

## Shaping and reshaping of salmonid genomes by amplification of tRNA-derived retroposons during evolution

YUKIHARU KIDO\*, MITSUKO AONO\*, TOSHIFUMI YAMAKI\*, KEN-ICHI MATSUMOTO\*, SHIGENORI MURATA\*, MINEO SANEYOSHI†, AND NORIHIRO OKADA\*‡

\*Institute of Biological Sciences, University of Tsukuba, Tsukuba Ibaraki 305, Japan; and †Department of Science and Engineering, Nishi-Tokyo University, Uenohara, Yamanashi 409-04, Japan

Communicated by Susumu Ohno, December 3, 1990

**ABSTRACT** Three families of tRNA-derived repeated retroposons in the genomes of salmonid species have been isolated and characterized. These three families differ in sequence, but all are derived from a tRNA<sup>Lys</sup> or from a tRNA species structurally related to tRNA<sup>Lys</sup>. The salmon *Sma* I family is present in the genomes of two species of the genus *Oncorhynchus* but not in other species, including five other species of the same genus. The charr *Fok* I family is present only in four species and subspecies of the genus *Salvelinus*. The third family, the salmonid *Hpa* I family, appears to be present in all salmonid species but is not present in species that are not members of the Salmonidae. Thus, the genome of proto-Salmonidae was originally shaped by amplification and dispersion of the salmonid *Hpa* I family and then reshaped by amplification of the *Sma* I and *Fok* I families in the more recently evolved species of salmon and charr, respectively. We speculate that amplification and dispersion of retroposons may have played a role in salmonid speciation.

Gene duplication is believed to be of major importance in creating genetic diversity (1). The genes for immunoglobulins, histocompatibility complexes, and globins are examples of this gene duplication. This mechanism operates at the DNA level and probably has as old a history as DNA genomes themselves. Another mechanism for maintaining the fluidity of eukaryotic genomes is that recently characterized retroposition, in which information in nonviral cellular RNA can flow back into the genome via cDNA intermediates (2, 3). Retroposition creates additional sequence combinations through dispersal of genetic information and can shape and reshape eukaryotic genomes in many different ways (3, 4). The precise mechanism of retroposition is at present speculative. Recently, Weiner and Meizels (5) presented an interesting hypothesis concerning the mechanism of generation of duplex DNA at the beginning of the DNA world, proposing that duplex DNA genomes may have been derived from earlier DNA genomes that replicated like retroviruses through an RNA intermediate. This suggests that the mechanism of retroposition might be closely linked to that of replication of retroviruses (6).

The highly repetitive sequences that are interspersed throughout eukaryotic genomes have been classified into two categories based on size: long interspersed repetitive elements (LINEs), which include L1 sequences, and short interspersed repetitive elements (SINEs), such as the primate *Alu* and rodent type 1 or 2 *Alu* families (7). Previously, highly repetitive and transcribable sequences have been found in the genome of the chum salmon (*Oncorhynchus keta*) (8, 9). Like all SINE families examined so far (10–14) other than *Alu* (15, 16), this *Sma* I family [formerly the salmon polymerase (Pol) III/SINE family] has been shown to be derived from a tRNA;

moreover the *Sma* I family has several of the characteristic features of retroposons and appears to be the youngest SINE family characterized to date.

The genus *Oncorhynchus* has many species, most of which are believed to have been generated recently on an evolutionary time scale. From a viewpoint of biogeographical aspects, Neave (17) has hypothesized that many of the Pacific salmon diverged from an ancestor of the cherry salmon (*Oncorhynchus masou*) in the Japan Sea during the glacial period of the late Pleistocene era (about one million years ago). Moreover, there are numerous species of charr and trout, and these show remarkable variation in sexual dimorphism, breeding shape, color, and life history. Because of the presence of these many young species and the youngest retroposon, the family Salmonidae provides an especially promising system for studying possible relationships between speciation and dispersion of retroposons. Here, we sketch the behavior of retroposon families, including one found in the chum salmon (9), during the evolution of the salmonid species.<sup>§</sup>

### MATERIALS AND METHODS

Experimental procedures were performed by standard methods (9, 18–21).

### RESULTS

**Classification of Salmonid Species and Strategy for Analyses of Retroposons.** The fish species examined in this study are listed in Table 1. The family Salmonidae consists mainly of the four genera, *Oncorhynchus*, *Salmo*, *Salvelinus*, and *Hucho*, in addition to a distantly related genus, *Coregonus*. Previous reports from our laboratory have shown that the genome of the chum salmon (*O. keta*) contains highly repetitive sequences derived from a tRNA<sup>Lys</sup> gene (9). The sequences of two genomic clones, *Sma*(OK)-2 and *Sma*(OK)-3 (9), are presented in Fig. 1. The characteristic features of the retroposon in this family are dispersion in the genome, an A+T-rich region at the 3' end of the repeat, and direct terminal repeats abutting what appear to be the boundaries of the transposed element. To examine the distribution of this family, we performed a dot hybridization experiment using as the probe a labeled RNA transcribed by T7 RNA polymerase from the tRNA-related region of the *Sma*(OK)-2 clone. As shown in Fig. 2a, DNA from the pink salmon (*O. gorbuscha*) gave a strong hybridization signal, suggesting that the genome of the pink salmon contains a repetitive family similar to that of the chum salmon. To our surprise, although the DNAs from other species of the genus *Oncorhynchus* did not give significant signals, the DNAs from several species of *Salvelinus* gave as strong hybridization signals as that of the

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "advertisement" in accordance with 18 U.S.C. §1734 solely to indicate this fact.

Abbreviation: SINE, short interspersed repetitive element.

<sup>‡</sup>To whom reprint requests should be addressed.

<sup>§</sup>The sequences reported in this paper have been deposited in the GenBank data base (accession nos. D90289–D90300).

Table 1. Fish species analyzed

Order	Family	Genus	Species	Common name		
Salmoniformes	Salmonidae	<i>Oncorhynchus</i>	<i>gorbuscha</i>	Pink salmon		
			<i>keta</i>	Chum salmon		
			<i>tshawytcsha</i>	Chinook salmon		
			<i>kisutch</i>	Coho salmon		
			<i>nerka adonis</i>	Kokanee		
			<i>masou</i>	Cherry salmon		
			<i>mykiss</i>	Steelhead trout		
			<i>Salmo</i>	<i>trutta</i>	Brown trout	
				<i>malma</i>	Dolly Varden	
			<i>Salvelinus</i>	<i>leucomaenis leucomaenis</i>	White-spotted charr	
				<i>leucomaenis pluvius</i>	Japanese common charr	
					<i>namaycush</i>	Lake trout
					<i>Hucho</i>	Japanese huohen
					<i>Coregonus</i>	Peled
	Plecoglossidae	<i>Plecoglossus</i>	Ayu fish			
	Osmeridae	<i>Hypomesus</i>	Pond smelt			
Perciformes	Channidae	<i>Channa</i>	<i>argus</i>	Snakehead		

chum salmon. These results suggest that the sequence of the repetitive family in the genomes of *Salvelinus* species is similar to that of the chum salmon or that it is less similar but is present in high copy number. The DNAs of species such as the coho salmon (*O. kisutch*), cherry salmon, kokanee (*O. nerka adonis*), and brown trout (*Salmo trutta*) did not hybridize with the probe significantly, but preliminary experiments using total genomic transcription suggested that they contained another highly repetitive and transcribable sequence. We have therefore determined the repetitive sequences present in the pink salmon, one species of *Salvelinus*, and two *Oncorhynchus* species that are less similar to that of the chum salmon.

**The Salmon *Sma* I Family Is Restricted to the Genomes of the Chum and Pink Salmon.** Using labeled RNA of the tRNA-related region of *Sma*(OK)-2 DNA, we isolated several phage clones from a genomic library of the pink salmon. As shown in Fig. 1, these sequences [clones *Sma*(OG)-5 and *Sma*(OG)-7] are almost identical to that in the chum salmon. To examine the distribution of this family among the salmonid species, we performed a PCR experiment using the DNAs of 17 fish species as templates. For this experiment we used 2 ng of template DNA and 20 cycles of synthesis, so that only repetitive sequences were detected. The results in Fig. 3a clearly indicate that this family is confined to two species, the chum salmon and pink salmon. This family is collectively designated the salmon *Sma* I family. A consensus sequence

was deduced from four sequences of the salmon *Sma* I family. The average sequence divergence of the salmon *Sma* I family is 0.7%, indicating that this family was amplified very recently; the comparable value for the human *Alu* family is 14%, suggesting that the *Sma* I family is younger than the *Alu* family by a factor of 20. As described previously (9), the tRNA-related region of the salmon *Sma* I family shows equal similarity to a tRNA<sup>Lys</sup> and an tRNA<sup>Ile</sup> (overall identity 74%).

**The Charr *Fok* I Family Is Restricted to *Salvelinus* Species.** Using the tRNA-related region of the salmon *Sma* I family as a probe, we isolated several phage clones from a genomic library of the white-spotted charr (*Salvelinus leucomaenis leucomaenis*). Four sequences are shown in Fig. 4. This family was named the charr *Fok* I family. Like the salmon *Sma* I family, this family consists of a tRNA-related region and tRNA-unrelated region. The tRNA-related region is similar to a tRNA<sup>Lys</sup> and an tRNA<sup>Ile</sup> (overall identity 75% and 72%, respectively). In particular, the aminoacyl stem of the tRNA-related sequence clearly resembles that of a tRNA<sup>Lys</sup>, suggesting this tRNA species as the parentage of the *Fok* I family. This is a remarkable feature of the charr *Fok* I family, because almost all other tRNA-derived retroposons show weak homology in this region (13). A CCA sequence at the 3' end of the tRNA is retained in the family, suggesting that the tRNA-related region was derived from tRNA itself, not from transfer DNA. Fig. 5 shows a comparison of the sequences and structures of the tRNA-related region of the charr *Fok* I

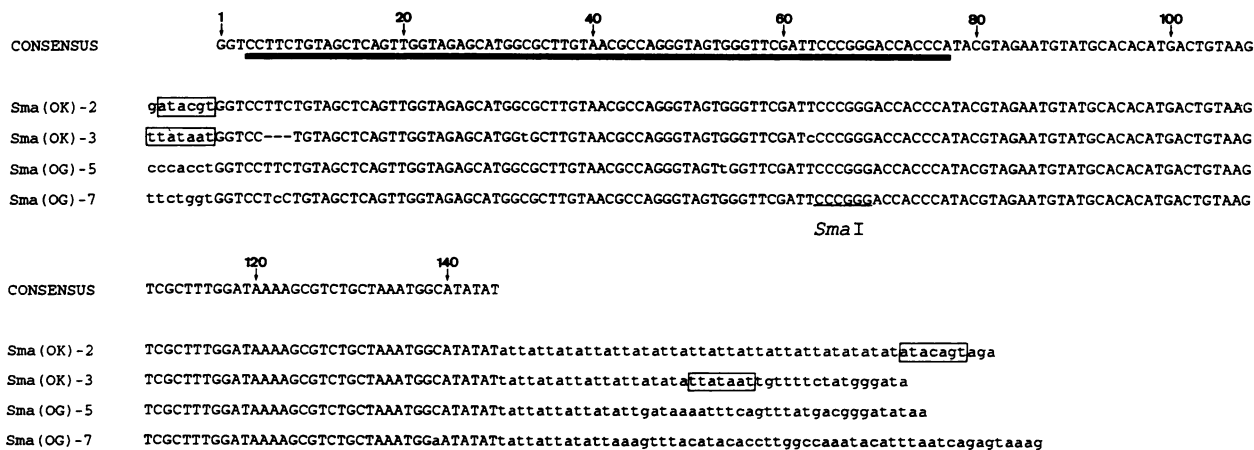


FIG. 1. Sequences and consensus sequence of the salmon *Sma* I family. *Sma*(OK)-2 and -3 [formerly Sm2 and Sm3, respectively] from the chum salmon have been described previously (9). *Sma*(OG)-5 and -7 are from the pink salmon. Direct terminal repeats are boxed, and the tRNA-related region of the family is underlined.

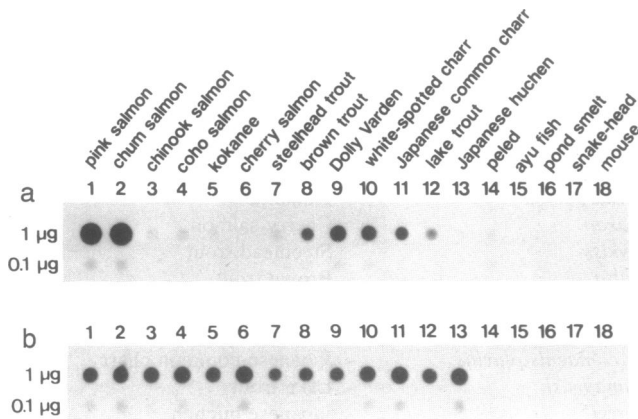


FIG. 2. Demonstration of retroposons in various salmonid species. Dot hybridization experiments were performed using T7 RNA polymerase-transcribed labeled RNA from the tRNA-related region of Sma(OK)-2 (residues 1-68) (a) or the tRNA-related region of Hpa(OM)-1 (residues 1-111) (b) as a probe.

family and a tRNA<sup>Lys</sup>. The charr *Fok* I family is the repetitive family that exhibits the highest degree of similarity with a specific tRNA molecule among the tRNA-derived retroposons thus far characterized (9-14), with a calculated average sequence divergence of 0.9%.

The distribution of the charr *Fok* I family was examined by PCR. As shown in Fig. 3b, this family is present in *Salvelinus* species only, although the genome of the lake trout appears to contain a subfamily in addition to this family, both of which probably have a small number of members. The charr *Fok* I family may have been amplified at the time of divergence of the genus *Salvelinus* or soon after.

**The Salmonid *Hpa* I Family Is Present in All Salmonid Species Examined.** Two of the salmonid species whose DNAs did not hybridize with Sma(OK)-2 DNA of the salmon *Sma* I family were selected for sequence analysis. One of these, the cherry salmon, was chosen because this species retains some primitive phenotypes and is thought to be a progenitor of other Pacific salmon (see the Introduction). The other species selected was the kokanee, which is phylogenetically closely related to the chum and pink salmon, although the exact relationships of these three species are unknown. Two sequences of DNA of the cherry salmon and two of the kokanee were determined (Fig. 6). These four sequences constitute one repetitive family, which we named the salmonid *Hpa* I family. The 5' half of the salmonid *Hpa* I family appears to be derived from a tRNA, but its extent of similarity

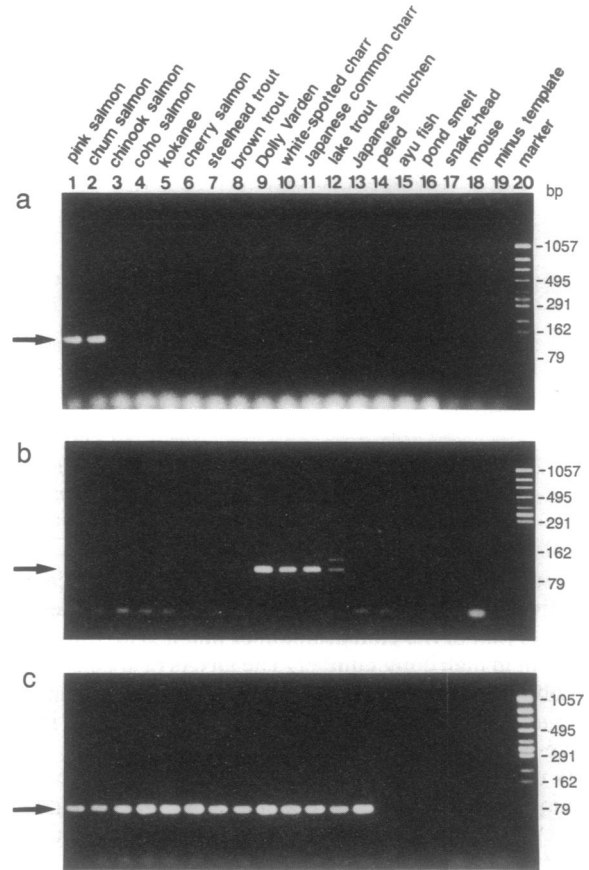


FIG. 3. PCR analysis of distribution of retroposons in salmonid species. Primers were as follows: a, specific for the salmon *Sma* I family (residues 1-20 and 95-76 of the consensus sequence); b, specific for the charr *Fok* I family [residues 5-23 and residues 106-87 of *Fok*(SLL)-3]; and c, specific for the salmonid *Hpa* I family (residues 11-30 and 84-64 of the consensus sequence). Lanes 19 and 20: results without template and for marker DNA (*Hinc*II-digested  $\phi$ X174 DNA). The fragment synthesized in each experiment is shown by an arrow. bp, Base pairs.

with known tRNA species is quite low; it is equally similar (55-60%) to tRNA<sup>Lys</sup>, tRNA<sup>Ile</sup>, tRNA<sup>Thr</sup>, and tRNA<sup>Tyr</sup> (data not shown). The average sequence divergence is 3.1%, indicating that this family is the oldest.

The PCR was used to investigate the distribution of this family. As shown in Fig. 3c, this family was found in all salmonid species belonging to *Oncorhynchus*, *Salmo*,

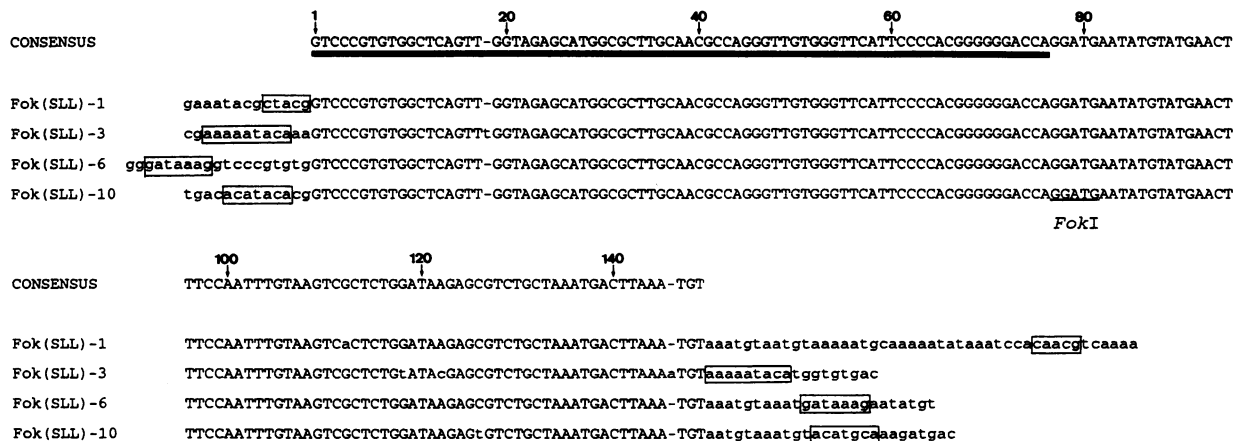


FIG. 4. Sequences and consensus sequence of the charr *Fok* I family. Direct terminal repeats are boxed, and the tRNA-related region is underlined.

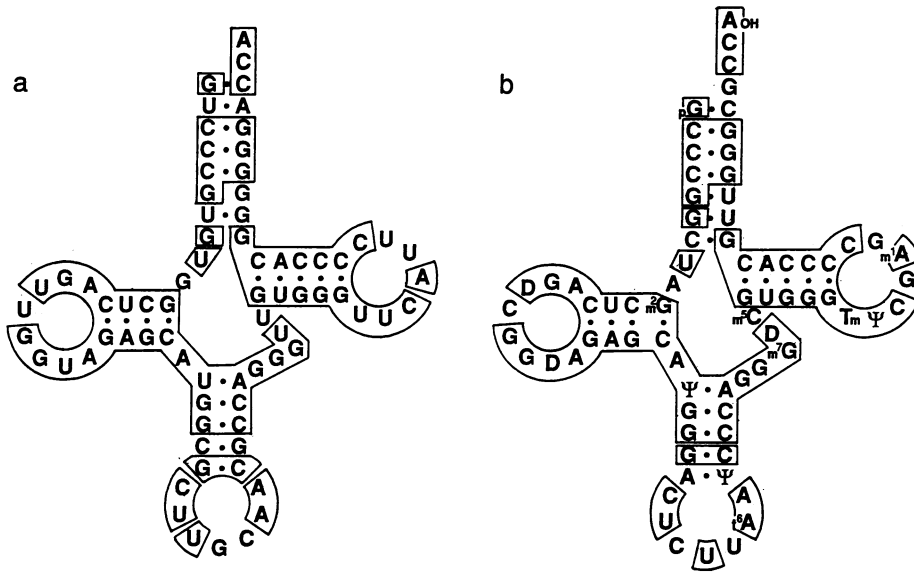


FIG. 5. Sequence and structural similarities of the tRNA-related segment of the charr *Fok I* family (a; see Fig. 4) and tRNA<sup>Lys</sup> (b). The sequence of tRNA<sup>Lys</sup> is taken from Raba *et al.* (22). Recently, it has been shown that sequences of the genes for tRNA<sup>Lys</sup> in the chum salmon are identical to those in rabbit (H. Sano, Y.K., and N.O., unpublished results). Identical sequences are boxed.

*Salvelinus*, and *Hucho*, consistent with the age of the family inferred by average sequence divergence. The distribution of the salmonid *Hpa I* family was also confirmed by a hybridization experiment with labeled RNA from the tRNA-related region of *Hpa*(OM)-1 DNA as a probe (Fig. 2b). The salmonid *Hpa I* family is apparently not present in the genomes of species such as snakehead (*Channa argus*), ayu fish (*P. altivelis*), and pond smelt (*Hypomesus transpacificus nipponensis*). As for a peled (*Coregonus peled*), ≈100 copies of the salmonid *Hpa I* family (about 1/100 members of the cherry salmon *Hpa I* family) are found to be present in its genome by PCR experiments (data not shown), so there must have been at least two amplification events of the salmonid *Hpa I* family in the major lineage of Salmonid evolution, which occurred at the time of establishment of the family Salmonidae and after divergence of the genus *Coregonus*, respectively (see Fig. 8 and Discussion).

DISCUSSION

**A tRNA<sup>Lys</sup> or a Structurally Related tRNA as the Origin of Retroposons.** An alignment of the consensus sequences of the *Sma I* family and the *Fok I* family is shown in Fig. 7. Surprisingly, the two sequences are remarkably similar not only in the tRNA-related region but also in the tRNA-unrelated region. This strongly suggests that these families have a common evolutionary lineage. Whether the putative

ancestor was horizontally transmitted as an agent like an RNA virus into the salmonid genomes followed by amplification or had resided within the genomes of the family Salmonidae long before amplification remains to be elucidated. Since the tRNA-related region of the *Fok I* family is likely to be derived from tRNA<sup>Lys</sup> (Fig. 5), it is presumed that the tRNA-related region of the *Sma I* family is also originated from tRNA<sup>Lys</sup>.

In the present work, we found that a tRNA<sup>Lys</sup> is the tRNA species most similar to all three families of salmonid retroposons. We have also found that tortoise Pol III/SINE shows closest similarity to a tRNA<sup>Lys</sup> and less similarity to a tRNA<sup>Thr</sup> (23). A tRNA<sup>Lys</sup>-like structure appears to be widespread among SINEs in the animal kingdom. The significance of this finding is not clear at present. Since cellular RNA must be copied into cDNA before being retrotransposed, one possible explanation is that a tRNA<sup>Lys</sup>-like structure within a retroposon can be preferentially copied by a reverse transcriptase that normally uses a tRNA<sup>Lys</sup> as a primer; in fact, tRNA primers are known to form stable binary complexes with appropriate reverse transcriptases (24, 25). This would imply that a retrovirus or reverse transcriptase that uses a tRNA<sup>Lys</sup> primer exists in salmonid species. Another explanation is that a tRNA<sup>Lys</sup>-like structure may confer a special selective advantage on their host or function as a regulatory element, as discussed elsewhere (6).

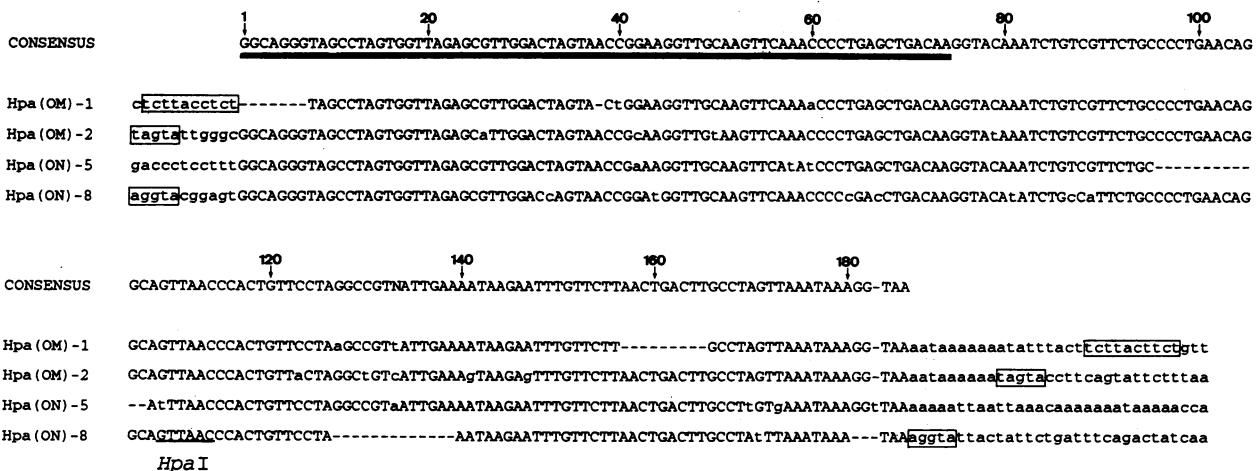


FIG. 6. Sequences and consensus sequence of the salmonid *Hpa I* family. Hpa(OM)-1 and -2 are from the cherry salmon and Hpa(ON)-5 and -8 are from the kokanee. Direct terminal repeats are boxed, and the tRNA-related region is underlined.

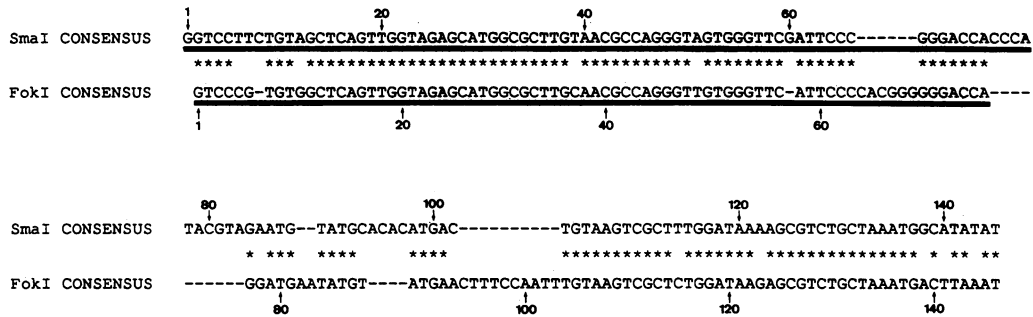


FIG. 7. Comparison of the consensus sequences of the *Sma* I family and the *Fok* I family. Identical nucleotides were indicated by stars. The tRNA-related regions are underlined.

**Shaping and Reshaping of the Salmonid Genomes During Evolution.** Fig. 8 shows a recently proposed phylogenetic tree of the salmonid species (26, 27), in which arrows indicate possible times of amplification of the three retroposons, the salmon *Sma* I family, the charr *Fok* I family, and the salmonid *Hpa* I family.

The aim of the present work was not to elucidate the relationships among salmonid species. However, several important conclusions on evolution of salmonid species can be drawn from our results. (i) The genus *Coregonus* diverged first from the other four genera of Salmonidae described here, since the number of members of the salmonid *Hpa* I family in the peled is about 1/100th that of the cherry salmon *Hpa* I family. (ii) The existence of the salmon *Sma* I family in the genomes of the chum and pink salmon but not in that of the kokanee (*O. nerka adonis*) appears to solve the problem of classification of these three species (26). These three species are believed to have deviated from other species of *Oncorhynchus*. The chum and pink salmon have reduced the freshwater phase of life to a minimum, fry emerging from riverbeds ready for downstream migration. A peculiar characteristic of the sockeye salmon (*O. nerka*) is that a small percentage of their fry go directly to sea, where only a few are presumed to survive. This probably reflects a natural experiment of this species to develop a new life cycle including an in-sea-water phase (26). Thus, the genetic relationships of these three species have been confusing. The present work

strongly suggests that the kokanee diverged first from the other two species not the reverse.

A more detailed analysis of salmonid SINEs may make it possible to ask whether amplification and dispersion of SINEs by retroposition is a cause or a consequence of speciation and whether retroposition can facilitate reproductive isolation once speciation has begun.

We thank Prof. Y. Watanabe, Dr. S. Nishimura, and Prof. S. Osawa for encouragement. We thank Dr. S. Ohno for suggesting to N.O. to analyze DNA of fishes other than the Salmonidae. We are grateful to Drs. A. Weiner, T. Okazaki, and T. Iwami for critical reading of the manuscript.

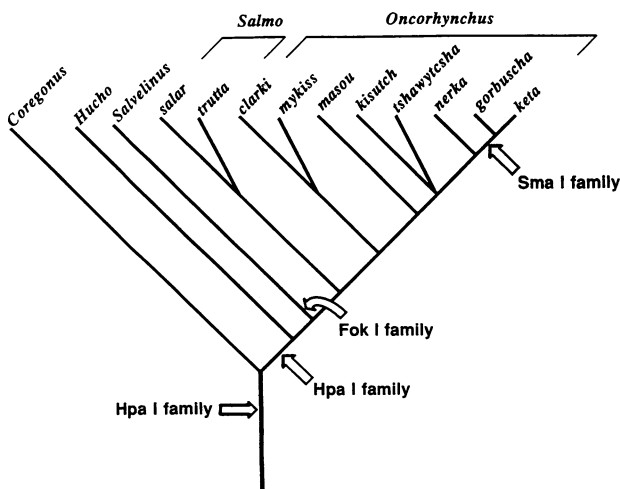


FIG. 8. Phylogenetic tree of salmonid species. The tree is according to Smith and Stearley (26) except for the relationships of the three species *nerka*, *gorbuscha*, and *keta*, which are according to Thomas et al. (27). Another tree concerning relationships among *Oncorhynchus* species has recently been proposed (28). Possible times of amplification of the three retroposons are indicated by arrows.

- Ohno, S. (1970) *Evolution by Gene Duplication* (Springer, Heidelberg).
- Rogers, J. (1985) *Int. Rev. Cytol.* **93**, 187-279.
- Weiner, A. M., Deininger, P. I. & Efstratiadis, A. (1986) *Annu. Rev. Biochem.* **55**, 631-661.
- Baltimore, D. (1985) *Cell* **40**, 481-482.
- Weiner, A. M. & Meizels, N. (1990) in *Proceedings of the International Symposium on Evolution of Life*, ed. Osawa, S. (Springer, Tokyo), in press.
- Okada, N. (1990) *J. Mol. Evol.* **31**, 500-510.
- Singer, M. F. (1982) *Cell* **28**, 433-434.
- Matsumoto, K., Murakami, K. & Okada, N. (1984) *Biochem. Biophys. Res. Commun.* **124**, 514-522.
- Matsumoto, K., Murakami, K. & Okada, N. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 3156-3160.
- Lawrence, C. B., McDonnell, D. P. & Ramsey, W. J. (1985) *Nucleic Acids Res.* **13**, 4239-4252.
- Daniels, G. R. & Deininger, P. L. (1985) *Nature (London)* **317**, 819-822.
- Okada, N., Endoh, H., Sakamoto, K. & Matsumoto, K. (1985) *Proc. Jpn. Acad. Ser. B* **61**, 363-367.
- Sakamoto, K. & Okada, N. (1985) *J. Mol. Evol.* **22**, 134-140.
- Endoh, H. & Okada, N. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 251-255.
- Weiner, A. M. (1980) *Cell* **22**, 209-218.
- Ullu, E. & Tschudi, C. (1984) *Nature (London)* **312**, 171-172.
- Neave, F. (1958) *Trans. R. Soc. Can.* **52**(III-5), 25-39.
- Blin, N. & Stafford, D. W. (1976) *Nucleic Acids Res.* **3**, 2303-2308.
- Manley, J. L., Fire, A., Cano, A., Sharp, P. A. & Gelfand, M. L. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3855-3859.
- Sanger, F., Nicklen, S. & Coulson, A. R. (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463-5467.
- Saiki, R. K., Gelfand, D. H., Stoffel, S., Scharf, S. J., Higuchi, R., Horn, G. T., Mullis, K. B. & Erlich, H. A. (1988) *Science* **239**, 487-491.
- Raba, M., Limburg, K., Burghagen, M., Katze, J. R., Simsek, M., Heckman, J. E., RajBhandary, U. L. & Gross, H. J. (1979) *Eur. J. Biochem.* **97**, 305-318.
- Endoh, H., Nagahashi, S. & Okada, N. (1990) *Eur. J. Biochem.* **189**, 25-31.
- Panet, A., Haseltine, W. A., Baltimore, D., Peters, G., Harada, F. & Dahlberg, J. E. (1975) *Proc. Natl. Acad. Sci. USA* **72**, 2535-2539.
- Barat, C., Lullien, V., Schatz, O., Keith, G., Nugeyre, M. T., Grüniger-Leitch, F., Barré-Sinoussi, F., LeGrice, S. F. J. & Darlix, J. L. (1989) *EMBO J.* **8**, 3279-3285.
- Smith, G. R. & Stearley, R. F. (1989) *Fisheries* **14**, 4-10.
- Thomas, W. K., Withler, R. E. & Beckenbach, A. T. (1986) *Can. J. Zool.* **64**, 1058-1064.
- Thomas, W. K. & Beckenbach, A. T. (1989) *J. Mol. Evol.* **29**, 233-245.