

## Splice junctions follow a 205-base ladder

(intron/exon/chromatin/higher order structure)

J. S. BECKMANN\* AND E. N. TRIFONOV†

\*Department of Plant Genetics and Breeding, Institute of Field and Garden Crops, Agricultural Research Organization, The Volcani Center, P.O. Box 6, Bet Dagan 50250, Israel; and †Department of Polymer Research, The Weizmann Institute of Science, Rehovot 76100, Israel

Communicated by Donald M. Crothers, December 19, 1990 (received for review March 26, 1990)

**ABSTRACT** Of all aspects of mRNA maturation the accuracy of intervening sequence excision and exon ligation is, perhaps, the most enigmatic. Attempts to identify the essential elements involved in this process have thus far not yielded any satisfactory answer as to what structural (sequence) features are prerequisite for the vital precision of this process. In our search for underlying structural orders we asked whether exons and introns had any positional preferences within a gene. This analysis led to the unexpected discovery that the DNA length is synchronized between successive 3' splicing sites as well as between successive 5' splicing sites, with a frame of  $\approx 205$  base pairs. This observation reveals additional organization of genes in eukaryotes and, perhaps, links gene splicing with chromatin structure.

Several properties differentiate eukaryote from prokaryote cell; notable among these differences are chromatin structure and gene splicing. Trying to find a rationale governing gene splicing in eukaryotes, we realized that chromatin structure and gene splicing might somehow be related—for example, the necessity for better protection of exons and exon–intron boundaries might predicate their involvement in the nucleosomes. On the other hand, giving priority to chromatin and, particularly, to its higher-order structure, some adjustments in gene sequences would be expected—e.g., introduction of introns to allow formation of appropriate higher-order structure (1). Nucleosome positions are known to be largely DNA sequence-dependent (2); their relative orientations in space appear dictated by sequence as well (3, 4). The sequences of uninterrupted structural genes with their coding capacity might not always be compatible with the complex set of rules underlying specific nucleosome positioning (5). We believe sequence requirements of the higher-order chromatin structure to be equally important in this respect. Thus introns, by interrupting the coding sequences at certain positions, could render the protein-coding and chromatin sequence requirements compatible. Once introduced, the intervening sequences might then serve other purposes as well. We, therefore, examined intragenic structural relationships of exon–intron boundaries to check whether these boundaries follow any positional preferences, perhaps reflecting a splicing–chromatin connection. In this communication we report that spliced genes do obey a regular spatial organization; such a connection is thus indicated.

### MATERIALS AND METHODS

The collection of human and murine sequences examined was retrieved from the European Molecular Biology Laboratory data base (release 15). The following intragenic distances were recorded: (i) from every 5' junction site to its downstream 3' junction site (intron length) and to the nearest and

all subsequent downstream 5' junction sites; (ii) from every 3' junction site to the next 5' junction site (exon length) and to the nearest and all subsequent downstream 3' junction sites. Entries were carefully examined to avoid duplications. All genes of simple reiterated motifs, such as collagen, elastase, as well as pseudogenes, were also discarded. The resulting human data base consisted of 155 genes and a total of 530 exons and 518 introns. Splice sites and their positions were used as listed in the feature tables of the sequences. Only a small proportion of splice sites in the data base are authentic—i.e., confirmed by mRNA or protein sequences, whereas others (consensus sites) are accepted as they were estimated by the corresponding authors. The data were analyzed in bulk (grouping authentic splice sites and questionable consensus splice sites together). Frequency distributions were obtained by counting the sequences in groups of increased lengths; successive groups differed from one another by 20 bases. Histograms were generated by “smoothing” (see below) the data.

### RESULTS AND DISCUSSION

Histograms indicate the measured distances separating like 5'–5' (Fig. 1A) and 3'–3' (Fig. 1B) ends of the human sequences. A rather regular profile is seen in both cases, which survives after summation of the two histograms (Fig. 1C), thus indicating a common periodical component. Notwithstanding the noise, maxima appear periodically, approximately every 200–210 base pairs (bp) as far as 2000 bp downstream. The first two peaks are most pronounced. The decreased amplitude as a function of increased length is but a trivial consequence of the fact that more short genes than long genes occur in the data base. Initially these results were obtained on a much smaller data set than that used for Fig. 1. Because of the unexpected character of the observed regularity, we examined a larger collection of sequences, which confirmed our first observation. Moreover, the same results were obtained when samples of the data were examined and counted by hand (to exclude any flaws in the computer program) and when the analysis involved batching the sequences at increments of 10 bases instead of 20 bases and smoothing every 10 bases, within windows of 60 and 90 bases, rather than 100 bases (data not shown).

A murine data set of comparable size was assembled and analyzed similarly; the same precautions were used as described for human sequences. Examination of the murine data set produced the results shown in Fig. 2 A–C, where, once again, regular underlying periodical patterns could be seen for 5'–5' and 3'–3' distances, with peaks occurring every 200–210 bases. The underlying periodicity, common for the human and mouse sequences, is seen most clearly when all histograms (human 5'–5', 3'–3' and mouse 5'–5', 3'–3') are summed (Fig. 3). Fourier analysis of these profiles yields periods of 205, 208, and 206 bases, indistinguishable within error bar ( $\pm 2$ ), for the human, murine, and total data sets, respectively. The calculated amplitude of this oscillation is

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.

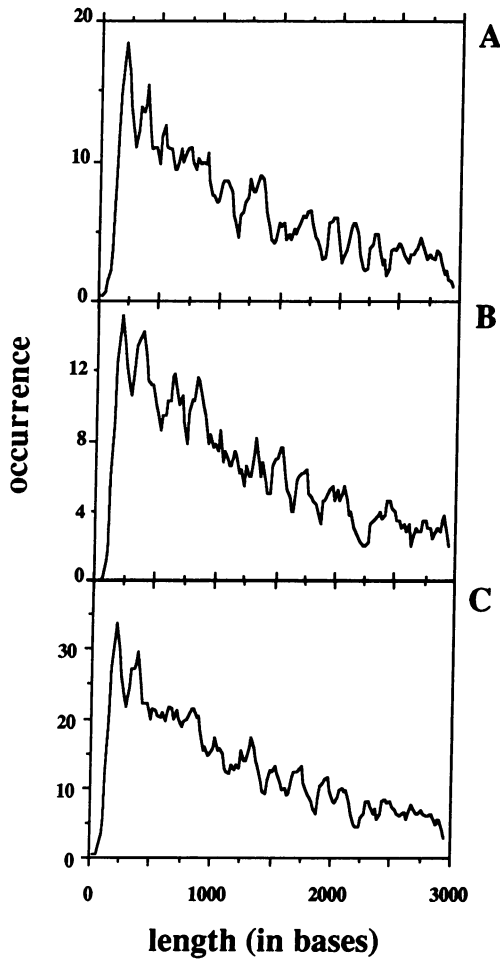


FIG. 1. Histograms of the distances separating like ends 5'-5' (A) and 3'-3' (B) and the sum of these two (C) from human gene sequences retrieved from European Molecular Biology Library; (release 15; see below). Distances within a gene were recorded from every position to its next downstream like position and subsequent like positions. Distances up to 3000 bases were used in these analyses; they were grouped in size clusters with increments of 20 and counted. The histograms were derived upon smoothing in 20-base steps with a window of 100 bases—thus, the figures along the ordinate correspond to a smoothed average of actual occurrences of the distances within 20 bases. Human entries used in analysis (EMBL release 15) include the following: HS242G2, HSA1ATP, HSACCYBB, HSACHR2, HSACTGA, HSACTH, HSADAG, HSAFL12, HSAFL34, HSAFPCP, HSAGL1, HSALBGC, HSALDC, HSALP5, HSAMYA2, HSANFG1, HSANTCD8, HSAPC3A, HSAPOA2G, HSAPOA4A, HSAPOA1A, HSAPOC2G, HSAPRT, HSB2M2, HSBGL3, HSBGPG, HSBLYM1, HSBNGF, HSBSF2, HSCATF, HSCATG7, HSCD14G, HSCD2G1, HSCEANCA, HSCFOS, HSCFVII, HSCFXII3, HSCG01, HSCG1A10, HSCRELA, HSCRPG, HSCRYGX1, HSCS1, HSCSFGMA, HSCYP450, HSDHFR01, HSERBBR, HSERPA, HSEV15VK, HSFABP, HSFBRGG, HSFERG2, HSFESFPS, HSFIBEDA, HSFIXG, HSFN3A, HSFUR1, HSFVII, HSGAS1, HSGASTA, HSGCSFG, HSLUCG2, HSGMCSFG, HSGO1, HSGSHPXG, HSGSTPIG, HSHDCB, HSHER2B, HSHL07, HSHLAB27, HSHLADPB, HSHLADZA, HSHLASBA, HSHLIA, HSHP201, HSHPARS1, HSHSC70, HSHSP27, HSHST, HSHIAIG2, HSHIAIG4, HSHIG05, HSHIGF2G, HSHIGGC3, HSHIGHAM, HSHIGHBC, HSHIGHTC2, HSHIGK15, HSHGLPAV, HSHIGV126, HSHIL05, HSHIL1AG, HSHIL1B, HSHIL2B, HSHINSR, HSHINT1G, HSKER672, HSKER673, HSKIN10, HSLACTG, HSLCATG, HSLMYC1, HSMDR1G, HSMED, HSMETIE, HSMHC1, HSMHCGE2, HSMHCP42, HSMHDC3B, HSMHSXA, HSMIS, HSMPO1, HSMT1B1, HSMYCC, HSNFLG, HSNMYC, HSOPS, HSOTNPI, HSPCRF, HSPLAPL, HSPLPSPC, HSPPPA, HSPRCA, HSPRPH1, HSREN03, HSREN04, HSREP10, HSRPBG1, HSRPS14, HSSAA, HSSB3B46, HSSBA2P, HSSLIPG, HS-

≈4 SDs over background, which can also be estimated from the actual occurrences presented in Fig. 3.

Both intron and exon lengths are known to vary extensively. The distributions of both show single maxima around 110 bp tailing up to as much as 1500 bp for exons and to several thousand bp for introns (data not shown; see also ref. 6). Nothing in these single-peak distributions would indicate the effect presented in Figs. 1–3. Even when single exon-intron pairs are taken (i.e., the nearest 5'-5' and 3'-3' pairs only), a second peak is observed (Fig. 4), which must reflect the basic periodic organization of the interrupted genes.

Remarkably, both 5' and 3' junctions, apparently independently from one another, follow the same ladder. Several peaks are seen in both cases at multiples of ≈205 bases, except for the peak at position three, which fails to appear in the total sum. Recall that all splice sites (both authentic and consensus) listed in the European Molecular Biology Laboratory data bank were used in this analysis. The fact that not all claimed sites might be *bona fide* splice sites would only blur the image and, hence, reduce the effect. Thus, the observed periodicity indicates that the actual effect is probably even stronger than reported here.

Apart from the periodicity, perhaps even more surprising is the persistence of the pattern over long distances, some 2000 bases away. These results seem to imply that splicing junctions, or some structural context thereof, have a memory of the presence and position of other neighboring junctions, being preferable in  $n \times 205$ -base phase with them. Yet, we emphasize that positioning of splice sites on an almost equidistant basis is far from the general rule for all junctions. Many short exons are followed by either a very short or a very long intron (and *vice versa*), thus not conforming to the pattern described. Thus, these facts did not lead to the *a priori* anticipation of such regular organization. In this context the preferential 205-base frame comes as a surprise.

These observations definitely deserve explanation, and several hypotheses come to mind. The most attractive hypothesis is suggested by the fact that the 205-base ladders are reminiscent of and quantitatively close to "nucleosomal repeats" (7). Projecting the observation on other eukaryotes, in general, is difficult because there is insufficient data to analyze genes of other organisms. Those organisms with atypical chromatin (yeast) or splicing organization (yeast, *Drosophila*) might lack the regularity seen in human and mouse genes. Should exon-intron junctions reflect a nucleosomal organization? One possible connection is the protective role nucleosomes might assume while packing the splice junctions. Splice sites differ from many other structural and regulatory DNA sequence elements by their stringent exactitude; misfiring of a single base in one or another of the canonical GU and AG ends of the introns would be disastrous for the entire molecule of mature mRNA. Splice sites are therefore extremely vulnerable to mutations. The great care through which these sequences are protected from mutations can also be seen in the elegant study of ref. 8. It thus seems likely that means to protect these sites from mutational hazards might have been devised throughout evolution, in particular, by encompassing these sites within a more protective environment, such as found in nucleosomes. This strategy, of course, would result in preferential positioning of the splice sites at the nucleosome repeat distance—i.e., avoiding the lengths that would put one of the junctions in a linker region.

SOMI, HSTCBCB, HSTCBV85, HSTCBYT, HSTCRAC, HSTCRT3D, HSTGFBG4, HSTHB, HSTHY1A, HSTHYR8, HSTKRA, HSTNFAB, HSTNFB, HSTPA, HSTUBAG, HSTUBB2, HSUK, HSUPA, HSVK02, HVPNP, HSYG01, and HSYUBGI.

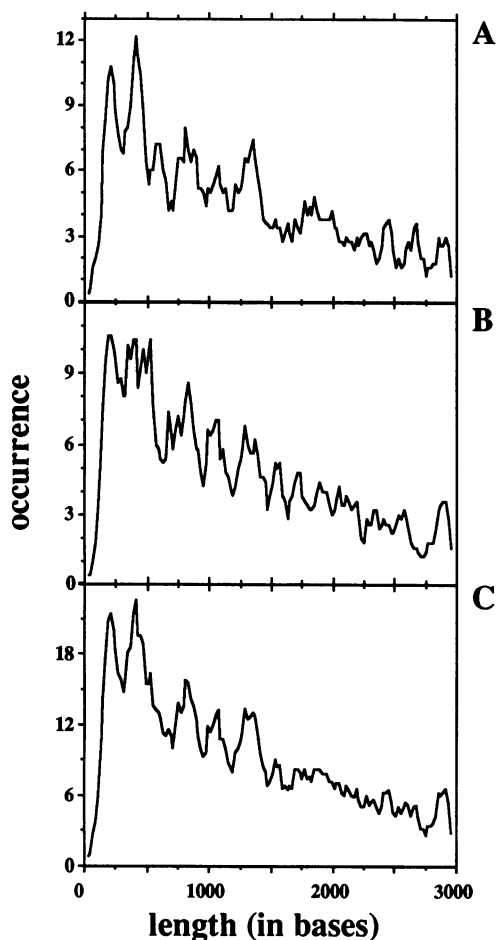


FIG. 2. Histograms of the distances separating like ends 5'-5' (A) and 3'-3' (B) and the sum of these two (C) from murine gene sequences retrieved from European Molecular Biology Library (release 15). Distances within a gene were recorded and analyzed as explained for Fig. 1. Mouse entries used in analysis (EMBL release 15) include the following: MMABLC1A, MMACASA, MMADIG, MMALDH1, MMAMY2, MMANF, MMANT12, MMAPOIVA, MMAPRT, MMASP, MMASPAT6, MMBAND31, MMCAIIN, MMCD42, MMCD43, MMCD46, MMCFOS, MMCRYG2D, MMCRYG3, MMCRYG4, MMCRYG5, MMCRYG6, MMCRYG7, MMCRYG8, MMCRYG9, MMCRYG10, MMCRYG11, MMCRYG12, MMCRYG13, MMCRYG14, MMCRYG15, MMCRYG16, MMCRYG17, MMCRYG18, MMCRYG19, MMCRYG20, MMCRYG21, MMCRYG22, MMCRYG23, MMCRYG24, MMCRYG25, MMCRYG26, MMCRYG27, MMCRYG28, MMCRYG29, MMCRYG30, MMCRYG31, MMCRYG32, MMCRYG33, MMCRYG34, MMCRYG35, MMCRYG36, MMCRYG37, MMCRYG38, MMCRYG39, MMCRYG40, MMCRYG41, MMCRYG42, MMCRYG43, MMCRYG44, MMCRYG45, MMCRYG46, MMCRYG47, MMCRYG48, MMCRYG49, MMCRYG50, MMCRYG51, MMCRYG52, MMCRYG53, MMCRYG54, MMCRYG55, MMCRYG56, MMCRYG57, MMCRYG58, MMCRYG59, MMCRYG60, MMCRYG61, MMCRYG62, MMCRYG63, MMCRYG64, MMCRYG65, MMCRYG66, MMCRYG67, MMCRYG68, MMCRYG69, MMCRYG70, MMCRYG71, MMCRYG72, MMCRYG73, MMCRYG74, MMCRYG75, MMCRYG76, MMCRYG77, MMCRYG78, MMCRYG79, MMCRYG80, MMCRYG81, MMCRYG82, MMCRYG83, MMCRYG84, MMCRYG85, MMCRYG86, MMCRYG87, MMCRYG88, MMCRYG89, MMCRYG90, MMCRYG91, MMCRYG92, MMCRYG93, MMCRYG94, MMCRYG95, MMCRYG96, MMCRYG97, MMCRYG98, MMCRYG99, MMCRYG100, MMCRYG101, MMCRYG102, MMCRYG103, MMCRYG104, MMCRYG105, MMCRYG106, MMCRYG107, MMCRYG108, MMCRYG109, MMCRYG110, MMCRYG111, MMCRYG112, MMCRYG113, MMCRYG114, MMCRYG115, MMCRYG116, MMCRYG117, MMCRYG118, MMCRYG119, MMCRYG120, MMCRYG121, MMCRYG122, MMCRYG123, MMCRYG124, MMCRYG125, MMCRYG126, MMCRYG127, MMCRYG128, MMCRYG129, MMCRYG130, MMCRYG131, MMCRYG132, MMCRYG133, MMCRYG134, MMCRYG135, MMCRYG136, MMCRYG137, MMCRYG138, MMCRYG139, MMCRYG140, MMCRYG141, MMCRYG142, MMCRYG143, MMCRYG144, MMCRYG145, MMCRYG146, MMCRYG147, MMCRYG148, MMCRYG149, MMCRYG150, MMCRYG151, MMCRYG152, MMCRYG153, MMCRYG154, MMCRYG155, MMCRYG156, MMCRYG157, MMCRYG158, MMCRYG159, MMCRYG160, MMCRYG161, MMCRYG162, MMCRYG163, MMCRYG164, MMCRYG165, MMCRYG166, MMCRYG167, MMCRYG168, MMCRYG169, MMCRYG170, MMCRYG171, MMCRYG172, MMCRYG173, MMCRYG174, MMCRYG175, MMCRYG176, MMCRYG177, MMCRYG178, MMCRYG179, MMCRYG180, MMCRYG181, MMCRYG182, MMCRYG183, MMCRYG184, MMCRYG185, MMCRYG186, MMCRYG187, MMCRYG188, MMCRYG189, MMCRYG190, MMCRYG191, MMCRYG192, MMCRYG193, MMCRYG194, MMCRYG195, MMCRYG196, MMCRYG197, MMCRYG198, MMCRYG199, MMCRYG200, MMCRYG201, MMCRYG202, MMCRYG203, MMCRYG204, MMCRYG205, MMCRYG206, MMCRYG207, MMCRYG208, MMCRYG209, MMCRYG210, MMCRYG211, MMCRYG212, MMCRYG213, MMCRYG214, MMCRYG215, MMCRYG216, MMCRYG217, MMCRYG218, MMCRYG219, MMCRYG220, MMCRYG221, MMCRYG222, MMCRYG223, MMCRYG224, MMCRYG225, MMCRYG226, MMCRYG227, MMCRYG228, MMCRYG229, MMCRYG230, MMCRYG231, MMCRYG232, MMCRYG233, MMCRYG234, MMCRYG235, MMCRYG236, MMCRYG237, MMCRYG238, MMCRYG239, MMCRYG240, MMCRYG241, MMCRYG242, MMCRYG243, MMCRYG244, MMCRYG245, MMCRYG246, MMCRYG247, MMCRYG248, MMCRYG249, MMCRYG250, MMCRYG251, MMCRYG252, MMCRYG253, MMCRYG254, MMCRYG255, MMCRYG256, MMCRYG257, MMCRYG258, MMCRYG259, MMCRYG260, MMCRYG261, MMCRYG262, MMCRYG263, MMCRYG264, MMCRYG265, MMCRYG266, MMCRYG267, MMCRYG268, MMCRYG269, MMCRYG270, MMCRYG271, MMCRYG272, MMCRYG273, MMCRYG274, MMCRYG275, MMCRYG276, MMCRYG277, MMCRYG278, MMCRYG279, MMCRYG280, MMCRYG281, MMCRYG282, MMCRYG283, MMCRYG284, MMCRYG285, MMCRYG286, MMCRYG287, MMCRYG288, MMCRYG289, MMCRYG290, MMCRYG291, MMCRYG292, MMCRYG293, MMCRYG294, MMCRYG295, MMCRYG296, MMCRYG297, MMCRYG298, MMCRYG299, MMCRYG300, MMCRYG301, MMCRYG302, MMCRYG303, MMCRYG304, MMCRYG305, MMCRYG306, MMCRYG307, MMCRYG308, MMCRYG309, MMCRYG310, MMCRYG311, MMCRYG312, MMCRYG313, MMCRYG314, MMCRYG315, MMCRYG316, MMCRYG317, MMCRYG318, MMCRYG319, MMCRYG320, MMCRYG321, MMCRYG322, MMCRYG323, MMCRYG324, MMCRYG325, MMCRYG326, MMCRYG327, MMCRYG328, MMCRYG329, MMCRYG330, MMCRYG331, MMCRYG332, MMCRYG333, MMCRYG334, MMCRYG335, MMCRYG336, MMCRYG337, MMCRYG338, MMCRYG339, MMCRYG340, MMCRYG341, MMCRYG342, MMCRYG343, MMCRYG344, MMCRYG345, MMCRYG346, MMCRYG347, MMCRYG348, MMCRYG349, MMCRYG350, MMCRYG351, MMCRYG352, MMCRYG353, MMCRYG354, MMCRYG355, MMCRYG356, MMCRYG357, MMCRYG358, MMCRYG359, MMCRYG360, MMCRYG361, MMCRYG362, MMCRYG363, MMCRYG364, MMCRYG365, MMCRYG366, MMCRYG367, MMCRYG368, MMCRYG369, MMCRYG370, MMCRYG371, MMCRYG372, MMCRYG373, MMCRYG374, MMCRYG375, MMCRYG376, MMCRYG377, MMCRYG378, MMCRYG379, MMCRYG380, MMCRYG381, MMCRYG382, MMCRYG383, MMCRYG384, MMCRYG385, MMCRYG386, MMCRYG387, MMCRYG388, MMCRYG389, MMCRYG390, MMCRYG391, MMCRYG392, MMCRYG393, MMCRYG394, MMCRYG395, MMCRYG396, MMCRYG397, MMCRYG398, MMCRYG399, MMCRYG400, MMCRYG401, MMCRYG402, MMCRYG403, MMCRYG404, MMCRYG405, MMCRYG406, MMCRYG407, MMCRYG408, MMCRYG409, MMCRYG410, MMCRYG411, MMCRYG412, MMCRYG413, MMCRYG414, MMCRYG415, MMCRYG416, MMCRYG417, MMCRYG418, MMCRYG419, MMCRYG420, MMCRYG421, MMCRYG422, MMCRYG423, MMCRYG424, MMCRYG425, MMCRYG426, MMCRYG427, MMCRYG428, MMCRYG429, MMCRYG430, MMCRYG431, MMCRYG432, MMCRYG433, MMCRYG434, MMCRYG435, MMCRYG436, MMCRYG437, MMCRYG438, MMCRYG439, MMCRYG440, MMCRYG441, MMCRYG442, MMCRYG443, MMCRYG444, MMCRYG445, MMCRYG446, MMCRYG447, MMCRYG448, MMCRYG449, MMCRYG450, MMCRYG451, MMCRYG452, MMCRYG453, MMCRYG454, MMCRYG455, MMCRYG456, MMCRYG457, MMCRYG458, MMCRYG459, MMCRYG460, MMCRYG461, MMCRYG462, MMCRYG463, MMCRYG464, MMCRYG465, MMCRYG466, MMCRYG467, MMCRYG468, MMCRYG469, MMCRYG470, MMCRYG471, MMCRYG472, MMCRYG473, MMCRYG474, MMCRYG475, MMCRYG476, MMCRYG477, MMCRYG478, MMCRYG479, MMCRYG480, MMCRYG481, MMCRYG482, MMCRYG483, MMCRYG484, MMCRYG485, MMCRYG486, MMCRYG487, MMCRYG488, MMCRYG489, MMCRYG490, MMCRYG491, MMCRYG492, MMCRYG493, MMCRYG494, MMCRYG495, MMCRYG496, MMCRYG497, MMCRYG498, MMCRYG499, MMCRYG500, MMCRYG501, MMCRYG502, MMCRYG503, MMCRYG504, MMCRYG505, MMCRYG506, MMCRYG507, MMCRYG508, MMCRYG509, MMCRYG510, MMCRYG511, MMCRYG512, MMCRYG513, MMCRYG514, MMCRYG515, MMCRYG516, MMCRYG517, MMCRYG518, MMCRYG519, MMCRYG520, MMCRYG521, MMCRYG522, MMCRYG523, MMCRYG524, MMCRYG525, MMCRYG526, MMCRYG527, MMCRYG528, MMCRYG529, MMCRYG530, MMCRYG531, MMCRYG532, MMCRYG533, MMCRYG534, MMCRYG535, MMCRYG536, MMCRYG537, MMCRYG538, MMCRYG539, MMCRYG540, MMCRYG541, MMCRYG542, MMCRYG543, MMCRYG544, MMCRYG545, MMCRYG546, MMCRYG547, MMCRYG548, MMCRYG549, MMCRYG550, MMCRYG551, MMCRYG552, MMCRYG553, MMCRYG554, MMCRYG555, MMCRYG556, MMCRYG557, MMCRYG558, MMCRYG559, MMCRYG560, MMCRYG561, MMCRYG562, MMCRYG563, MMCRYG564, MMCRYG565, MMCRYG566, MMCRYG567, MMCRYG568, MMCRYG569, MMCRYG570, MMCRYG571, MMCRYG572, MMCRYG573, MMCRYG574, MMCRYG575, MMCRYG576, MMCRYG577, MMCRYG578, MMCRYG579, MMCRYG580, MMCRYG581, MMCRYG582, MMCRYG583, MMCRYG584, MMCRYG585, MMCRYG586, MMCRYG587, MMCRYG588, MMCRYG589, MMCRYG590, MMCRYG591, MMCRYG592, MMCRYG593, MMCRYG594, MMCRYG595, MMCRYG596, MMCRYG597, MMCRYG598, MMCRYG599, MMCRYG600, MMCRYG601, MMCRYG602, MMCRYG603, MMCRYG604, MMCRYG605, MMCRYG606, MMCRYG607, MMCRYG608, MMCRYG609, MMCRYG610, MMCRYG611, MMCRYG612, MMCRYG613, MMCRYG614, MMCRYG615, MMCRYG616, MMCRYG617, MMCRYG618, MMCRYG619, MMCRYG620, MMCRYG621, MMCRYG622, MMCRYG623, MMCRYG624, MMCRYG625, MMCRYG626, MMCRYG627, MMCRYG628, MMCRYG629, MMCRYG630, MMCRYG631, MMCRYG632, MMCRYG633, MMCRYG634, MMCRYG635, MMCRYG636, MMCRYG637, MMCRYG638, MMCRYG639, MMCRYG640, MMCRYG641, MMCRYG642, MMCRYG643, MMCRYG644, MMCRYG645, MMCRYG646, MMCRYG647, MMCRYG648, MMCRYG649, MMCRYG650, MMCRYG651, MMCRYG652, MMCRYG653, MMCRYG654, MMCRYG655, MMCRYG656, MMCRYG657, MMCRYG658, MMCRYG659, MMCRYG660, MMCRYG661, MMCRYG662, MMCRYG663, MMCRYG664, MMCRYG665, MMCRYG666, MMCRYG667, MMCRYG668, MMCRYG669, MMCRYG670, MMCRYG671, MMCRYG672, MMCRYG673, MMCRYG674, MMCRYG675, MMCRYG676, MMCRYG677, MMCRYG678, MMCRYG679, MMCRYG680, MMCRYG681, MMCRYG682, MMCRYG683, MMCRYG684, MMCRYG685, MMCRYG686, MMCRYG687, MMCRYG688, MMCRYG689, MMCRYG690, MMCRYG691, MMCRYG692, MMCRYG693, MMCRYG694, MMCRYG695, MMCRYG696, MMCRYG697, MMCRYG698, MMCRYG699, MMCRYG700, MMCRYG701, MMCRYG702, MMCRYG703, MMCRYG704, MMCRYG705, MMCRYG706, MMCRYG707, MMCRYG708, MMCRYG709, MMCRYG710, MMCRYG711, MMCRYG712, MMCRYG713, MMCRYG714, MMCRYG715, MMCRYG716, MMCRYG717, MMCRYG718, MMCRYG719, MMCRYG720, MMCRYG721, MMCRYG722, MMCRYG723, MMCRYG724, MMCRYG725, MMCRYG726, MMCRYG727, MMCRYG728, MMCRYG729, MMCRYG730, MMCRYG731, MMCRYG732, MMCRYG733, MMCRYG734, MMCRYG735, MMCRYG736, MMCRYG737, MMCRYG738, MMCRYG739, MMCRYG740, MMCRYG741, MMCRYG742, MMCRYG743, MMCRYG744, MMCRYG745, MMCRYG746, MMCRYG747, MMCRYG748, MMCRYG749, MMCRYG750, MMCRYG751, MMCRYG752, MMCRYG753, MMCRYG754, MMCRYG755, MMCRYG756, MMCRYG757, MMCRYG758, MMCRYG759, MMCRYG760, MMCRYG761, MMCRYG762, MMCRYG763, MMCRYG764, MMCRYG765, MMCRYG766, MMCRYG767, MMCRYG768, MMCRYG769, MMCRYG770, MMCRYG771, MMCRYG772, MMCRYG773, MMCRYG774, MMCRYG775, MMCRYG776, MMCRYG777, MMCRYG778, MMCRYG779, MMCRYG780, MMCRYG781, MMCRYG782, MMCRYG783, MMCRYG784, MMCRYG785, MMCRYG786, MMCRYG787, MMCRYG788, MMCRYG789, MMCRYG790, MMCRYG791, MMCRYG792, MMCRYG793, MMCRYG794, MMCRYG795, MMCRYG796, MMCRYG797, MMCRYG798, MMCRYG799, MMCRYG800, MMCRYG801, MMCRYG802, MMCRYG803, MMCRYG804, MMCRYG805, MMCRYG806, MMCRYG807, MMCRYG808, MMCRYG809, MMCRYG810, MMCRYG811, MMCRYG812, MMCRYG813, MMCRYG814, MMCRYG815, MMCRYG816, MMCRYG817, MMCRYG818, MMCRYG819, MMCRYG820, MMCRYG821, MMCRYG822, MMCRYG823, MMCRYG824, MMCRYG825, MMCRYG826, MMCRYG827, MMCRYG828, MMCRYG829, MMCRYG830, MMCRYG831, MMCRYG832, MMCRYG833, MMCRYG834, MMCRYG835, MMCRYG836, MMCRYG837, MMCRYG838, MMCRYG839, MMCRYG840, MMCRYG841, MMCRYG842, MMCRYG843, MMCRYG844, MMCRYG845, MMCRYG846, MMCRYG847, MMCRYG848, MMCRYG849, MMCRYG850, MMCRYG851, MMCRYG852, MMCRYG853, MMCRYG854, MMCRYG855, MMCRYG856, MMCRYG857, MMCRYG858, MMCRYG859, MMCRYG860, MMCRYG861, MMCRYG862, MMCRYG863, MMCRYG864, MMCRYG865, MMCRYG866, MMCRYG867, MMCRYG868, MMCRYG869, MMCRYG870, MMCRYG871, MMCRYG872, MMCRYG873, MMCRYG874, MMCRYG875, MMCRYG876, MMCRYG877, MMCRYG878, MMCRYG879, MMCRYG880, MMCRYG881, MMCRYG882, MMCRYG883, MMCRYG884, MMCRYG885, MMCRYG886, MMCRYG887, MMCRYG888, MMCRYG889, MMCRYG890, MMCRYG891, MMCRYG892, MMCRYG893, MMCRYG894, MMCRYG895, MMCRYG896, MMCRYG897, MMCRYG898, MMCRYG899, MMCRYG900, MMCRYG901, MMCRYG902, MMCRYG903, MMCRYG904, MMCRYG905, MMCRYG906, MMCRYG907, MMCRYG908, MMCRYG909, MMCRYG910, MMCRYG911, MMCRYG912, MMCRYG913, MMCRYG914, MMCRYG915, MMCRYG916, MMCRYG917, MMCRYG918, MMCRYG919, MMCRYG920, MMCRYG921, MMCRYG922, MMCRYG923, MMCRYG924, MMCRYG925, MMCRYG926, MMCRYG927, MMCRYG928, MMCRYG929, MMCRYG930, MMCRYG931, MMCRYG932, MMCRYG933, MMCRYG934, MMCRYG935, MMCRYG936, MMCRYG937, MMCRYG938, MMCRYG939, MMCRYG940, MMCRYG941, MMCRYG942, MMCRYG943, MMCRYG944, MMCRYG945, MMCRYG946, MMCRYG947, MMCRYG948, MMCRYG949, MMCRYG950, MMCRYG951, MMCRYG952, MMCRYG953, MMCRYG954, MMCRYG955, MMCRYG956, MMCRYG957, MMCRYG958, MMCRYG959, MMCRYG960, MMCRYG961, MMCRYG962, MMCRYG963, MMCRYG964, MMCRYG965, MMCRYG966, MMCRYG967, MMCRYG968, MMCRYG969, MMCRYG970, MMCRYG971, MMCRYG972, MMCRYG973, MMCRYG974, MMCRYG975, MMCRYG976, MMCRYG977, MMCRYG978, MMCRYG979, MMCRYG980, MMCRYG981, MMCRYG982, MMCRYG983, MMCRYG984, MMCRYG985, MMCRYG986, MMCRYG987, MMCRYG988, MMCRYG989, MMCRYG990, MMCRYG991, MMCRYG992, MMCRYG993, MMCRYG994, MMCRYG995, MMCRYG996, MMCRYG997, MMCRYG998, MMCRYG999, MMCRYG1000.

In fact, the extreme relative vulnerability of splice junctions leads one to wonder why apparently useless introns were not lost during evolution. The contrary seems to have occurred. We are then forced to conclude that introns are far from ancillary and probably occupy some essential functions. We believe that this function may also be related to a structural requirement for chromatin folding. Both nucleosome positions along DNA and their folding in space appear largely dictated by the nucleotide sequences (2, 4). Consequently, the arrangement of nucleosomes, organized by the

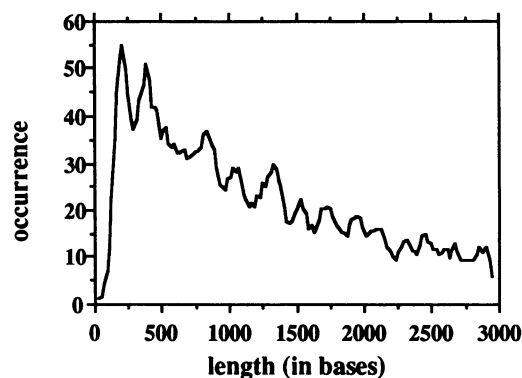


FIG. 3. Cumulative integrated histogram of the distances separating all human and murine intron ends (5'-5' and 3'-3' together).

genomic sequence with alternating exons and introns, has to be substantially different from that expected were noninterrupted protein-coding DNA sequence to organize the chromatin structure. Being spliced, the mRNA sequence carries neither the originally encoded chromatin design nor any positionally and spatially consistent chromatin design at all, serving primarily in its protein-coding capacity. Thus, gene splicing can be viewed as a device for simultaneously satisfying conflicting protein-coding and chromatin-organization sequence requirements, as originally suggested in earlier studies (1, 5). Nucleosomes and their higher-order organization could constitute one of the synchronizing factors for the observed preferentially regular spacing of splice sites along a gene, as well as the *raison d'être* of introns.

Other phasing mechanisms could also be involved. Units of  $\approx 200$  bp of DNA are seen in other instances: Okazaki pieces are synthesized in segments of about this size before being end-ligated (9). mRNA processing could also structurally involve a kind of processive scanning "measure" that would accommodate  $\approx 200$  nucleotides at a time. The packaging of the mRNA precursors in the heterogeneous nuclear RNP particles might also follow some structural regularity (10) to position the splice junctions properly for efficient and accurate processing. And, of course, the observed regularity could result from joint pressures of all the above size preferences.

Whatever mechanism is considered, it also must account for the fact that rather long-distance peaks are seen as well. We assume the presence of some responsible phasing elements in the DNA sequences, but preliminary analyses have failed to recognize any simple nucleotide signal; perhaps more complex periodically organized patterns await detection.

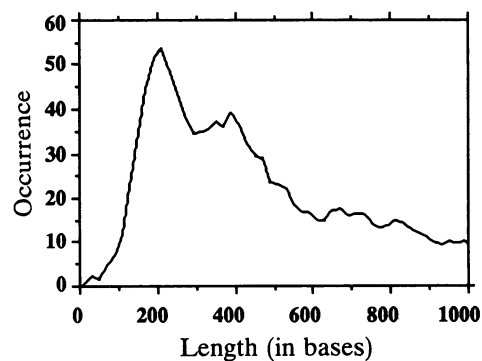


FIG. 4. Histogram of the distances separating single exon-intron pairs, nearest like ends (5'-5' and 3'-3') from the human and murine data sets. Distances within each gene were recorded and analyzed as explained for Fig. 1.

The expert help in data assembly by Dr. Philip Bucher is highly appreciated. This work is contribution no. 2840.E (1989 series) from the Agricultural Research Organization, The Volcani Center, Bet Dagan, Israel.

1. Zuckerkandl, E. (1981) *Mol. Biol. Rep.* **7**, 149–158.
2. Mengeritsky, G. & Trifonov, E. N. (1983) *Nucleic Acids Res.* **11**, 3833–3851.
3. Noll, M., Zimmer, S., Engel, A. & Dubochet, J. (1980) *Nucleic Acids Res.* **8**, 21–42.
4. Ulanovsky, L. E. & Trifonov, E. N. (1986) in *Biomolecular Stereodynamics III*, eds. Sarma, R. H. & Sarma, M. H. (Adenine, Guilderland, NY), pp. 35–44.
5. Soloviev, V. V. & Kolchanov, N. A. (1986) *Dokl. Biochem. (Engl. Transl.)* **284**, 286–290.
6. Hawkins, J. D. (1988) *Nucleic Acids Res.* **16**, 9893–9908.
7. van Holde, K. E. (1988) *Chromatin* (Springer, New York).
8. Smith, C. W. J., Porro, E. B., Patton, J. G. & Nadal-Ginard, B. (1989) *Nature (London)* **242**, 243–247.
9. Ogawa, T. & Okazaki, T. (1980) *Annu. Rev. Biochem.* **49**, 421–457.
10. Barnett, S. F., Friedman, D. L. & LeStourgeon, W. M. (1989) *Mol. Cell. Biol.* **9**, 492–498.