



Molecular mechanism and history of non-sense to sense evolution of antifreeze glycoprotein gene in northern gadids

Xuan Zhuang^{a,1,2}, Chun Yang^a, Katherine R. Murphy^a, and C.-H. Christina Cheng^{a,1}

^aSchool of Integrative Biology, University of Illinois at Urbana-Champaign, Urbana, IL 61801

Edited by David M. Hillis, The University of Texas at Austin, Austin, TX, and approved January 17, 2019 (received for review October 9, 2018)

A fundamental question in evolutionary biology is how genetic novelty arises. De novo gene birth is a recently recognized mechanism, but the evolutionary process and function of putative de novo genes remain largely obscure. With a clear life-saving function, the diverse antifreeze proteins of polar fishes are exemplary adaptive innovations and models for investigating new gene evolution. Here, we report clear evidence and a detailed molecular mechanism for the de novo formation of the northern gadid (codfish) antifreeze glycoprotein (AFGP) gene from a minimal non-coding sequence. We constructed genomic DNA libraries for AFGP-bearing and AFGP-lacking species across the gadid phylogeny and performed fine-scale comparative analyses of the AFGP genomic loci and homologs. We identified the noncoding founder region and a nine-nucleotide (9-nt) element therein that supplied the codons for one Thr-Ala-Ala unit from which the extant repetitive AFGP-coding sequence (c_{ds}) arose through tandem duplications. The latent signal peptide (SP)-coding exons were fortuitous non-coding DNA sequence immediately upstream of the 9-nt element, which, when spliced, supplied a typical secretory signal. Through a 1-nt frameshift mutation, these two parts formed a single read-through open reading frame (ORF). It became functionalized when a putative translocation event conferred the essential *cis* promoter for transcriptional initiation. We experimentally proved that all genic components of the extant gadid AFGP originated from entirely non-genic DNA. The gadid AFGP evolutionary process also represents a rare example of the proto-ORF model of de novo gene birth where a fully formed ORF existed before the regulatory element to activate transcription was acquired.

de novo gene | proto-ORF | adaptive evolution | noncoding origin | codfish AFGP

Evolutionary innovation of new genetic elements is recognized as a key contributor to organismal adaptation. For decades, the treatise of Ohno (1) shaped the paradigm of new gene creation in that it relies on the duplication of a preexisting protein gene. When subjected to selection, adaptive sequence changes in one copy may occur from which a gene with a novel function may emerge (1, 2). Creating new protein-coding genes de novo from noncoding DNA sequences was considered extremely rare. In recent years, however, examples of de novo genes have been reported in diverse animals and plants (see review ref. 3 and studies referenced therein). De novo gene births were generally deduced using a combination of phylogenetics and comparative genomic/transcriptomic analyses or the phylostratigraphy approach (4), which revealed evidence for lineage- or species-specific gene transcripts, whereas the orthologous sequences in sister species were nongenic. These revelations have spurred considerable interest and hypotheses of how de novo genes arise and evolve as well as questions regarding their functional importance (5, 6). Validating new genes identified from sequence-based comparisons is complicated by uncertainties around how comprehensive the genome assemblies and gene expression data are (7, 8). More challenging yet is identifying the selective pressures and molecular mechanisms

that created these putative new genes, and the adaptive functions and species fitness they may confer.

In contrast, antifreeze protein genes of polar teleost fishes are unequivocal new genes that confer a clear life-saving function and fitness benefit. The selective pressure that compelled their evolution is also clear. They evolved in direct response to polar marine glaciations, preventing death of fish from inoculative freezing by environmental ice crystals in subzero waters (9, 10). Such a strong life-or-death selective pressure has driven the independent evolution of multiple structurally distinct types of antifreezes: antifreeze peptide (AFP) types I, II, and III and antifreeze glycoprotein (AFGP) in diverse fish lineages where they perform the same ice-growth inhibition function (10). The structural differences lie in their distinct genetic ancestry. Thus, fish antifreezes as a group can richly inform on the diversity of molecular origins and evolutionary mechanisms that produced a vital function.

The well-known mechanism of evolution by gene duplication from a preexisting ancestor as diverse as C-type lectin and sialic acid synthase followed by sequence tinkering by natural selection produced AFP II (11) and AFP III (12), respectively. AFGPs

Significance

The diverse antifreeze proteins enabling the survival of different polar fishes in freezing seas offer unparalleled vistas into the breadth of genetic sources and mechanisms that produce crucial new functions. Although most new genes evolved from preexisting genic ancestors, some are deemed to have arisen from noncoding DNA. However, the pertinent mechanisms, functions, and selective forces remain uncertain. Our paper presents clear evidence that the antifreeze glycoprotein gene of the northern codfish originated from a noncoding region. We further describe the detailed mechanism of its evolutionary transformation into a full-fledged crucial life-saving gene. This paper is a concrete dissection of the process of a de novo gene birth that has conferred a vital adaptive function directly linked to natural selection.

Author contributions: X.Z. and C.-H.C.C. designed research; X.Z., C.Y., K.R.M., and C.-H.C.C. performed research; X.Z. and C.Y. analyzed data; and X.Z. and C.-H.C.C. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

This open access article is distributed under [Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 \(CC BY-NC-ND\)](https://creativecommons.org/licenses/by-nc-nd/4.0/).

Data deposition: The sequences reported in this paper have been deposited in the GenBank database (accession nos. [MK011258](https://accession.csiro.au/MK011258)–[MK011272](https://accession.csiro.au/MK011272), [MH992395](https://accession.csiro.au/MH992395)–[MH992397](https://accession.csiro.au/MH992397), and [MK011291](https://accession.csiro.au/MK011291)–[MK011308](https://accession.csiro.au/MK011308)).

¹To whom correspondence may be addressed. Email: zhuangxuan@gmail.com or c-cheng@illinois.edu.

²Present address: Department of Ecology and Evolution, The University of Chicago, Chicago, IL 60637.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1817138116/-DCSupplemental.

Published online February 14, 2019.

have evolved independently in two unrelated fish lineages at opposite poles: the Antarctic notothenioid fishes (Notothenioidei) and the Arctic/northern codfishes (Gadidae), providing a striking example of protein sequence convergence (9, 13). In both lineages, AFGPs occur as a family of size isoforms composed of varying numbers of repeats of a basic tripeptide unit (Thr-Ala-Ala) with each Thr glycosylated with a disaccharide (10, 14). They are encoded by a family of polyprotein genes, each of which produces a large polyprotein precursor consisting of many tandemly linked AFGP molecules that are then post-translationally cleaved to yield mature AFGPs (13, 15). The Antarctic notothenioid *AFGP* evolved through a more innovative process than gene duplication and sequence divergence. It exemplifies partial de novo gene evolution. A functionally unrelated ancestral trypsinogenlike protease (TLP) gene provided the secretory signal and a 3' untranslated sequence of the incipient *AFGP*. The large repetitive AFGP polyprotein-coding region was generated de novo from duplications of a partly non-sense 9-nt sequence that straddled an intron-exon junction in the *TLP*, which happened to comprise the three codons for one Thr-Ala-Ala unit (15, 16).

Where from and how the northern gadid *AFGP* evolved have remained lasting enigmas. Despite voluminous collections of genes and genome sequences available in databases, there are no meaningful homologs to any part of the gadid *AFGP* to hint at ancestry. This peculiar absence of related genes suggests that the gadid *AFGP* gene may have originated from nonprotein-coding DNA. Gadid *AFGP* presumably evolved very recently, in response to the cyclic northern hemisphere glaciation that commenced in the late Pliocene about 3 Mya. We reason it is unlikely that mutational processes could completely obscure even noncoding sequences within such a short evolutionary time such that the extant form of the *AFGP* nongenic ancestor should remain identifiable. We, therefore, decided to track the AFGP genotype and its homologs within the gadid phylogeny to pinpoint the ancestral DNA site of origin and reconstruct the gadid *AFGP* evolutionary path. Here, we report the identification of the noncoding founder sequence and the mechanism by which it gave rise to a new functional gadid AFGP gene. Our results also show that the gadid *AFGP* evolutionary process likely represents a rare example of the proto-ORF model of de novo gene birth (6, 17) where the noncoding founder ORF existed well before the novel gene arose.

Results and Discussion

At the minimum, de novo formation of a functional protein gene requires the acquisition of an ORF encoding the new protein and the basic *cis*-regulatory elements to activate its transcription and translation. AFGPs are secreted plasma proteins, thus, a signal peptide (SP) will also be needed to instruct cellular export of AFGP molecules into the blood circulation. Thus, to reconstruct the formation of the gadid AFGP gene requires elucidating how these essential genic components were generated and became properly linked into a functional whole gene. We began with precise delineation of these components that make up the structure of functional AFGPs in AFGP-bearing gadids. We then juxtaposed them against the structures of the *AFGP* homologs in more basal non-AFGP-bearing species representing progressively more ancestral states. This enabled us to decipher the essential molecular steps and timing in the de novo formation of the AFGP gene in the gadid lineage.

Phylogenetic Context for the Selected Gadid Species. The phylogenetic tree in Fig. 1 (detailed in *SI Appendix, Fig. S1*) depicts the relationships of the northern cod species used in this paper. We characterized the AFGP genes or noncoding homologs of seven species (gene structures to the right of the tree, Fig. 1), which were chosen for the strategic positions they occupy in the gadid tree. The monophyletic Gadidae family [*sensu* (18)] includes both AFGP-bearing and non-AFGP-bearing species; the former occurs in two subclades within the subfamily Gadinae (Fig. 1).

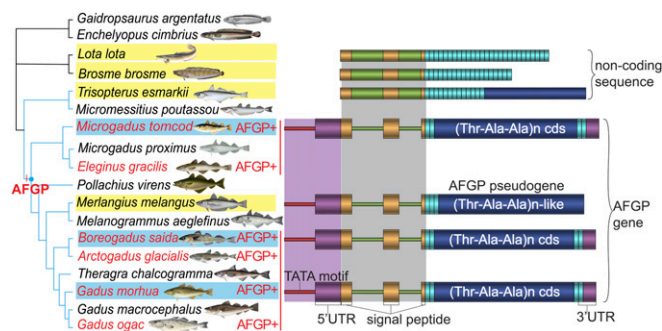


Fig. 1. Gadid phylogeny and AFGP gene/homolog structures. The phylogenetic tree of Gadidae is a congruent cladogram derived from Bayesian and maximum likelihood trees using complete ND2 gene sequences (*SI Appendix, Fig. S1*). Light blue branches indicate lineages of the Gadinae subfamily. The two gadine subclades containing AFGP-bearing species (red vertical bars), their most recent common ancestor (blue dot), and the emergence of the AFGP trait are as indicated. The three AFGP-bearing species (AFGP+) and four AFGP-lacking species analyzed in this paper are shaded in blue and yellow, respectively. The structure of their AFGP gene or nongenic homolog is shown to the right. Gray and purple shaded areas indicate homologous regions. Cyan segments are sequence repeats. The dark blue segment is a repetitive AFGP cds or AFGP-like sequence.

The selected AFGP-bearing *Boreogadus saida* (polar cod) (19) and *Gadus morhua* (Atlantic cod) (20) represent one gadine subclade, and *Microgadus tomcod* (Atlantic tomcod) (21) represents the other. The four AFGP-lacking gadids were chosen for their evolutionary distances from the AFGP bearers. Two of them are gadines; *Merlangius merlangus* (whiting) nests within the AFGP-bearing subclade containing *B. saida* and *G. morhua* and thus shares the last common ancestor with all AFGP-bearing species (Fig. 1, blue dot), whereas *Trisopterus esmarkii* (Norway pout) is basal to the two AFGP-bearing subclades. The other two, *Brosme brosme* (cusk) and the freshwater *Lota lota* (burbot) belong to the subfamily Lotinae and are basal species that serve as ancestral proxies before the AFGP trait emerged (Fig. 1).

Gadid AFGP Genomic Regions and AFGP Gene Structure. We isolated AFGP-positive large-insert genomic DNA clones from the respective Bacterial Artificial Chromosome (BAC) library of *B. saida* and *M. tomcod* by screening with a probe specific to the AFGP (Thr-Ala-Ala)_n cds and sequenced the minimal tiling path clones spanning the *AFGP* genomic region. For *G. morhua*, a draft genome was available, but the *AFGP* genomic region was incomplete (22). We bioinformatically deduced the pertinent BAC clones (*SI Appendix, Fig. S2*) and obtained them (available from the vendor) for sequencing. The reconstructed genomic regions contained 12 functional and four pseudo AFGP genes in *B. saida*, five and two in *G. morhua*, and three and one in *M. tomcod*, spanning ~510, 190, and 80 kbp, respectively, in the three species (*SI Appendix, Fig. S3*).

To determine the gene structure of functional AFGPs we used supporting transcript sequences obtained by 5' rapid amplification of cDNA ends (RACE) (*SI Appendix, Fig. S4*). A functional AFGP consists of three exons and two introns (Fig. 2). The first two small exons (E1 and E2) and the first two nts of the large third exon (E3) encode the SP, and the rest of E3 encodes a short pro-peptide and the long AFGP polyprotein. These demarcations differ from the only known full-length gadid *AFGP* gene sequence to date (13). In that sequence, the predicted SP contained many atypical hydrophilic residues, indicating inaccurate assignment of splice junctions and reading frames, hence, the need for reassessment. In this paper, the predicted SP contained the requisite stretch of hydrophobic residues followed by a putative cleavage site with high prediction scores (*SI Appendix, Fig. S5*), which is characteristic of a secretory signal. It is conserved in all functional AFGPs

existed in the founder genomic site and became functionalized when the promoter region was acquired.

Nongenic Origin of SP and Promoter Region. We experimentally verified that the AFGP SP cds and promoter sequence did not originate from any existing protein-coding genes in the gadid genome. We hybridized the genomic BAC library macroarrays of *B. saida* and *M. tomcod* with the AFGP 5' probe that is specific to this region. The hybridized clones were exactly the same clones that hybridized to the (Thr-Ala-Ala)_n cds probe (SI Appendix, Fig. S10). This strongly supports that no homologs of the SP, 5' UTR, and promoter regions of AFGP exist outside of the AFGP genomic loci. Thus, the promoter and SP cds of functional AFGP also originated de novo, unassociated with any preexisting protein gene.

Further Verifications of Nongenic Origin of AFGP. Recently evolved genes and the extant homologs of their genetic ancestor often remain as near neighbors in the genome. For example, the AFGP gene family of the Antarctic notothenioids closely clusters with its ancestral homologs: the trypsinogenlike protease genes along with the broader trypsin gene family within an ~400 kbp region (27). In contrast, we found none of the neighboring genes (e.g., *MAK16* and *RAB14*) in the AFGP genomic regions of the three AFGP-bearing gadids *B. saida*, *G. morhua*, and *M. tomcod* share any sequence similarity with AFGPs (SI Appendix, Fig. S3), and, thus, they are evolutionarily unrelated to AFGP. The absence of a potential protein gene ancestor nearby is consistent with gadid AFGP having evolved de novo.

We further reasoned that an absence of transcription of the AFGP homologs in the AFGP-lacking gadids would provide compelling support that they are nonfunctional or nongenic DNA. Thus, we performed Northern blot hybridizations of RNA from pancreatic tissue [the site of AFGP synthesis (28)] of the four AFGP-lacking species using their respective species-specific AFGP-homolog sequence as probes and included *B. saida* for comparison. No transcripts of AFGP homologs were detectable in any of the three AFGP-lacking gadids. Only *B. saida* pancreatic RNA showed hybridization with strong intensity to its own AFGP cds probe and in varying intensity to the AFGP homolog probes from the other species due to various degrees of nt sequence identity (SI Appendix, Fig. S11). Since *L. lota*, *B. brosme*, and *T. esmarkii* are basal to the AFGP-bearing clade, their AFGP homologs must represent the ancestral transcriptionally inactive noncoding form. The AFGP-lacking *M. merlangus* is nested within the AFGP-bearing clade, and its AFGP homolog most closely resembles a functional AFGP except for inactivating mutations in the (Thr-Ala-Ala)_n cds. Thus, it represents a subsequent non-functionalization into a nontranscribed pseudogene after the emergence of AFGP in the common ancestor of the AFGP-bearing clade. The loss of function relates to the nonfreezing water (Tromsø fjord in this study) *M. merlangus* inhabits today where antifreeze protection is not needed.

Gadid AFGP Evolved from Entirely Nongenic DNA. Fig. 4 summarizes the forgoing deductions on the noncoding origins of the essential AFGP sequence components and the possible molecular steps in the evolutionary transformation of these components into a complete new functional AFGP. The AFGP founder structure (Fig. 4A) existed in the gadid ancestor as a short noncoding genomic sequence comprising a segment (~240 nt) with latent-coding exons (bronze segments) that have the potential to form a peptide sequence with properties for a secretory signal. The adjoining 27-nt GCA(Ala)-rich sequence (cyan segment) contained multiple nested 9-nt elements, any of which could become the three codons for the AFGP tripeptide (Thr-Ala-Ala) building block through a 1-nt substitution. Chance duplications of this ancestral 27-nt GCA-rich sequence produced four tandem copies (Fig. 4B). One of the 9-nt AFGP tripeptide-coding elements in the midst of

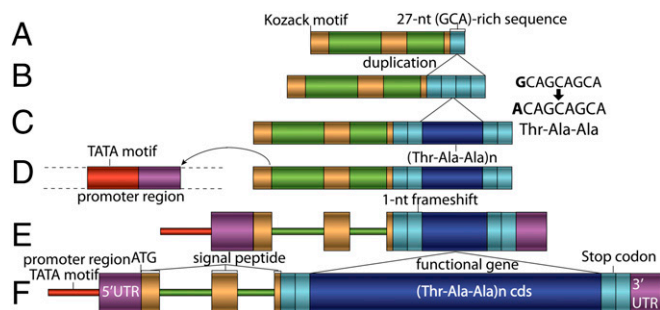


Fig. 4. Evolutionary mechanism of the gadid AFGP gene from noncoding DNA. The color codes of the sequence components follow Fig. 1. (A) The ancestral noncoding DNA contained latent signal peptide-coding exons with a 5' Kozack motif, adjacent to a duplication-prone 27-nt GCA-rich sequence. (B) The 27-nt GCA(Ala)-rich sequence duplicated forming four tandem copies. (C) A 9-nt in the midst of the four 27-nt duplicates became the three codons for one AFGP Thr-Ala-Ala unit and underwent microsatellitellike duplication forming a proto-ORF. (D) A proximal upstream regulatory region acquired through a putative translocation event. (E) A 1-nt frameshift led to a contiguous SP, a propeptide, and a Thr-Ala-Ala-like cds in a read-through ORF. (F) Intragenic (Thr-Ala-Ala)_n cds amplification, fulfilling the antifreeze function under natural selection.

the four copies likely underwent microsatellitellike duplications producing a budding ORF for the repetitive AFGP tripeptide cds, which began spreading the two pairs of 27-nt GCA-rich duplicates apart to the flanking positions (Fig. 4C). A putative translocation event in the last common ancestor of AFGP-bearing gadids moved the hitherto unexpressed AFGP precursor to a new genomic location that fortuitously contained a TATA motif thereby enabled transcription (Fig. 4D). Concurrently or subsequently, a 1-nt frameshift deletion in the second 5' 27-nt duplicate likely occurred and served to link the latent cds for the SP and the downstream AFGP (Thr-Ala-Ala)_n repeats in a single read-through ORF. Expression and secretion of the nascent antifreeze protein became possible (Fig. 4E). The smallest (and often the most abundant) functional AFGP isoform (AFGP8) comprises only four tripeptide repeats (10), which could be achieved through only two tandem duplications. The fledgling antifreezing protection could, therefore, augment fitness in the individual at the onset of northern hemisphere marine glaciation. Subsequent intensification of environmental selection pressures likely drove the intragenic (Thr-Ala-Ala)_n cds expansion forming large AFGP polyprotein genes (Fig. 4F) as well as additional whole gene duplications. The result manifests in the multigene family of AFGP polyproteins (SI Appendix, Fig. S3) and the robust antifreeze activities the AFGP-bearing gadids possess today (10, 13).

Proto-ORF Model of de Novo Evolution of Gadid AFGP. The deduced evolutionary process of the gadid AFGP gene from non-sense DNA adds valuable insights into how adaptive functional genes could arise “from scratch.” The birth of de novo genes involves two fundamental events: the formation of an ORF and the acquisition of regulatory signals for transcription. In principle, these events could occur in either order. This prompted two major competing models: the protogene versus the proto-ORF model (3, 17, 29–31). The occurrence of a protogene is generally easier to detect as the de novo gene has a noncoding ortholog with demonstrable transcripts in the out-group species. Thus, the model has found ample support in studies that showed transcription preceded the emergence of an ORF (6, 17, 30, 32). The proto-ORF model states that an ORF was present before regulatory signals for expression were acquired. The existence of proto-ORFs is challenging to prove as they likely accumulate mutations that would interrupt the ORF before they could become transcribed for selection to act upon (3, 29, 30, 33). The history and mechanism of gadid AFGP evolution deduced in this paper (Fig. 4)

fits the proto-ORF model. This is because the (Thr-Ala-Ala)_n-like repeats and SP cds were formed in the basal lineage-lacking AFGP (represented by *T. esmarkii*) before the regulatory signal for transcription appeared in the more derived gadids in the AFGP-bearing clade (Fig. 1). Thus, we suggest that the recently evolved gadid AFGP serves as a clear and rare supporting example of the proto-ORF model of de novo gene birth.

Although the emergence of de novo genes has been well documented, the selective pressure and functional necessity that compelled their birth remain largely unknown (32). Most de novo genes are deemed unlikely to gain or retain function before their genelike properties decay (5, 34) unless a timely major shift in the fitness landscape allows them to be sufficiently useful for selection to take hold. The de novo gadid AFGP is a rare example where the affecting fitness shift is clear. Its emergence correlated with strong risk-of-death selection pressures from environmental changes in the form of plunging ocean temperatures and formation of ice in the water column during the Pliocene/Pleistocene northern hemisphere glaciation. The gene multiplied in gadid species that remained in frigid habitats (e.g., *B. saida*, *M. tomcod*, and *G. morhua*) but gradually decayed in species that no longer experience the threat of freezing (e.g., *M. merlangus*).

Conclusion

We have characterized the evolutionary process and the details of the underlying molecular mechanisms through which all of the essential genic components of the northern cod AFGP gene could have developed from noncoding DNA and the union of these emerging coding parts into a new functional whole gene. We provide evidence that latent-coding components existed before the acquisition of the necessary *cis*-regulatory region for transcriptional activation. Thus, the gadid AFGP evolutionary history is a rare example supporting the proto-ORF hypothesis of de novo gene birth. With this paper, we fully resolved the lasting question of how two unrelated groups of fish at opposite poles: the Antarctic notothenioid fishes and the northern codfishes, invented a near-identical AFGP. The notothenioid AFGP evolved

within the structural framework of a preexisting gene ancestor but constructed a new AFGP cds from de novo expansion of a rudimentary partly non-sense tripeptide-coding element. Northern gadid displayed even greater evolutionary ingenuity, constructing all parts of a functional AFGP gene entirely from noncoding DNA.

Materials and Methods

Detailed materials and methods are given in the *SI Appendix, Materials and Methods*. Briefly, we constructed large genomic DNA-insert BAC libraries for two AFGP-bearing gadids *B. saida* and *M. tomcod* and the AFGP-lacking basal *B. brosmie* and smaller-insert phage libraries for three other AFGP-lacking species *M. merlangus*, *T. esmarkii*, and *L. lota*. The libraries were screened with radiolabeled probes derived from (Thr-Ala-Ala)_n cds or the 5' sequence of the AFGP gene to isolate clones containing AFGP or AFGP homologs, respectively. For the AFGP-bearing *G. morhua*, we deduced the AFGP-positive BAC clones from published genome data and obtained them from a commercial vendor. The relevant positive clones or clone fragments were sequenced using various sequencing strategies, and the assembled sequences were analyzed. To correctly determine the gene structure of the functional AFGP, we obtained 5' RACE map intron-exon junctions. To verify that the SP and promoter region evolved de novo, we rescreened the *B. saida* and *M. tomcod* BAC libraries with probes specific to this region to detect whether they hybridized elsewhere outside of the AFGP loci. We conducted Northern blot hybridizations to test for mRNA expression of AFGP homologs using species-specific probes to verify the hypothesis that they are untranscribed noncoding DNA. Fish collection and sampling followed University of Illinois at Urbana-Champaign institutional approved protocol as described in *SI Appendix*. All sequence data have been deposited in National Center for Biotechnology Information.

ACKNOWLEDGMENTS. We sincerely thank our colleagues Kim Praebel, Svein-Erik Fevolden, Arthur DeVries, Howard Reisman, Kevin Bilyk, Shannon Zellerhoff, as well as the Cornell University Biological Field Station for their kind assistance in collecting the gadid species in this study. We thank Jørgen Christiansen for the opportunity to participate in the TUNU cruises on the R/V Helmer Hanssen to the Svalbard and East Greenland coasts to collect high Arctic species. We also thank Dr. Chris Amemiya and his previous lab member Andrew Stuart for their insightful advice on the BAC library construction and for making the lab facility available for our use. Special thanks go to Melody Clark, Lloyd Peck, Konrad Meister, and Arthur DeVries for their help in editing the paper. This work was supported by the US National Science Foundation Grant DEB 0919496 (to C.-H.C.C.).

- Ohno S (1970) *Evolution by Gene Duplication* (George Allen & Unwin Ltd., London; Springer, New York).
- Jacob F (1977) Evolution and tinkering. *Science* 196:1161–1166.
- McLysaght A, Guerzoni D (2015) New genes from non-coding sequence: The role of de novo protein-coding genes in eukaryotic evolutionary innovation. *Phil Trans R Soc B* 370: 20140332.
- Domazet-Lošo T, Brajković J, Tautz D (2007) A phylostratigraphy approach to uncover the genomic history of major adaptations in metazoan lineages. *Trends Genet* 23:533–539.
- Tautz D, Domazet-Lošo T (2011) The evolutionary origin of orphan genes. *Nat Rev Genet* 12:692–702.
- McLysaght A, Hurst LD (2016) Open questions in the study of de novo genes: What, how and why. *Nat Rev Genet* 17:567–578.
- Guerzoni D, McLysaght A (2016) De novo genes arise at a slow but steady rate along the primate lineage and have been subject to incomplete lineage sorting. *Genome Biol Evol* 8:1222–1232.
- Moyers BA, Zhang J (2015) Phylostratigraphic bias creates spurious patterns of genome evolution. *Mol Biol Evol* 32:258–267.
- Cheng C-HC (1998a) Evolution of the diverse antifreeze proteins. *Curr Opin Genet Dev* 8:715–720.
- DeVries AL, Cheng C-HC (2005) Antifreeze proteins and organismal freezing avoidance in polar fishes. *The Physiology of Polar Fishes*, eds Farrell AP, Steffensen JF (Elsevier Academic Press, San Diego), Vol 22, pp 155–201.
- Liu Y, et al. (2007) Structure and evolutionary origin of Ca²⁺-dependent herring type II antifreeze protein. *PLoS One* 2:e548.
- Deng C, Cheng C-HC, Ye H, He X, Chen L (2010) Evolution of an antifreeze protein by neofunctionalization under escape from adaptive conflict. *Proc Natl Acad Sci USA* 107:21593–21598.
- Chen L, DeVries AL, Cheng C-HC (1997a) Convergent evolution of antifreeze glycoproteins in Antarctic notothenioid fish and Arctic cod. *Proc Natl Acad Sci USA* 94:3817–3822.
- DeVries AL (1971) Glycoproteins as biological antifreeze agents in antarctic fishes. *Science* 172:1152–1155.
- Chen L, DeVries AL, Cheng C-HC (1997b) Evolution of antifreeze glycoprotein gene from a trypsinogen gene in Antarctic notothenioid fish. *Proc Natl Acad Sci USA* 94:3811–3816.
- Cheng C-HC, Chen L (1999) Evolution of an antifreeze glycoprotein. *Nature* 401:443–444.
- Schlötterer C (2015) Genes from scratch—The evolutionary fate of de novo genes. *Trends Genet* 31:215–219.
- Teletchea F, Laudet V, Hänni C (2006) Phylogeny of the Gadidae (sensu Svetovidov, 1948) based on their morphology and two mitochondrial genes. *Mol Phylogenet Evol* 38:189–199.
- Osuga DT, Feeney RE (1978) Antifreeze glycoproteins from Arctic fish. *J Biol Chem* 253:5338–5343.
- Hew CL, Slaughterd D, Fletcher GL, Joshi SB (1981) Antifreeze glycoproteins in the plasma of Newfoundland Atlantic cod (*Gadus morhua*). *Can J Zool* 59:2186–2192.
- Fletcher GL, Hew CL, Joshi SB (1982) Isolation and characterization of antifreeze glycoproteins from the frostfish, *Microgadus tomcod*. *Can J Zool* 60:348–355.
- Zhuang X, Yang C, Fevolden S-E, Cheng CH (2012) Protein genes in repetitive sequence-antifreeze glycoproteins in Atlantic cod genome. *BMC Genomics* 13:293.
- O'Grady SM, Schrag JD, Raymond JA, DeVries AL (1982) Comparison of antifreeze glycopeptides from Arctic and Antarctic fishes. *J Exp Zool* 224:177–185.
- Kashi Y, King DG (2006) Simple sequence repeats as advantageous mutators in evolution. *Trends Genet* 22:253–259.
- Clark MB, et al. (2011) The reality of pervasive transcription. *PLoS Biol* 9:e1000625.
- Kozak M (1987) An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Res* 15:8125–8148.
- Nicodemus-Johnson J, Silic S, Ghigliotti L, Pisano E, Cheng C-HC (2011) Assembly of the antifreeze glycoprotein/trypsinogen-like protease genomic locus in the Antarctic toothfish *Dissostichus mawsoni* (Norman). *Genomics* 98:194–201.
- Cheng CC, Cziko PA, Evans CW (2006) Nonhepatic origin of notothenioid antifreeze reveals pancreatic synthesis as common mechanism in polar fish freezing avoidance. *Proc Natl Acad Sci USA* 103:10491–10496.
- Andersson DI, Jerlström-Hultqvist J, Näsvall J (2015) Evolution of new functions de novo and from preexisting genes. *Cold Spring Harb Perspect Biol* 7:a017996.
- Reinhardt JA, et al. (2013) De novo ORFs in *Drosophila* are important to organismal fitness and evolved rapidly from previously non-coding sequences. *PLoS Genet* 9:e1003860.
- Tautz D (2014) The discovery of de novo gene evolution. *Perspect Biol Med* 57:149–161.
- Carvunis A-R, et al. (2012) Proto-genes and de novo gene birth. *Nature* 487:370–374.
- Zhao L, Saelao P, Jones CD, Begun DJ (2014) Origin and spread of de novo genes in *Drosophila melanogaster* populations. *Science* 343:769–772.
- Palmieri N, Kosiol C, Schlötterer C (2014) The life cycle of *Drosophila* orphan genes. *eLife* 3:e01311.