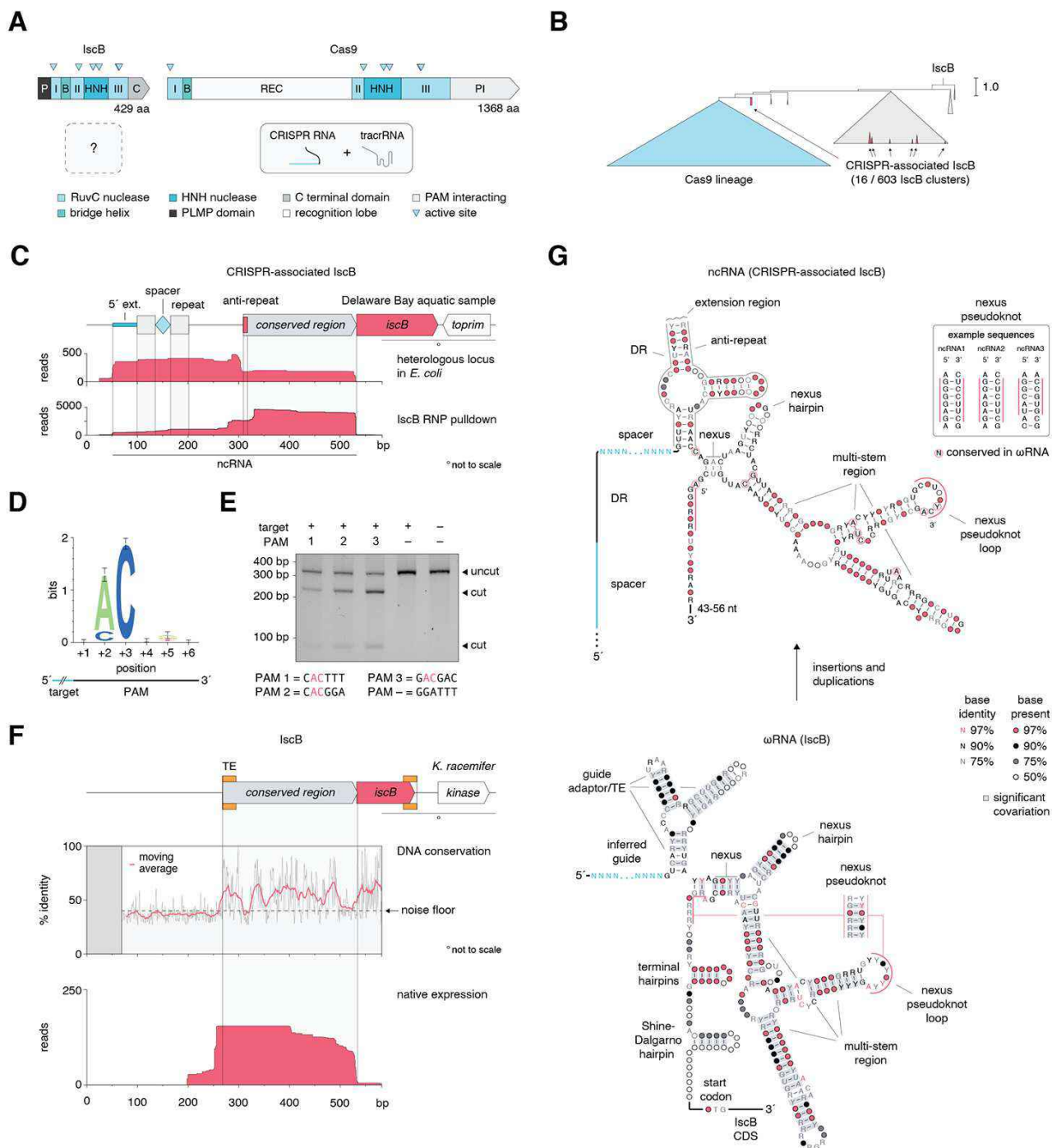
Data S1 to S4

References (25–65)

26 May 2021; accepted 9 August 2021

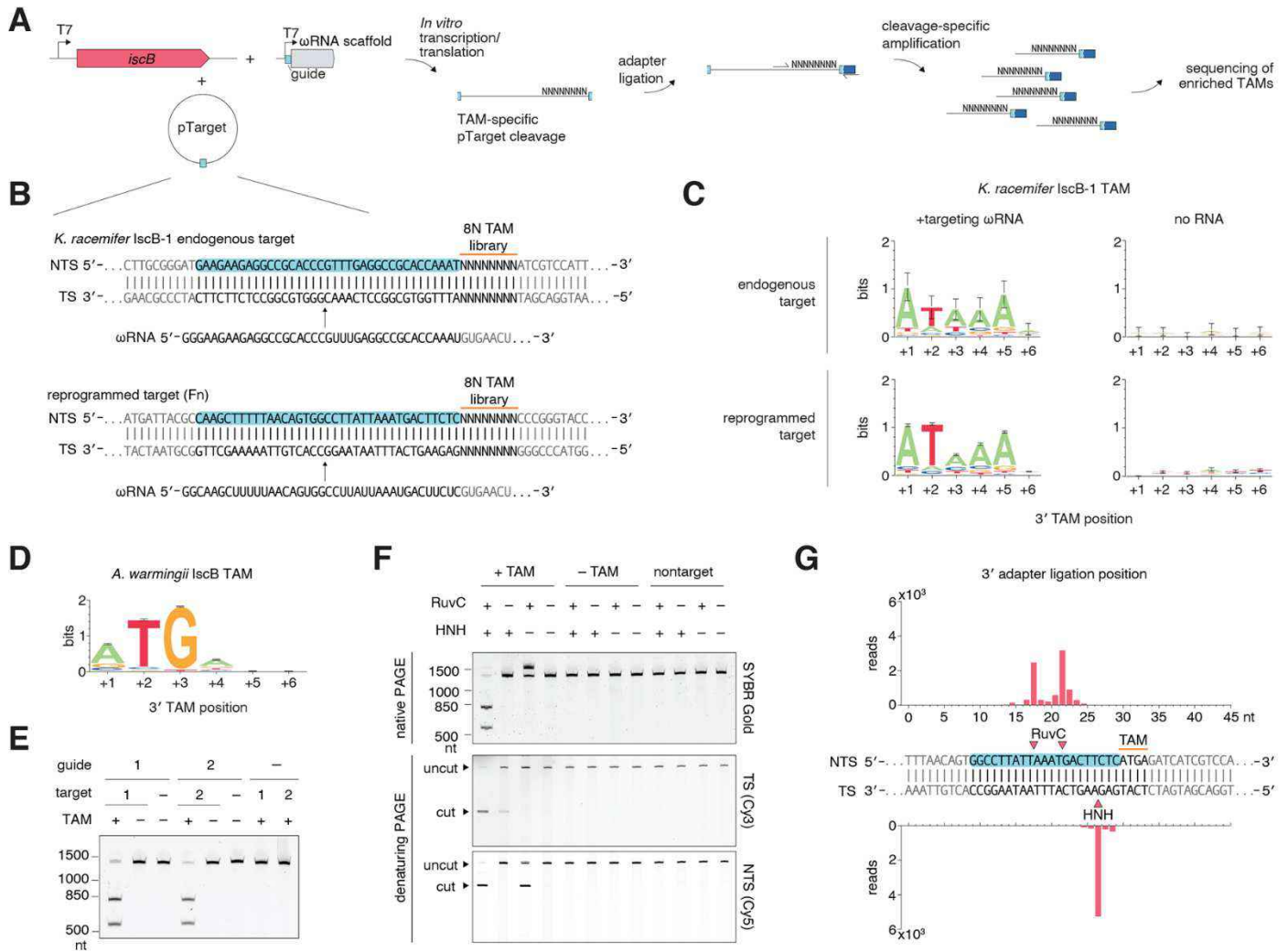
Published online 9 September 2021

10.1126/science.abj6856

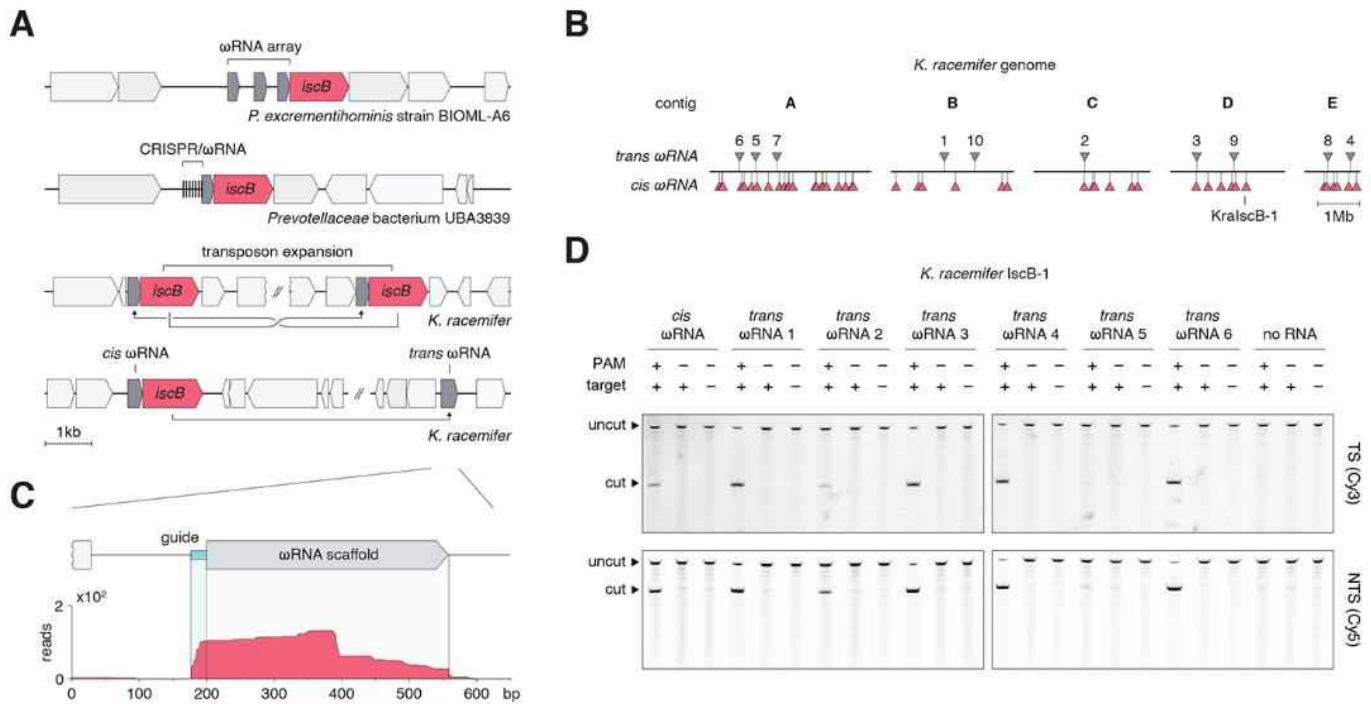


**Fig. 1. IscBs are associated with ncRNAs of unknown function.** (A) Comparison of IscB and Cas9 domains and previously described ncRNAs. (B) Phylogenetic analysis of the RuvC, BH, and HNH domains of Cas9 and IscB clusters using IQ-Tree 2. Genomic association shows 16/603 IscB clusters have strong association to CRISPR, occurring independently in multiple clades. (C) Small RNA-seq of a heterologously expressed CRISPR-associated IscB locus (top) and RNP pulldown (bottom). (D) Sequence logo for the PAM as determined by a plasmid depletion assay. (E) In vitro cleavage by IscB-single guide RNA RNP complex. (F) (Top) Conservation analysis of regions upstream of  $N=563$  non-redundant IscB loci. (Bottom) Small RNA-seq of an IscB locus in *K. racemifer* strain SOSPI-21. (G) Secondary structure predictions of CRISPR-associated IscB ncRNA and IscB  $\omega$ RNA. Guiding function of  $\omega$ RNAs was inferred by comparison of the two structures. TE: transposon end.

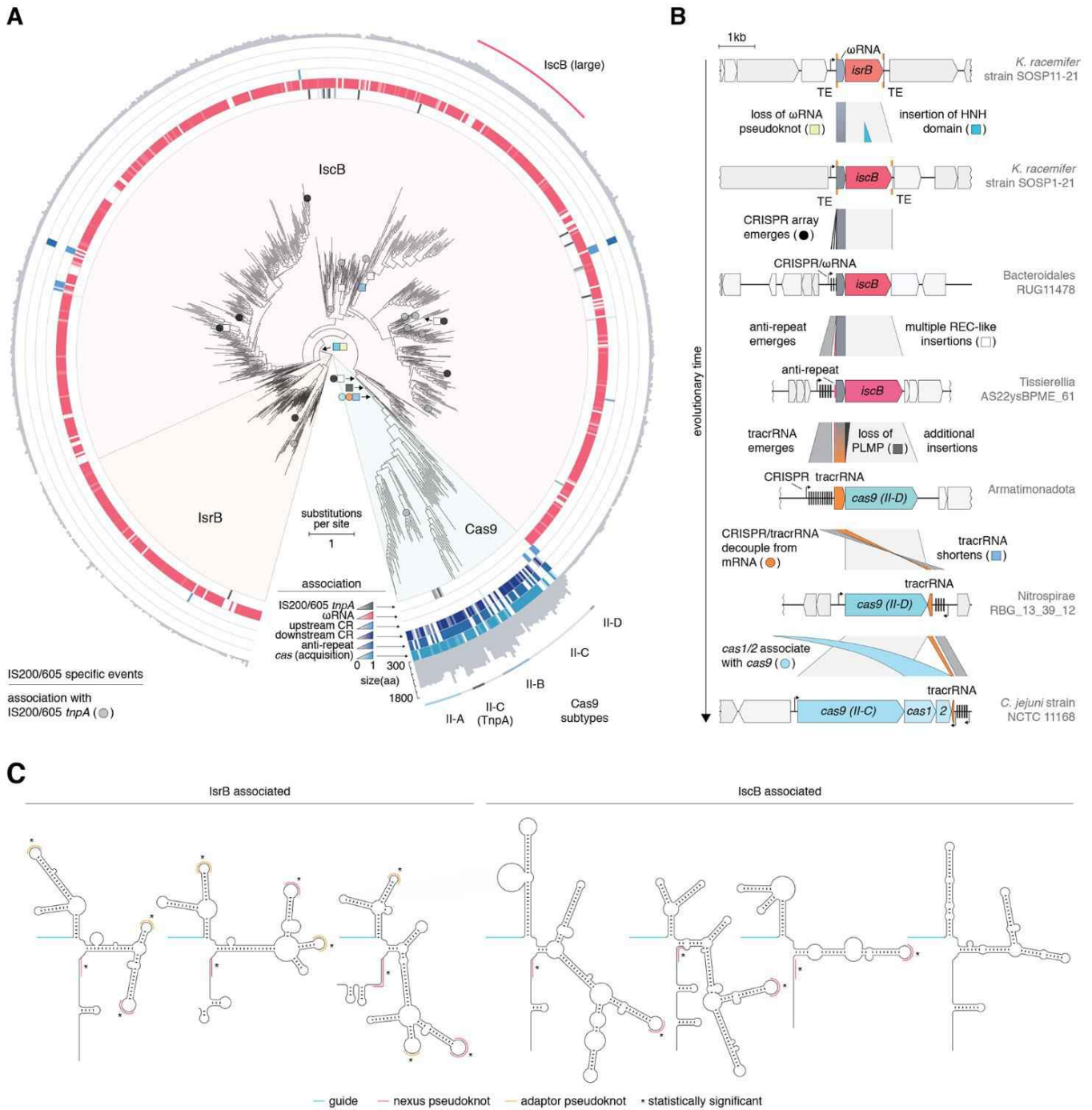




**Fig. 2. IscB is an RNA-guided DNA endonuclease.** (A) Design of an IVTT-based TAM screen. (B) *Kra*IscB-1 endogenous target and reprogrammed target sequences used in IVTT TAM screens. (C) dsDNA cleavage by *Kra*IscB-1 and ωRNA targeting sequence flanked by ATAAA 3' TAM. (D) dsDNA cleavage by *Awa*IscB and ωRNA targeting sequence flanked by ATGA 3' TAM. (E) In vitro-reconstituted *Awa*IscB-ωRNA RNP cleavage of dsDNA substrates in the presence or absence of a target and/or TAM. TS: target strand; NTS: non-target strand; nt: nucleotides. (F) In vitro dsDNA cleavage of *Awa*IscB with selectively inactivated nuclease domains. (G) Sequencing of cleavage products generated by *Awa*IscB.



**Fig. 3. Guide-encoding mechanisms of IscB.** (A) Example loci for each major mechanism of encoding multiple guides: entire  $\omega$ RNA arrays associate with IscB,  $\omega$ RNAs duplicate or insert into CRISPRs, transposon expansion results in multiple nearly identical loci that each express different guides, and standalone *trans*-acting  $\omega$ RNAs form independently of adjacent IscBs. (B) *K. racemifer* encodes 48 IscB loci with *cis*  $\omega$ RNAs and 10 standalone *trans*-acting  $\omega$ RNAs. (C) Small RNA-seq of a standalone  $\omega$ RNA locus in *K. racemifer*. (D) KraliscB-1, in complex with *cis* or *trans*  $\omega$ RNAs with the same guide sequence, mediate cleavage of dsDNA in a TAM- and target-dependent manner. Reactions were performed in IVTT using 5' strand-specific labeled linear targets. TS: Target strand; NTS: non-target strand. Contig accession and position information for all displayed loci are listed in table S6.



**Fig. 4. Diversity and evolution of IscB.** (A) Phylogenetic tree of IsrB, IscB, and Cas9. Associations with IS200/605 TnpA,  $\omega$ RNA, CRISPR arrays, anti-repeats (where applicable), and Cas acquisition genes. ORF size of cluster representative is shown on the second outermost ring. Notable groups are shown as colored arcs on the outermost ring. First occurrences of evolutionary events in each clade are marked by colored circles/squares, as described in (B). CR: CRISPR array. (B) Parsimonious evolutionary timeline linking IsrB to Cas9 with exemplifying loci. Colors of protein of interest indicate distinct stages in the evolution of IsrB to Cas9. (C) Structural diversity and evolution of  $\omega$ RNAs in IsrB and IscB systems.



Fig. 5 (preceding page). Exploration of the diversity of IS200/605 superfamily nucleases. (A) Evolution between IS200/605 transposon superfamily-encoded nucleases and associated RNAs. Dashed lines reflect tentative/unknown relationships. (B) Locations of *IscB* loci and fragments in the *I. tetrasporus* genome. Intact locus is labeled as "ChlorIscB." (C) Small RNA-seq of *I. tetrasporus*. (D) Weblogo of ChlorIscB cleavage TAM using a reprogrammed guide in an IVTT TAM screen. (E) Weblogo of OgeulscB TAM using a reprogrammed guide in an IVTT TAM screen. (F) Targeted OgeulscB-mediated indel formation at the *VEGFA* locus in HEK293FT cells ordered by abundance, with indel size on the left. (G) OgeulscB-mediated indel formation at multiple sites in HEK293T cells, \**P* < 0.05). (H) Small RNA-seq of  $\omega$ RNA from *IscB* locus in *K. racemifer* strain SOSP1-21. (I) Weblogo of *Desulfovibrio thermocuniculi* (*DthIscB*) TAM using a reprogrammed guide in an IVTT TAM screen. (J) *DthIscB* mediates  $\omega$ RNA-guided non-target strand nicking in a TAM- and target-dependent manner in an IVTT cleavage assay using 5' strand-specific labeled targets. (K) Small RNA-seq of  $\omega$ RNA from *TnpB* locus in *K. racemifer* strain SOSP1-21. (L) Comparison of  $\omega$ RNAs from *K. racemifer* *IscB* and *TnpB* loci. (M) Secondary structure prediction of *KraTnpB*-associated  $\omega$ RNA. (N) Weblogo of *A. macrosporangiidus* *TnpB* (*AmaTnpB*) TAM using a reprogrammed guide in an IVTT TAM screen. (O) In vitro-reconstituted *AmaTnpB* cleavage of dsDNA substrates in the presence or absence of  $\omega$ RNA, target, and/or TAM. (P) *AmaTnpB* performs  $\omega$ RNA-guided TAM-independent target-dependent cleavage of 3' Cy5.5-labeled ssDNA substrates. (Q) *AmaTnpB* cleaves a 3' Cy5.5-labeled collateral ssDNA substrate in the presence of TAM- and target-containing dsDNA or target-containing ssDNA substrates. Contig accession and position information for all displayed loci are listed in table S6.

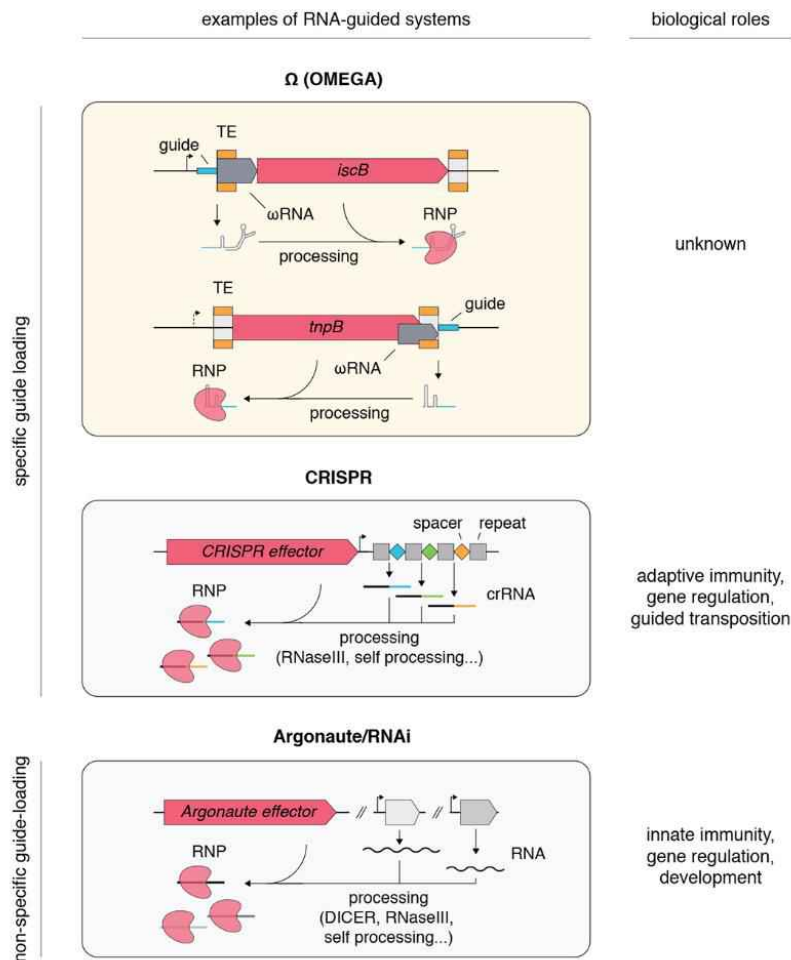


Fig. 6. Naturally occurring RNA-guided DNA-targeting systems. Comparison of  $\Omega$  (OMEGA) systems with other known RNA-guided systems. In contrast to CRISPR systems, which capture spacer sequences and store them in the locus within the CRISPR array,  $\Omega$  systems may transpose their loci (or *trans*-acting loci) into target sequences, converting targets into  $\omega$ RNA guides.

## The widespread IS200/605 transposon family encodes diverse programmable RNA-guided endonucleases

Han Altae-Tran, Soumya Kannan, F. Esra Demircioglu, Rachel Oshiro, Suchita P. Nety, Luke J. McKay, Mensur Dlaki#, William P. Inskeep, Kira S. Makarova, Rhiannon K. Macrae, Eugene V. Koonin, and Feng Zhang

*Science*, **Ahead of Print** • DOI: 10.1126/science.abj6856

### View the article online

<https://www.science.org/doi/10.1126/science.abj6856>

### Permissions

<https://www.science.org/help/reprints-and-permissions>