

The rat serum albumin gene: Analysis of cloned sequences

(genome libraries/cDNA clones/R-loops/Southern blots)

THOMAS D. SARGENT, JUNG-RUNG WU, JOSE M. SALA-TREPAT*, R. BRUCE WALLACE†, ANTONIO A. REYES, AND JAMES BONNER

Division of Biology, California Institute of Technology, Pasadena, California 91125

Contributed by James F. Bonner, April 27, 1979

ABSTRACT The rat serum albumin gene has been isolated from a recombinant library containing the entire rat genome cloned in the λ phage Charon 4A. Preliminary R-loop and restriction analysis has revealed that this gene is split into at least 14 fragments (exons) by 13 intervening sequences (introns), and that it occupies a minimum of 14.5 kilobases of genomic DNA.

Recent advances in recombinant DNA technology have made it possible to obtain virtually any desired single-copy genomic sequence in cloned form, provided an appropriate probe is available. We have used these techniques to isolate the rat serum albumin gene. Serum albumin synthesis is one of the major characteristics of vertebrate liver. Observation of the activity and state of this gene during development and in adult tissues should be informative as to the process of terminal differentiation. Albumin synthesis is essentially constitutive, but does respond significantly to a variety of stimuli (1). It is also expressed to variable extents in different hepatoma cell lines (2). The availability of cloned albumin genomic DNA will greatly facilitate the study of this variable expression, particularly at the level of transcript processing.

Determination of the sequence organization of the albumin gene is also of interest, especially with regard to the disposition of repetitive elements and intervening sequences. Although regulatory (3) and evolutionary (4) significance has been postulated, the functional role, if any, of these striking features of eukaryotic genomes remains unknown. The comparative studies that will be possible as other genes are extracted from the rat and related species can be expected to provide considerable insight into this fascinating problem.

MATERIALS AND METHODS

Rat Genome Library. High molecular weight liver DNA was extracted from an adult male Sprague-Dawley rat (Simonsen Labs, Gilroy, CA) by the method of Blin and Stafford (5) and aliquots were digested with *Eco*RI (Boehringer Mannheim) under conditions adjusted to cleave either one-third or one-fifth of the *Eco*RI sites in an equivalent amount of bacteriophage λ DNA. The fragments resulting from this partial digestion were sedimented through a 10–30% sucrose gradient; the material between 10 and 20 kilobases (kb) was recovered by ethanol precipitation. A sample of this rat DNA (2.5 μ g) was ligated with 8.5 μ g of a preparation of Charon 4A "cloning fragments" (6, 7). This recombinant DNA was packaged *in vitro* by using extracts from defective λ lysogens provided by N. Sternberg (6). The method used was that of Hohn and Murray (8). Approximately 2,000,000 independent clones were obtained. The library was amplified 100,000-fold by subconfluent plating on *Escherichia coli* strain DP50SupF (9).

cDNA Clones. cDNA was synthesized from purified albumin

mRNA as described (10). This cDNA contained a small amount of full-length material and had a number average size of approximately 1000 nucleotides. It was rendered double-stranded by sequential treatment with *E. coli* DNA polymerase I and S1 nuclease (11). The resulting DNA had a number average size of 600 nucleotide pairs. An average of 10 dCMP residues were polymerized per 3' end by terminal transferase (ref. 12; W. Rowekamp and R. Firtel, personal communication). Forty nanograms of the tailed albumin cDNA was mixed with 200 ng of pBR322 DNA that had been cleaved with *Pst* I and similarly polymerized with dGMP residues at the 3' ends (a gift of W. Rowekamp) and induced to cocircularize by incubation at 42°C for 4 hr followed by 16 hr of slow cooling to 4°C. *E. coli* strain χ 1776 was transformed with this mixture (13, 14). Several hundred colonies were obtained on appropriate selective media. These were screened by the filter colony hybridization method of Grunstein and Hogness (15), with ³²P-labeled albumin cDNA as a probe. The most intensely reacting clones were selected and plasmid DNA was prepared. Proof of their identity was obtained by comparison of partial nucleotide sequences to the existing amino acid sequence of rat serum albumin (ref. 16; T. D. Sargent and J. Posakony, unpublished data).

Screening. Approximately 1,000,000 plaques from the rat library were screened by a modification of the method of Benton and Davis (17), using as a probe nick-translated albumin cDNA clones. Nine different genome clones were obtained, three of which, λ RSA14, λ RSA30, and λ RSA40, are the subject of the present report.

R-Loops. Fifty nanograms of recombinant phage DNA that had been digested with either *Eco*RI or *Hind*III was mixed with 40 ng of purified albumin mRNA in a total volume of 20 μ l of 70% recrystallized formamide/0.4 M NaCl/5 mM EDTA/80 mM Pipes (1,4-piperazinediethanesulfonic acid), pH 7.4 (18) and incubated at 50–53°C for 24 hr. The hybrids were spread for electron microscopy by the modified Kleinschmidt method (19). The grids were rotary shadowed with platinum and palladium (80/20) and viewed in a Phillips 300 electron microscope. DNA contour lengths were measured with a Hewlett-Packard digitizer.

Blots. Rat liver DNA from the same preparation used to make the genomic library was digested three times with a 6-fold excess of either *Eco*RI or *Hind*III, extracted with phenol, and precipitated with ethanol. Ten micrograms of the digested DNA was fractionated on an 8-mm-thick 0.8% agarose slab gel buffered with 50 mM Tris/18 mM NaCl/2 mM EDTA/20 mM sodium acetate, pH 7.4. Electrophoresis was at 2.5 V/cm for 16 hr. Clone DNA was digested once with a 10-fold excess of restriction endonuclease and 0.2 μ g was electrophoresed on a 2-mm-thick gel. Transfer to nitrocellulose (Millipore) was es-

Abbreviation: kb, kilobase.

* Present address: Laboratoire d'Enzymologie, Centre National de la Recherche Scientifique, 91190 Gif-sur-Yvette, France.

† Present address: Department of Biology, City of Hope National Medical Center, Duarte, CA 91010.

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.

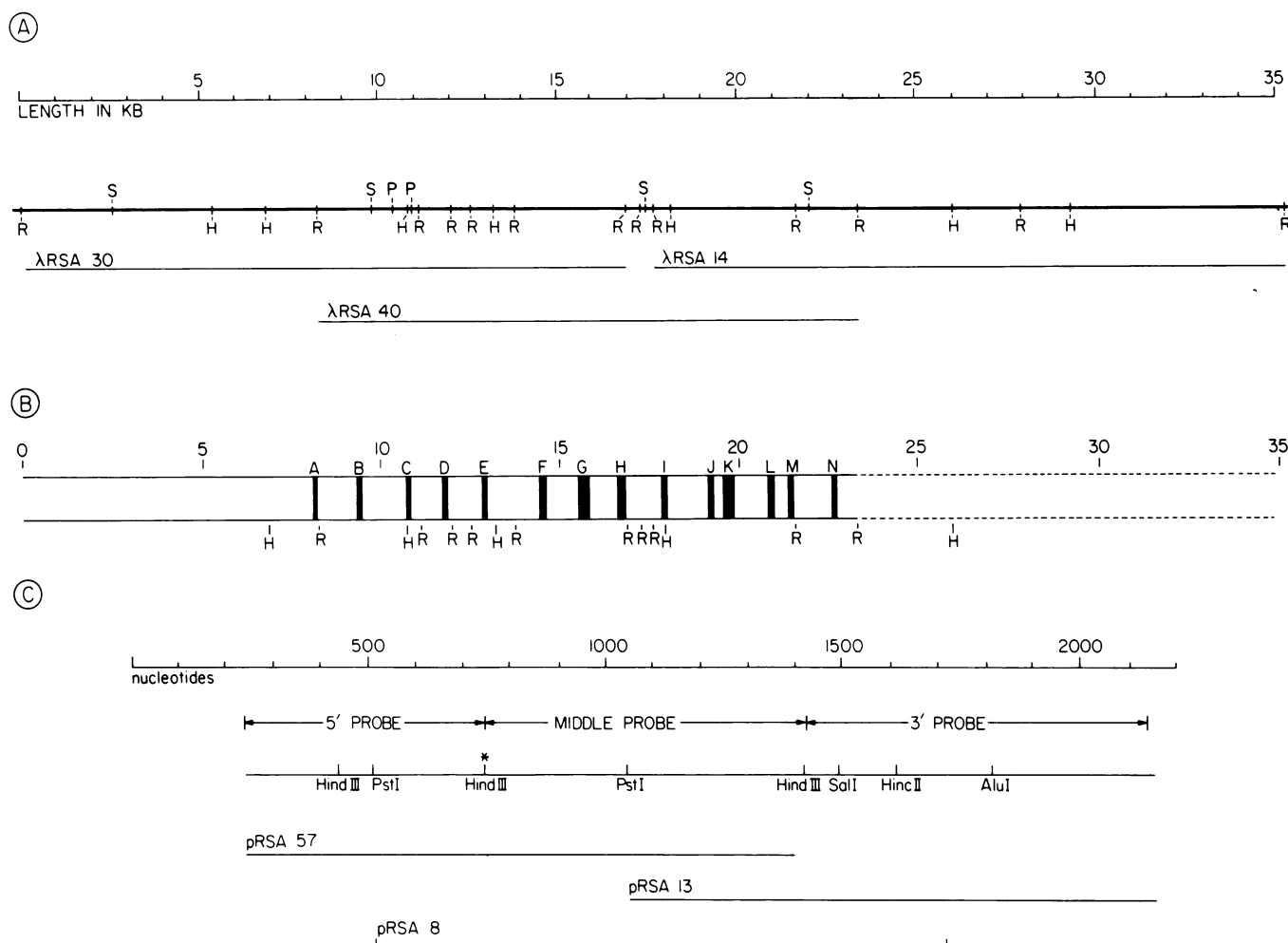


FIG. 1. (A) Restriction site map for three albumin genome clones, λ RSA30, λ RSA40, and λ RSA14. R, *EcoRI*; H, *HindIII*; S, *Sac I*; P, *Pst I*. Molecular weights were estimated by comparison with phage PM2 DNA digested with *HindIII*. (B) Intron/exon map of the albumin gene deduced from R-loops. Black bars are exons; white bars are introns. Dashed line indicates region that does not react with albumin cDNA and was not tested for R-loop formation. Scale same as in A. (C) Restriction site map for three albumin cDNA clones, pRSA57, pRSA8, and pRSA13. Asterisk indicates *HindIII* site not found in genomic DNA (see text).

entially as described by Southern (20). Washing after hybridization was done at 62–63°C with a descending series of salt concentrations from 1.0 to 0.1 M NaCl in Denhardt's solution (21).

RESULTS

The recombinant DNA methodology that we have used involves considerable manipulation of DNA, including ligation of a mixture of restriction fragments and several rounds of replication in the bacterial host. This creates the potential for two particularly serious artifactual modifications of genomic sequences: ligation of noncontiguous restriction fragments and genetic rearrangement during propagation. The partial restriction library approach used in this study provides an effective mechanism for detecting the former artifact— independently generated clones confirm the legitimacy of the restriction map of shared DNA since the probability of any given spurious ligation event occurring more than once in the production of a recombinant genome library is negligible. Fig. 1A shows the map of restriction sites for *EcoRI*, *HindIII*, and *Sac I*. The region included in clone λ RSA14 and λ RSA30.

Verification of nonoverlapping regions and exclusion of gross genetic artifacts such as deletions or rearrangements depend upon comparisons made between the cloned sequences and the rat genome, which are accomplished by use of "Southern blots"

of genomic and cloned rat DNA. Fig. 2 shows the pattern obtained when rat DNA is digested with *EcoRI* (lane A) or *HindIII* (lane B), fractionated by electrophoresis, transferred to nitrocellulose, and driven by cloned albumin cDNA labeled with 32 P by nick translation (22). The *EcoRI* digestion results in seven bands complementary to the albumin probe. Band d is quite faint, but is clearly visible in the original and in blots probed with albumin cDNA (unpublished data). This is presumably due to the absence of some 3'-terminal mRNA sequence from the cDNA clones. A mixture of two different plasmid clones, pRSA13 and pRSA57, which includes approximately 85% of the albumin mRNA sequence complexity, is used as a probe in these experiments (unpublished data). The restriction site map for these cDNA clones is shown in Fig. 1C. Since there are no *EcoRI* sites present in the probe sequence, the genome blot suggests that either the gene exists in multiple divergent copies or is interrupted by sequences not present in the mRNA or conceivably both. Evidence will be presented that shows that the albumin gene is in fact interrupted. Fig. 2, lane B, shows the result of a similar experiment with rat DNA digested with *HindIII*. A total of seven bands can be visualized. Four of these, labeled a, b, c, and d, are consistent with the restriction map obtained from the genome clones, as are all but the largest band in Fig. 2, lane A. The remaining bands, indicated by asterisks, are unexpected. These anomalous bands appear with variable intensity in different experiments and are

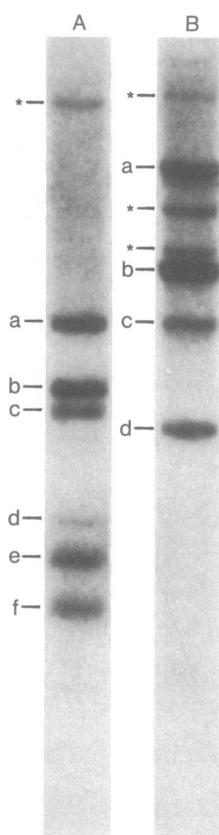


FIG. 2. Rat genome blots. *Eco*RI-digested (lane A) and *Hind*III-digested (lane B) rat liver DNA. Electrophoresis, filter transfer, and hybridization conditions are given in *Materials and Methods*. The sizes of the material in the lettered bands (in kb) are as follows. Lane A: a = 3.9, b = 3.1, c = 2.7, d = 1.6, e = 1.3, and f = 1.0. Lane B: a = 7.9, b = 4.9, c = 3.9, and d = 2.5. Bands marked by asterisks are not present in the genome clones.

probably due to partial digestion. However, their appearance also is consistent with the hypothesis that there are multiple, slightly divergent albumin genes in the rat.

Fig. 3 illustrates analogous experiments performed with DNA from the various genomic clones. In these sets of blots the albumin cDNA probe has been cleaved with *Hind*III and electrophoretically fractionated into a 5', "middle," and 3' probe (Fig. 1C). The patterns generated are entirely consistent with that seen with whole rat liver DNA. The sizes are given in the legend to Fig. 2. This suggests that there has been little or no disruption of individual restriction fragments during the cloning procedure although the resolution limit of the genome blots is probably about 100 nucleotides and small modifications or point mutations would most likely be undetected. A curious aspect of these results is revealed in Fig. 3, gels F and G. Both the 5' and middle albumin mRNA probes react with the same *Hind*III fragment (band d), which implies that the *Hind*III site separating these two probes (indicated by an asterisk in Fig. 1C) does not exist in the genome. The trivial explanation of probe cross contamination is unlikely because there is no analogous reaction with *Eco*RI-digested λ RSA40 (Fig. 3, gels B–D). Sequence rearrangement is not responsible because the phenomenon is observed with two different genome clones, λ RSA30 and λ RSA40, and two different cDNA clones, pRSA13 and pRSA8 (data not shown). Furthermore, this unexpected pattern is seen when 5' and middle albumin probes are used to drive *Hind*III genome blots (unpublished data). Polymorphism cannot be excluded, but seems unlikely because Sprague–Dawley rats are highly inbred. It is conceivable that this restriction site is created in the mRNA by the fusion of two adjacent exons. Final resolution of this question must await determination of the nucleotide sequence of the relevant genomic DNA.

The pattern of hybridizations with the various probes shown in Fig. 3 also establishes colinearity of the genomic and mRNA sequences. The polarity of the cDNA clones was established by

Table 1. Intron and exon lengths from R-loops formed with albumin mRNA and albumin genome clone restriction fragments

Exon	Length (mean \pm SD)	n	Intron	Length (mean \pm SD)	n
A	102 \pm 29	12	AB	927 \pm 64	6
B	111 \pm 31	11	BC	1370 \pm 37	21
C	95 \pm 35	10	CD	946 \pm 145	10
D	148 \pm 34	10	DE	978 \pm 42	10
E	108 \pm 32	12	EF	1458 \pm 88	19
F	163 \pm 34	13	FG	920 \pm 65	13
G	245 \pm 39	15	GH	778 \pm 66	14
H	189 \pm 53	13	HI	908 \pm 62	10
I	133 \pm 26	10	IJ	1161 \pm 90	11
J	131 \pm 27	16	JK	297 \pm 88	11
K	239 \pm 31	14	KL	1011 \pm 74	16
L	146 \pm 28	13	LM	434 \pm 94	13
M	108 \pm 29	14	MN	1054 \pm 163	13
N	125 \pm 58	10			

Sizes in nucleotide pairs of exons and introns in the albumin genomic clones. Introns are designated by two letters, indicating adjacent exons. Lengths were determined by comparison with simian virus 40 DNA contour lengths, as measured from electron micrographs. These data were used to construct the schematic diagram shown in Fig. 1B. n, Number of samples measured for each intron or exon.

nucleotide sequence determination (T. Sargent and J. Posakony, unpublished data) and, thus, the direction of transcription, from left to right in Fig. 1, can be inferred. These blot hybridization data, in conjunction with the restriction maps, suggest that within the cloned locus there exists only one albumin mRNA sequence which is interrupted at least five times and is dispersed over a minimum of 14 kb.

To verify these observations and improve the resolution beyond the limits of a cursory restriction analysis, we digested DNA from the genome clones with either *Eco*RI or *Hind*III, mixed it with purified albumin mRNA, and incubated the mixture under conditions suitable for "R-loop" formation (18). The hybrid structures that formed were then visualized by electron microscopy. The result of these experiments was the identification of a total of 13 interruptions in the albumin mRNA sequence. Selected electron micrographs along with interpretive drawings of the R-loops are shown in Fig. 4. The various hybrids that formed could be identified by their contour lengths, determined by comparison with simian virus 40 DNA spread on the same grids. This permits ordering of R-loop structures by consulting the genome clone restriction maps, but does not specify the correct orientation of the fragments. In most cases the polarity can be deduced from the *Hind*III fragment R-loops by identification of the molecules containing the left or partial right arms of the Charon 4A vector. In one case, shown in Fig. 4, second from left, it was necessary to assume that the end lacking a visible DNA "branch" corresponds to the *Hind*III site located in the mRNA sequence, that is, within exon C. If this assumption is incorrect, then the albumin gene is subdivided into 15 rather than 14 exons, and the order of exons C, D, and E is reversed. The results obtained from *Eco*RI-digested clones are more difficult to interpret due to the larger number of fragments and their similar sizes. However, the structures visualized are consistent with the *Hind*III-digested R-loops. Certain hybrids do not form efficiently, apparently due to close proximity to a restriction cleavage, and these tend to appear with one but not the other digestion. Also, the temperature optima for hybrid formation with the various fragments differed over a few degrees and had to be adjusted accordingly. The data from which the exon/intron map was deduced are summarized in Table 1. These values are in some

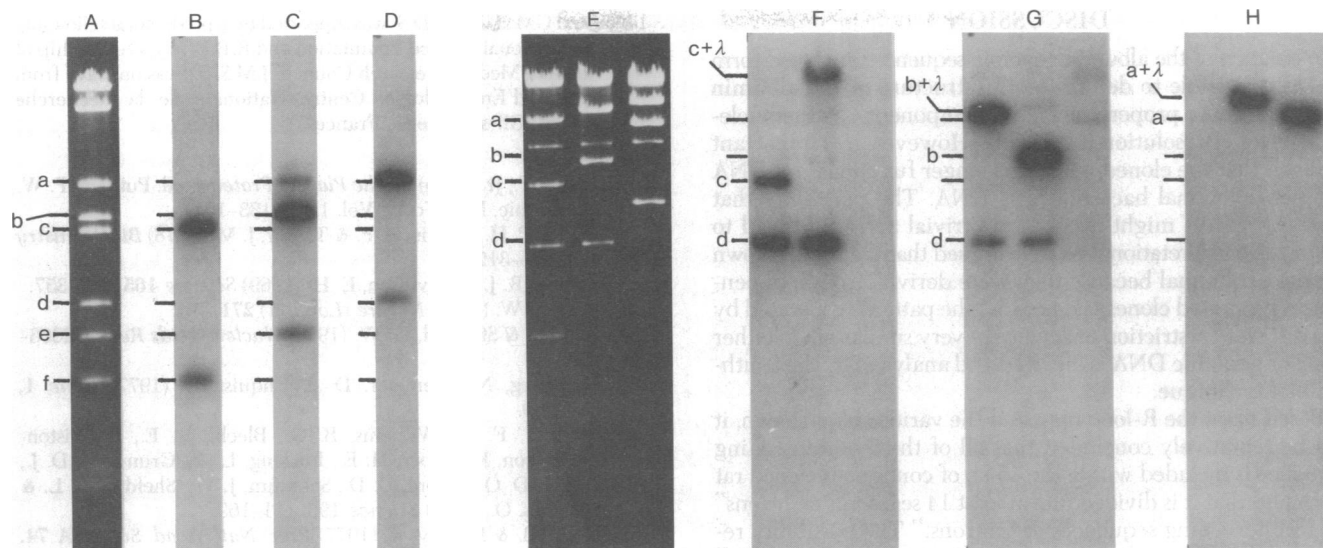


FIG. 3. Albumin genome clone blots. Gel A, photograph of gel of *Eco*RI-digested λ SA40 DNA. Gels B–D, triplicate blots of DNA shown in gel A probed with 5', middle, and 3' fragments of albumin cDNA clones, respectively (see Fig. 1C). Gel E, photograph of gel of *Hind*III-digested λ SA30 (left lane), λ SA40 (middle lane), and λ SA14 (right lane). Gels F–H, triplicate blots of DNA shown in gel E hybridized with 5' (gel F), middle (gel G), and 3' (gel H) probes as above. The sizes of the material in the lettered bands are the same as that given for the corresponding bands in Fig. 2. The order of fragments in the map (Fig. 1A) is as follows, from left to right: *Eco*RI, c, f, e, b, a, d; *Hind*III, c, d, b, a.

cases the sum of two measurements which terminate on opposite sides of the same restriction site. All of the structures represented by the diagram in Fig. 1B are supported by a minimum of 10 measurements made from different unambiguous examples. However, it is important to bear in mind the uncertainty inherent in data of this type. Aside from possible interpretive errors, small exons and introns might not be detectable by electron microscopy. Furthermore, very small hybrids might not be stable under the conditions used for hybridization or spreading of the DNA. DNA displacement loops alongside exons were not visible in most cases. Presumably this is due to the collapse of single-stranded DNA, which may be a property of smaller R-loops prepared by the methods we have used (23).

Close examination of the RNA involved in the R-loop structures reveals a continuous translocation of duplex regions from the 5' end to the 3' end of the mRNA, supporting the conclusion drawn from the blot experiments that there is only

one albumin gene in the cloned complex and that it is colinear with the mRNA sequence, although interrupted. Another interesting observation is that exons A and N seem very near the termini of the mRNA. Very little, if any, RNA extends beyond the 5' side of exon A. This implies that unless a tiny "leader" exon is separated from the rest of the gene by a huge intron or there is more than 8 kb of RNA cleaved from the 5' end of a primary transcript, the initiation of transcription must be located somewhere on the 5' end of clone λ SA30. Similarly, the small "whisker" of RNA visible on the 3' side of exon N is probably mostly poly(A) (10), suggesting that the albumin mRNA coding sequence terminates at this point in the genome. Neither of the *Eco*RI fragments located 3' to this terminal exon reacts with albumin cDNA (unpublished data). Since these fragments represent over 10 kb of genomic DNA, it is unlikely that there are any albumin exons not contained on the genomic clones.

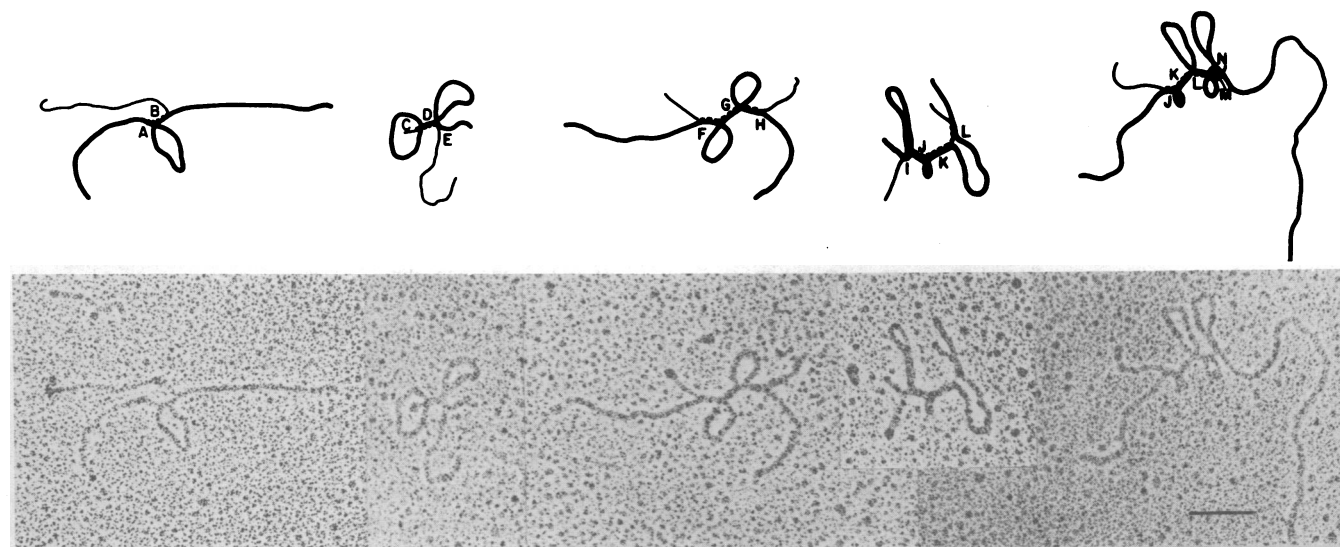


FIG. 4. Selected electron micrographs along with interpretive drawings of R-loops. The restriction fragments involved in the R-loops are, from left to right, *Hind*III fragments c, d, and b, *Eco*RI fragment a, and *Hind*III fragment a + λ (from λ SA40). In the tracings at the top of the figure, heavy solid lines indicate DNA and thin lines represent albumin mRNA. Bar at bottom right represents 500 nucleotide pairs.

DISCUSSION

The isolation of the albumin genomic sequences in cloned form makes it possible to determine the structure of the albumin gene—or, more properly, its DNA component—to the nucleotide level of resolution if desired. However, it is important to recall that the cloned gene is no longer functional rat DNA but nonfunctional bacteriophage DNA. The possibility that rearrangements might occur is not trivial and could lead to major misinterpretation. We have argued that the results shown are not artifactual because they were derived from independently generated clones and because the patterns generated by two different restriction enzymes are very similar when either clone or genomic DNA is digested and analyzed by the Southern blot technique.

Based upon the R-loop map and the various blots shown, it can be tentatively concluded that all of the albumin coding sequence is included within the 35 kb of contiguous cloned rat DNA and that it is divided into at least 14 segments or “exons” by 13 intervening sequences or “introns.” The possibility remains that additional exons or introns exist that are too small to be easily detectable by electron microscopy. The possibility of multiple variants of the albumin gene within the individual rat used to generate the library has not been excluded. However, we have isolated a total of nine albumin clones from the library and, with the exception of one obvious ligation artifact, all contain linear permutations of the same rat DNA restriction fragments. It seems unlikely that other, different albumin genes would have escaped inclusion in the rat library or detection by plaque hybridization.

There seems to be no simple pattern to the size or position of the exons in this gene. The interruptions occur along the entire length of the mRNA sequence. The middle *Hind*III site in the mRNA has been aligned with amino acid number 181 by nucleotide sequence determination (T. Sargent and J. Posakony, unpublished data). This would mean that approximately 200–300 nucleotides of untranslated RNA must reside at the 3′ end of the albumin mRNA (Fig. 1C). If exon 14 does in fact contain the 3′ end of the mRNA, then one and possibly two introns may be located within the untranslated sequence (Table 1).

Assuming that we have not inadvertently cloned an inactive variant of the albumin gene, the arrangement of introns and exons has obvious ramifications regarding the processing of the albumin mRNA precursors. First, the transcription unit is at least 14.5 kb in length and, by analogy to other systems (24, 25), the primary transcript can be expected to be of the same size. This is 3 times larger than the value reported by Shafritz and coworkers (26), suggesting that the 26S species they detected represents accumulated processing intermediates. Second, there should be a minimum of 13 different nuclear intermediates that have been processed to some extent. If multiple pathways exist, then there could be many more species.

We thank Merrie Jo Johnson, Marie Carter, and Leila Gonzalez for excellent technical assistance, and Mark Davis, Doug Engel, Tom Maniatis, and James Posakony for their advice and comments. This work was supported in part by U.S. Public Health Service Grants GM

13262 and GM 20927. T.D.S. was supported by a predoctoral fellowship from the National Science Foundation and R.B.W. by a fellowship of the Canadian Medical Research Council. J.M.S.-T. was on leave from Laboratoire d'Enzymologie, Centre Nationale de la Recherche Scientifique, Gif-sur-Yvette, France.

1. Peters, T., Jr. (1975) in *The Plasma Proteins*, ed. Putnam, P. W. (Academic, New York), Vol. 1, pp. 133–181.
2. Tse, T. P. H., Morris, H. P. & Taylor, J. M. (1978) *Biochemistry* 17, 3121–3128.
3. Britten, R. J. & Davidson, E. H. (1969) *Science* 165, 349–357.
4. Gilbert, W. (1978) *Nature (London)* 271, 501.
5. Blin, N. & Stafford, D. W. (1976) *Nucleic Acids Res.* 3, 2303–2308.
6. Sternberg, N., Tiemeier, D. & Enquist, L. (1977) *Gene* 1, 255–280.
7. Blattner, F. R., Williams, B. G., Blechl, A. E., Denniston-Thompson, K., Faber, H. E., Furlong, L.-A., Grunwald, D. J., Kiefer, D. O., Moore, D. D., Schumm, J. W., Sheldon, E. L. & Smithies, O. (1977) *Science* 196, 161–169.
8. Hohn, B. & Murray, K. (1977) *Proc. Natl. Acad. Sci. USA* 74, 3259–3263.
9. Maniatis, T., Hardison, R. C., Lacy, E., Lauer, J., O'Connell, C., Quon, D., Sim, G. K. & Efstratiadis, A. (1978) *Cell* 15, 687–701.
10. Sala-Trepat, J. M., Dever, J., Sargent, T. D., Thomas, K., Sell, S. & Bonner, J. (1979) *Biochemistry*, in press.
11. Higuchi, R., Paddock, G. V., Wall, R. & Salser, W. (1976) *Proc. Natl. Acad. Sci. USA* 73, 3146–3150.
12. Roychoudhury, R., Jay, E. & Wu, R. (1976) *Nucleic Acids Res.* 3, 101–116.
13. Curtiss, R., III, Pereira, D. A., Hsu, J. C., Hull, S. C., Clarke, J. E., Maturin, L. J., Sr., Goldsmith, R., Moody, R., Inoue, M. & Alexander, L. (1977) in *Proceedings of the 10th Miles International Symposium*, eds. Beers, R. F., Jr. & Bassett, E. G. (Raven, New York), pp. 45–56.
14. Villa-Komaroff, L., Efstratiadis, A., Broome, S., Lomedico, P., Tizard, R., Naber, S. P., Chick, W. L. & Gilbert, W. (1978) *Proc. Natl. Acad. Sci. USA* 75, 3727–3731.
15. Grunstein, M. & Hogness, D. S. (1975) *Proc. Natl. Acad. Sci. USA* 72, 3961–3965.
16. Isemura, S. & Ikenaka, T. (1978) *J. Biochem.* 83, 35–48.
17. Benton, W. D. & Davis, R. W. (1977) *Science* 196, 180–182.
18. White, R. L. & Hogness, D. S. (1977) *Cell* 10, 177–192.
19. Davis, R. W., Simon, M. & Davidson, N. (1971) *Methods Enzymol.* 21D, 413–428.
20. Southern, E. (1975) *J. Mol. Biol.* 98, 503–517.
21. Denhardt, D. T. (1966) *Biochem. Biophys. Res. Commun.* 23, 641–646.
22. Maniatis, T., Jeffrey, A. & Kleid, D. G. (1975) *Proc. Natl. Acad. Sci. USA* 72, 1184–1188.
23. Dugaiczky, A., Woo, S. L. C., Lai, E. C., Mace, M. L., Jr., McReynolds, L. & O'Malley, B. W. (1978) *Nature (London)* 274, 328–333.
24. Tilghman, S. M., Curtiss, P. J., Tiemeier, D. C., Leder, P. & Weissman, C. (1978) *Proc. Natl. Acad. Sci. USA* 75, 1309–1313.
25. Catterall, J. F., Stein, J. P., Lai, E. C., Woo, S. L. C., Dugaiczky, A., Mace, M. L., Means, A. R. & O'Malley, B. W. (1979) *Nature (London)* 278, 323–327.
26. Strair, R. K., Yap, S. H., Nadal-Ginard, B. & Shafritz, D. A. (1977) *J. Biol. Chem.* 253, 1328–1331.