
Concordance of experimentally mapped or predicted Z-DNA sites with positions of selected alternating purine-pyrimidine tracts

Andrzej K.Konopka, Johannes Reiter^{1*}, Manfred Jung², David A.Zarling⁺ and Thomas M.Jovin

Abteilung Molekulare Biologie, Max-Planck-Institut für biophysikalische Chemie, D-3400 Göttingen, FRG, ¹Institut für theoretische Chemie und Strahlchemie der Universität, A-1090 Wien, Austria, and ²Abteilung Biochemische Kinetik, Max-Planck-Institut für biophysikalische Chemie, D-3400 Göttingen, FRG

Received 7 December 1984; Revised and Accepted 4 February 1985

ABSTRACT

The recent electronmicroscopic and biochemical mapping of Z-DNA sites in ϕ X174, SV40, pBR322 and PM2 DNAs has been used to determine two sets of criteria for identification of potential Z-DNA sequences in natural DNA genomes. The prediction of potential Z-DNA tracts and corresponding statistical analysis of their occurrence have been made on a sample of 14 DNA genomes.

Alternating purine and pyrimidine tracts longer than 5 base pairs in length and their clusters (quasi alternating fragments) in the 14 genomes studied are under-represented compared to the expectation from corresponding random sequences. The fragments $[d(G\cdot C)]_n$ and $[d(C\cdot G)]_n$ ($n \geq 3$) in general do not occur in circular DNA genomes and are under-represented in the linear DNAs of phages λ and T7, whereas in linear genomes of adenoviruses they are strongly over-represented. With minor exceptions, potential Z-DNA sites are also under-represented compared to random sequences.

In the 14 genomes studied, predicted Z-DNA tracts occur in non-coding as well as in protein coding regions. The predicted Z-DNA sites in ϕ X174, SV40, pBR322 and PM2 correspond well with those mapped experimentally. A complete listing together with a compact graphical representation of alternating purine-pyrimidine fragments and their Z-forming potential are presented.

INTRODUCTION

The alternation of purines and pyrimidines in DNA sequences constitutes one of the most important factors potentiating the transition from the right-handed B to the left-handed Z helical conformation *in vitro* (1-4).

Topological stress in the form of negative supercoiling promotes the B to Z transition of protein-free covalently closed circular DNA (ccc DNA) (5-12). Studies of anti-Z-DNA-IgG binding to chromosomal DNA (11-16) have established the existence of potential Z-DNA tracts *in vivo*. It is probable that the combined effects of nucleotide sequence, topological stress, and interactions with ions, proteins and polyamines (5, 11, 12, 17) determine the physiological distribution and functions of left-handed DNA (for a review see 4).

The biological significance of Z-DNA is unknown. It has been suggested that some potential Z-DNA loci in the SV40 genome can play a role in the

control of transcription or in genetic recombination (5, 11, 12 18). Other studies with cytological material have emphasized potential structural roles for left-handed DNA in chromosomal organization (4, 11, 12, 15, 16).

Different alternating purine-pyrimidine sequences in linear synthetic polymers exhibit a hierarchy in the potential for undergoing the B-Z transition (4, 20). The minimum length of a linear alternating oligonucleotide required for the establishment of the Z form in solution has been evaluated as 6 base pairs (20). Although the precise sequence-dependence of the transition equilibrium remains to be experimentally established, it is clear that the G·C basepair is much more effective than the A·T basepair in stabilizing the Z conformation (4, 20). Thus, we can identify two primary factors which determine the Z-forming potential of a natural sequence: length and base composition. In addition, studies of anti-Z-DNA-Ig binding sites in pBR322 (8, 19) indicate that DNA fragments including bases out of alternation may also assume the left-handed conformation at high superhelix density. It follows that a favorable (clustered) distribution of potential Z-forming tracts may lead to a cooperative and collective behaviour.

The mapping of anti-Z-DNA immunoglobulin binding sites by immunoelectron microscopy provides several examples of naturally occurring Z-DNA tracts in ϕ X174 (23), PM2 (24-26) and SV40 (18, 27) DNAs. Corresponding data also exist for the cloning vector pBR322 (8). On the basis of the available data, the minimal length of an alternating purine-pyrimidine fragment required for stabilization of the Z conformation in natural sequences is on the order of 8 base-pairs, although shorter tracts composed exclusively of G and C may also be effective (23). The results of these experimental studies have been used to define empirical criteria for identifying potential Z-DNA tracts on the basis of nucleotide sequence (in the next section two working definitions of sequences with the potential for adopting the left-handed conformation are presented). The algorithms based on these definitions have been applied to several viral and episomal DNA genomes. The results of the search together with a corresponding statistical analysis are presented and discussed in this paper.

BASIC PRINCIPLES

Terminology. DNA sequences can be analysed for simple dinucleotide repeating units. In particular, we will consider alternating repetitions of a purine (R) and a pyrimidine (Y). A sequence of purines and pyrimidines in alternation is referred to as an *alternating fragment* (AF). Formally we can

consider two kinds of AFs: those which are alternations of only two bases and those which consist of more than two bases. An AF of the first kind will be referred to as a *uniform alternating fragment* (uAF). Examples of uAFs are the sequences: GTGTGTGTGTG, ACACACAC, CGCGCGCGCG, ATATATATAT. The other kind of AF will be referred to as a *mixed alternating fragment* (mAF). Examples of mAFs are the sequences: GCGTACGT, GCACATGTA, ACACGTACATG, ACGTACGTACGT.

In a long DNA tract, AFs can be separated or clustered. The obvious criterion for establishing whether a block of AFs constitute a cluster is based on the distances between the AFs. From the viewpoint of Z-forming potential, we regard one base-pair as a reasonable maximum distance between AFs in a cluster. A cluster of AFs will be referred to as a *clustered alternating fragment* (cAF). An example is the sequence:

ATACGT TGTGTGT T CGATCGTG. (Here and elsewhere a space will be used to denote the separation between AFs constituting a cAF.) This cluster consists of three AFs. The first two are contiguous (distance = 0). The second and third AFs are separated by one base (distance = 1). Thus, we consider a cAF as equivalent to an AF with a few bases out of perfect alternation. The length of a cAF will be taken as the difference between the positions of the first base of the first AF and of the last base of the last AF in the cluster. Under the assumption that the Z-forming potential can be a property of AFs as well as cAFs, we define a *quasi-alternating fragment* (qAF) as a DNA sequence which is either an AF or a cAF.

Potential Z-DNA. Taking into account the facts briefly described in the Introduction, we should not expect that every qAF has the potential for adopting the Z conformation. In light of available experimental data, the two following definitions of a potential Z-DNA fragment seem to be appropriate. The first definition evaluates the criteria of length and composition for a qAF as a whole, whereas the second considers the length of the longest subfragment of a given qAF.

DEFINITION 1: A potential Z-DNA fragment is a qAF fulfilling the following conditions:

- i. The total length (base-pair units) is $> a$.
- ii. The fraction of the sequence consisting of A and T in alternating repetition is $\leq b$.
- iii. If the fragment is a cAF, the constituent AFs have a length $> c$.

DEFINITION 2: A potential Z-DNA fragment is a qAF fulfilling conditions (i) and (iii) from the previous definition and containing a subfragment which fulfills conditions (i) and (ii).

For the reasons given previously, we have applied the search algorithm with

the parameters $a = 7$, $b = 0.3$, and $c = 4$. These values lead to the identification of the binding sites for anti-Z DNA immunoglobulins reported to date (with minor exceptions; see Discussion).

Let us consider as examples the following qAFs:

Sequence 1: ATATCG TGTGTG GCATATATAT

Sequence 2: GTATATAT TATATAT T CACAC

Sequence 3: CGCGCG T CATGTG ACACACAT

Sequence 4: CACGTATGTGTATATGTGCA

Sequence 5: GTGTA

Sequence 1 is a potential Z-DNA fragment according to definition 2 but not definition 1, due to violation of condition (ii). Sequences 2 and 5 are not potential Z-DNA fragments according to both definitions [sequence 5 violates all conditions whereas sequence 2 violates condition (ii)]. Sequences 3 and 4 are potential Z-DNA fragments according to both definitions.

DNA sequences studied. The DNA sequences chosen from the EMBL Nucleotide Sequence Library are as follows (ssc, single-stranded circular; dsc, double-stranded circular; dsl, double-stranded linear. Lengths are in base or base-pair units): cloning vector pBR322 (dsc 4362); bacteriophages: ϕ X174 (ssc 5386), M13 (ssc 6407), T7(dsl 39936), λ (dsl 48502); papovaviruses: SV40 (simian; dsc 5243), BKV (human, strain Dunlop, dsc 5153), Polyoma-A2 (strain A2; dsc 5292); adenoviruses: Adeno-7·l (type 7; dsl 6707; left 0-18.5%), Adeno-2·l (type 2; dsl 11600; left 0-32%) Adeno-2·r (type 2; dsl 10305; right 70.7-100%); mitochondria: bovine (*Bos taurus*, dsc 16338), murine (*Mus musculus*, dsc 16295), human (dsc 16569). The sequence of a purine-pyrimidine rich region of phage PM2 (dsc 1757) has been taken from reference 26. The single-stranded DNA genomes have, of course, double-stranded DNA replicative intermediates. Linear genomes can circularize during replication.

Occurrence and average length of qAFs in a long DNA sequence. Let us consider the Y/R tracts found in a fragment of the SV40 genome (Fig. 1a-c). Figure 1a shows all AFs not shorter than 4 bases. There are two such fragments in the sequence studied. Both of them are separated AFs. The AFs which are not shorter than 3 bases are shown in Figure 1b. There are four such fragments, two of which make a cluster with two bases out of perfect alternation. Figure 1c shows all possible AFs (including doublets of alternating Ys and Rs) in the DNA fragment studied. In this case 12 AFs are found. Two of them constitute a cluster with two bases out of alternation and another six are involved in a cluster of total length 14.

These examples suggest that in general, a tendency of AFs to cluster is stronger in the case of short AFs (Figure 1c) compared to long AFs (Figure 1a

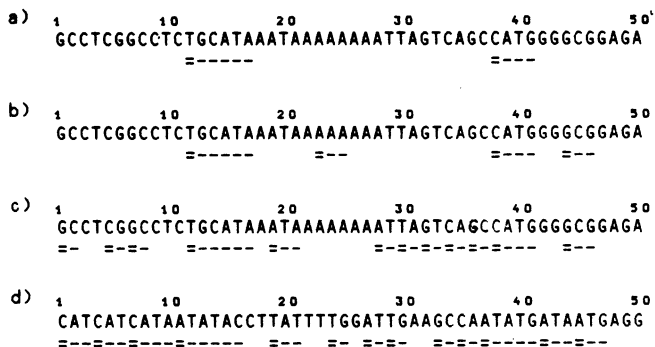


Fig. 1. Alternating Y and R repetitions in the first 50 bases of SV40 (a-c) and Adeno-2·l (d). a) AFs of length ≥ 4 ; b) AFs of length ≥ 3 ; c) and d) AFs of length ≥ 2 . Every AF is underlined. The = indicates the first base of an AF.

and 1b, which have a smaller number of underlined bases than Figure 1c). When we analyse all possible AFs from different natural DNAs, the tendency to cluster differs. An example is provided by a comparison of the SV40 genome fragment (Figure 1c) with a corresponding Adeno-2·l fragment (Figure 1d). There are 13 AFs in the first 50 bases of the Adeno-2·l but only one is clearly isolated. Another 12 AFs are involved in three clusters. We know from the previous example that there are 12 AFs in the first 50 bases of the SV40 genome; two of them are separated and the other 10 occur in three clusters. Thus, the tendency of short AFs to cluster is greater in the Adeno-2·l than in the SV40 fragment. The clusters in Adeno-2·l are, in general, longer than in SV40 (the lengths of the clusters are 16, 7 and 15 in Adeno-2·l whereas these lengths in SV40 are equal to 4, 10 and 14 bases).

We require a quantitative measure of the clustering tendency. Such a measure consists of the average length of qAFs in a given DNA tract. In order to define this quantity let us assume that a sequence under consideration is of length N bases, and that it contains a number k of qAFs. In addition, let n_2, n_3, \dots, n_m be the numbers of qAFs of length 2, 3, ..., m bases, respectively. The average length of a qAF is defined by the following general expression:

$$\langle L \rangle = \sum_{i=2}^m i \cdot p_i = (1/k) \cdot \sum_{i=2}^m i \cdot n_i \quad (1)$$

where n_i is the number of qAFs of length i and the probabilities p_i are equal to n_i/k .

Let us return to our example of the 50 bp regions of SV40 and Adeno-2. We have already pointed out that the tendency of AFs to cluster is greater in the adenovirus than in the SV40 fragment. We compute the $\langle L \rangle$ values for both these fragments by using (1) and distinguish the $\langle L \rangle$ values for AFs as $\langle L_{AF} \rangle$ and the values for qAFs as $\langle L_{qAF} \rangle$. For the SV40 fragment $\langle L_{AF} \rangle = 2.7$ and $\langle L_{qAF} \rangle = 6.6$, whereas the corresponding values for the adenovirus fragment are equal to 3.1 and 10.3. It appears from this calculation that although the two $\langle L_{AF} \rangle$ values are similar, the $\langle L_{qAF} \rangle$ values are considerably different. This result suggests that the quantity $\langle L_{qAF} \rangle - \langle L_{AF} \rangle$ is a good measure of the tendency of AFs to cluster.

Random DNA sequence. Random DNA sequences have been generated and representative examples chosen for Y/R searches. In order to verify our definitions of potential Z-DNA we have generated three categories of random sequences, i.e. those with an equiprobable base composition, those rich in A and T (30% each) and those rich in G and C (30% each).

The expected number of AFs longer or equal to $2k$ bases in a fragment of length L and a given base composition (N_A adenines, N_C cytosines, N_G guanines and N_T thymines) can be calculated in the following way: The frequencies of the bases are: $p_A = N_A/L$, $p_C = N_C/L$, $p_G = N_G/L$, $p_T = N_T/L$. Let α be the frequency of a fragment RY (R = purine, Y = Pyrimidine). If $p(R) = p_A + p_G$ and $p(Y) = p_T + p_C$, we have $\alpha = p_{RY} = p_{YR} = p(R) \cdot p(Y)$. Then the probability of an AF not shorter than $2k$ bases is equal to $P = [2\alpha + p^3(Y) + p^3(R)] \cdot \alpha^k / (1 - \alpha)$ and the expected number A_{exp} of such fragments equals $L \cdot P$.

The expected number of uAFs of a given kind is calculated in a similar way: Let $\beta(AT) = p_A \cdot p_T$, $\beta(AC) = p_A \cdot p_C$, $\beta(GC) = p_G \cdot p_C$ and $\beta(GT) = p_G \cdot p_T$. Then the probabilities of uAFs are:

$$\begin{aligned} P(AT) &= [2\alpha + p^2(R) \cdot p_A + p^2(Y) \cdot p_T] \cdot \beta^k(AT) \cdot [1 - \beta(AT)]^{-1} \\ P(AC) &= [2\alpha + p^2(R) \cdot p_A + p^2(Y) \cdot p_C] \cdot \beta^k(AC) \cdot [1 - \beta(AC)]^{-1} \\ P(GC) &= [2\alpha + p^2(R) \cdot p_G + p^2(Y) \cdot p_C] \cdot \beta^k(GC) \cdot [1 - \beta(GC)]^{-1} \\ P(GT) &= [2\alpha + p^2(R) \cdot p_G + p^2(Y) \cdot p_T] \cdot \beta^k(GT) \cdot [1 - \beta(GT)]^{-1} \end{aligned} \quad (2)$$

The expected numbers of uAFs are then equal to: $A_{exp}(AT) = L \cdot P(AT)$, $A_{exp}(AC) = L \cdot P(AC)$, etc.

Occurrence of given fragments. Comparison between natural and random sequences. Let the number of fragments of a given kind (for example AFs, qAFs or potential Z-DNA sequences) found in a natural sequence be equal to A . The quantity $F = (A - A_{exp}) \cdot (A_{exp})^{-1/2}$ (analogous to a coefficient of variation) measures the degree to which the frequency of a fragment in a natural sequence differs from that calculated for the corresponding random sequence. If $F < 0$, we state that the fragment in the natural sequence is

F-fold under-represented. If $F > 0$ we state that a fragment is F-fold over-represented.

RESULTS

Occurrence of AFs. Table 1 shows the frequencies of occurrence of AFs longer than 5 bases. In each case we also show the corresponding values in the random sequences (second row of every case listed). It is evident that AFs

TABLE 1. Occurrence of uAFs and mAFs longer than 5 bases.^a

GENOME	LENGTH	GC+CG	uAFs			mAFs	All number	AFs F ^b
			AC+CA	GT+TG	AT+TA			
SV40	5243	0	3	2	0	42	47	-3.9
		0.2	0.8	0.8	2.5	77.5	81.8	
BKV	5153	0	2	0	8	40	51	-3.3
		0.2	0.7	0.7	2.8	76.0	80.5	
Polyoma	5292	0	1	2	2	54	59	-2.6
		0.6	0.9	0.8	1.2	79.1	82.7	
Human mito.	16569	0	15	1	12	143	171	-4.9
		0.7	10.8	0.3	5.0	231.2	248.1	
Murine mito.	16295	0	22	0	29	111	162	-5.7
		0.3	6.8	0.5	11.6	232.5	251.6	
Bovine mito.	16338	0	14	0	21	131	166	-5.5
		0.4	7.5	0.5	8.7	235.1	253.	
pBR322	4363	1	0	1	2	46	50	-2.2
		1.1	0.7	0.7	0.4	64.9	67.9	
Adeno-2· l	11600	36	1	9	3	113	162	-1.4
		4.4	0.9	3.2	0.7	171.2	180.3	
Adeno-2· r	10305	3	10	3	3	93	112	-3.8
		1.4	3.4	0.7	1.8	153.2	160.6	
Adeno-7· l	6707	12	1	3	3	33	42	-6.1
		1.0	0.6	1.9	1.1	99.2	103.8	
M13	6407	0	0	0	5	55	60	-3.8
		0.3	0.5	1.6	2.8	92.0	97.2	
φX174	5386	1	0	2	0	48	51	-3.6
		0.4	.5	1.4	1.6	79.5	83.	
T7	39936	3	16	13	3	459	494	-7.9
		6.4	6.3	8.2	8.1	673.6	702.7	
λ	48502	7	20	16	15	579	637	-4.3
		7.7	6.8	9.1	8.0	723.1	754.	

^a The second row of every case listed shows the values expected for random sequences. The first rows show the numbers found in natural sequences.

^b The last column lists the F-values defined in the text.

TABLE 2. Tendency of alternating fragments to cluster within DNA genomes^a

GENOME	$\langle L_{AF} \rangle$	$\langle L_{qAF} \rangle$	$\langle L_{qAF} \rangle - \langle L_{AF} \rangle$
SV40	2.85	9.48	6.61
BKV	2.90	8.70	5.80
Polyoma-A2	2.91	9.69	6.78
Human Mito.	2.87	10.97	8.10
Murine Mito.	2.85	11.48	8.63
Bovine Mito.	2.88	11.60	8.72
pBR322	3.00	13.60	10.60
Adeno-2·l	3.01	11.70	8.69
Adeno-2·r	2.99	11.55	8.56
Adeno-7·l	2.95	10.72	7.37
M13	2.82	11.60	8.78
φX174	2.90	12.79	9.89
T7	2.93	12.38	9.45
λ	3.01	14.02	11.01
RANDOM 8	2.97	12.96	9.99
RANDOM 4	3.02	13.69	10.67

^a $\langle L_{AF} \rangle$ and $\langle L_{qAF} \rangle$ are the average lengths of AFs and qAFs, respectively.

are under-represented compared to the expectation for random sequence (see the last column of Table 1 where all F values are negative). The same observation appears for mAFs which constitute about 90% of all AFs.

The above conclusion does not hold in the comparison of the occurrences of uAFs in the genomes studied compared to the corresponding random sequences. Thus, $[d(G-C)]_n$ and $[d(C-G)]_n$ fragments are strongly over-represented in adenoviruses ($F = 15$ for Adeno-2·l and $F = 11$ for Adeno-7·l), whereas in all circular DNAs they are generally absent (which is in agreement with the expectation for random sequences: $|F|$ close to 0). However, in phages T7 and λ these uAFs are under-represented ($F = -1.3$ and -2.5 , respectively). This suggests that in circular DNAs, long $[d(G-C)]_n$ and $[d(C-G)]_n$ tracts are avoided and that this circumstance arises at least in part from the base composition (correlation coefficient between the fractions of G and C and the F value is equal to -0.73 in these genomes).

Very different patterns of occurrence are displayed by the $[d(A-C)]_n$ and $[d(C-A)]_n$ sequences (and the complementary $[d(G-T)]_n$ and $[d(T-G)]_n$). The occurrence of these fragments seems to vary from one genome to another. Even

TABLE 3. Occurrence of qAFs and potential Z-DNA sites.^a

GENOME	qAFs		Potential Z-DNA			
	Number	% bases	Definition 1		Definition 2	
			Number	% bases	Number	% bases
SV40	13	2	10	1	8	1
	34	5	19	2	18	2
BKV	17	3	7	1	6	1
	32	6	19	3	19	3
Polyoma-A2	22	5	12	2	14	3
	33	8	19	3	18	4
Human Mito.	44	3	28	2	25	2
	102	7	61	4	61	5
Murine Mito.	51	4	21	1	23	2
	100	8	58	3	56	2
Bovine Mito.	53	4	21	2	22	2
	101	8	58	6	58	5
pBR322	22	5	17	4	16	4
	26	6	16	4	16	4
Adeno-2· l	52	5	46	4	46	4
	72	7	43	4	42	4
Adeno-2· r	39	4	28	3	29	3
	64	7	38	4	38	4
Adeno-7· l	32	6	22	4	21	4
	41	8	25	5	24	5
M13	14	2	6	1	4	0.7
	39	6	23	4	23	4
φX174	15	3	12	2	11	2
	33	7	20	3	20	3
T7	159	5	107	4	110	4
	246	8	151	5	151	5
λ	215	5	163	4	149	3
	299	7	184	5	184	5

^aThe first row of every case corresponds to a natural sequence whereas the second row corresponds to random sequences.

in the papovavirus genomes (SV40, BKV and Polyoma) the F values vary between 0.1 and 2.5. The same observations apply to the sequences of $[d(A-T)]_n$ and $[d(T-A)]_n$.

We can also see from Table 1 that mAFs are about ten-fold more frequent

than uAFs in both natural and random sequences. This means that although AFs are under-represented in natural DNAs, the expected proportion of uAFs to mAFs (about 1:10) is conserved in natural DNA tracts.

Tendency of AFs to cluster. Table 2 shows values for the average lengths for AFs ($\langle L_{AF} \rangle$) and qAFs ($\langle L_{qAF} \rangle$). The values of $\langle L_{AF} \rangle$ are almost the same for all genomes studied (i.e. about 3 bases/fragment). In contrast, $\langle L_{qAF} \rangle$ varies from 8.7 in BKV to 14.2 in the case of phage λ . It appears from Table 2 that the $\langle L_{qAF} \rangle$ values are higher in procaryotic than in eucaryotic systems ($\langle L_{qAF} \rangle$ is highest for pBR322 and phage λ , and lowest for eucaryotic papovaviruses SV40, BKV and Polyoma). It is also evident that in the case of single-stranded circular genomes (M13 and ϕ X174) the tendency of AFs to cluster is stronger than in double-stranded circular eucaryotic viruses. The same observations appear from an analysis of the difference $\langle L_{qAF} \rangle - \langle L_{AF} \rangle$ (last column of Table 2). In random sequences, there is a very strong tendency of short AFs to cluster (last two rows of Table 2).

The conclusion which can be drawn in this section is that the tendency of AFs to cluster decreases in the order: random sequence > procaryotic DNAs > mitochondrial and linear eucaryotic DNA > circular DNAs of eucaryotic viruses.

Distribution of qAFs and potential Z-DNA within genomes. The frequencies of occurrence of qAFs and potential Z-DNA sites in the DNAs studied (including the corresponding random sequences) are shown in Table 3. The linear maps of qAFs and potential Z-DNA sites for the genomes studied and the two random sequences are shown in the Figures 2 and 3. The plots for random sequences illustrate that our criteria for potential Z-sites restrict the amount of A and T in qAFs.

The general conclusion which appears from Figures 2-4 and Table 3 is that qAFs are under-represented in every case studied. This under-representation is not so universal for potential Z-DNA tracts. Except for the cloning vector pBR322, the circular DNAs studied display strong under-representation of potential Z-DNA sites (particularly in the case of mitochondrial DNAs). The same observation holds for the linear DNAs of phages T7 and λ , whereas potential Z-DNA sites occur with almost random frequency in the linear genomes of adenoviruses (Adeno-2·l and Adeno-7·l). Complete listings of all qAFs longer than 7 bases found in the DNAs studied are presented in Figures 5 - 8.

It is interesting that potential Z-DNA occurs both within coding and non-coding regions of the genomes studied. This finding is in agreement with existing experimental data (see Discussion).

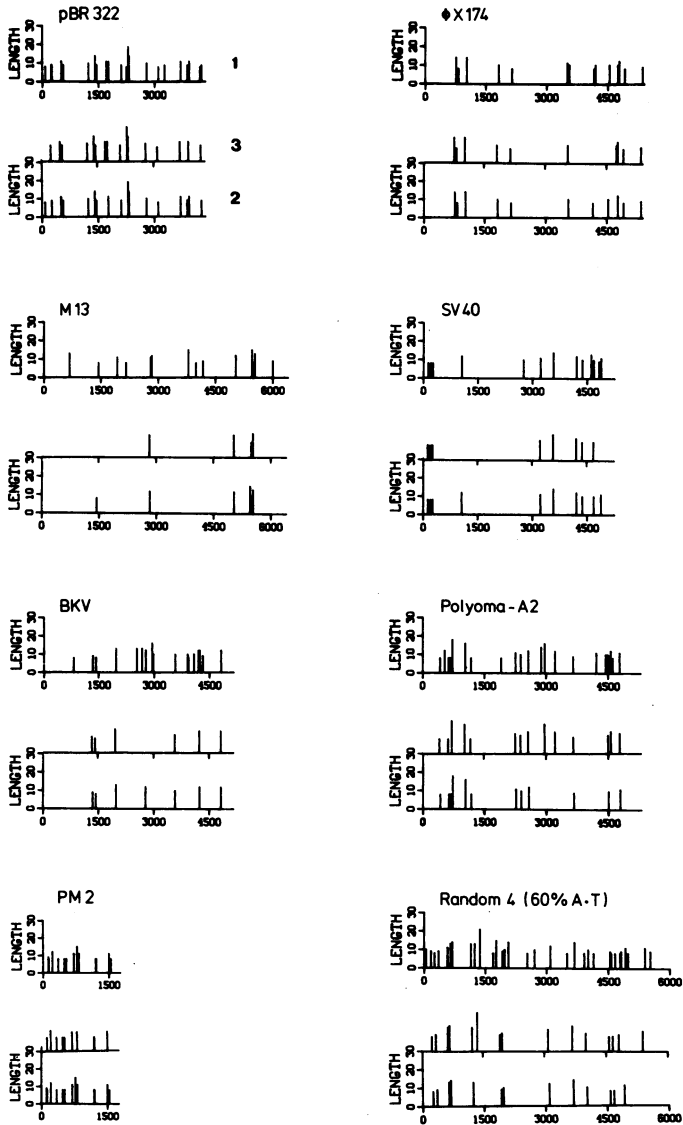


Fig. 2. Linear maps of qAFs and potential Z-DNA sequences in cloning vector pBR322, papovaviruses, dsc-phage PM2 (partial sequence), and ssc-DNA phages. The maps are rendered linearly starting at the origin defined in the Data Bank (EMBL Nucleotide Sequence Library) and extending to the right. The fragments are shown as vertical lines with lengths in bases given by the ordinate. The top horizontal line (e.g. number 1 in the case of pBR322) corresponds to qAFs, the bottom line (number 2) to potential Z-DNA tracts fulfilling definition 1 and the middle line (number 3) to potential Z-DNA fulfilling definition 2. Random 4 is a random sequence rich in A and T (30% each).

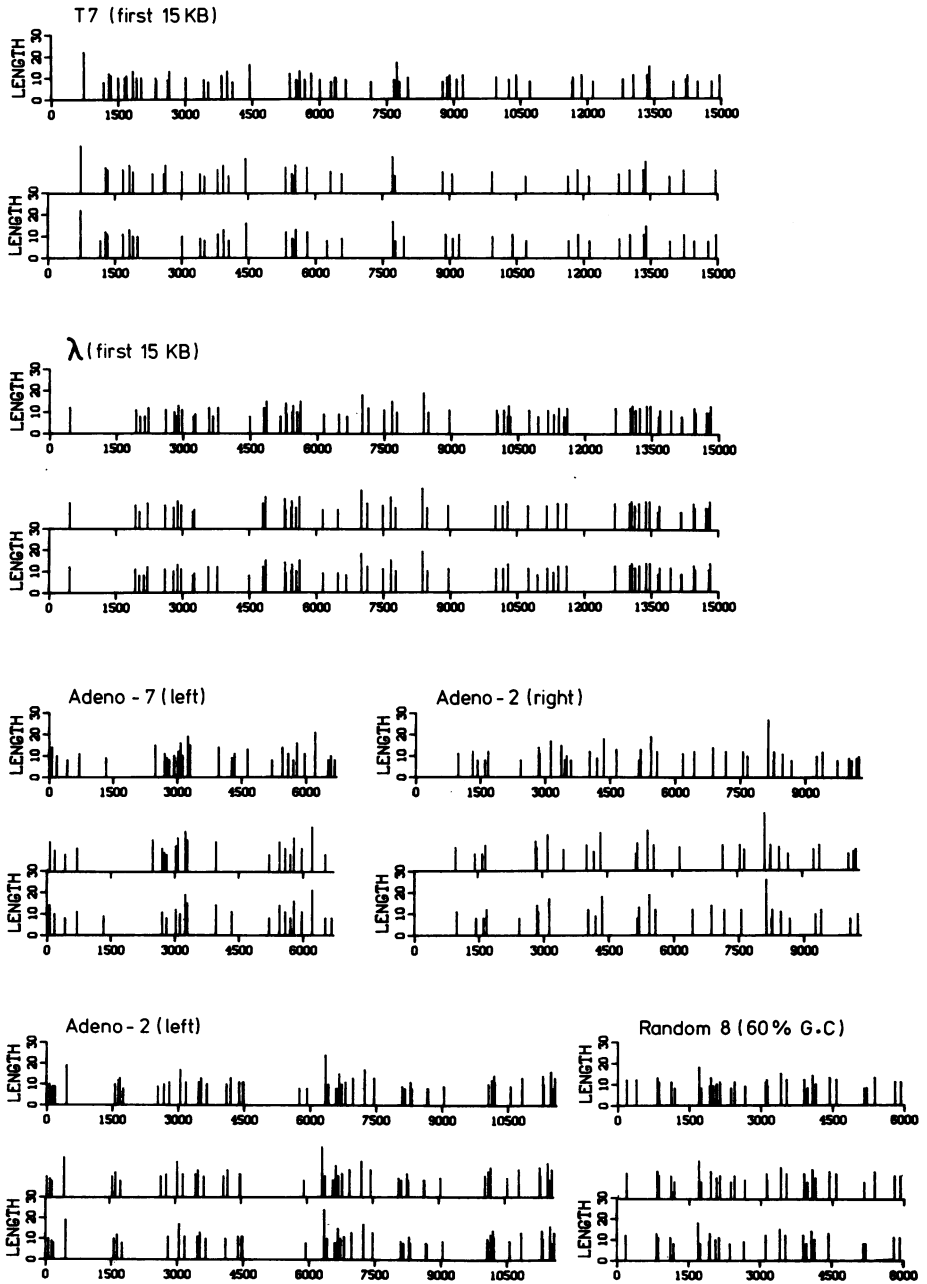


Fig. 3. Linear maps of qAFs and potential Z-DNA sites in dsL-DNA of phages T7 and λ and adenoviruses. Conventions as in Fig. 2. Random 8 is a random sequence rich in G and C (30% each).

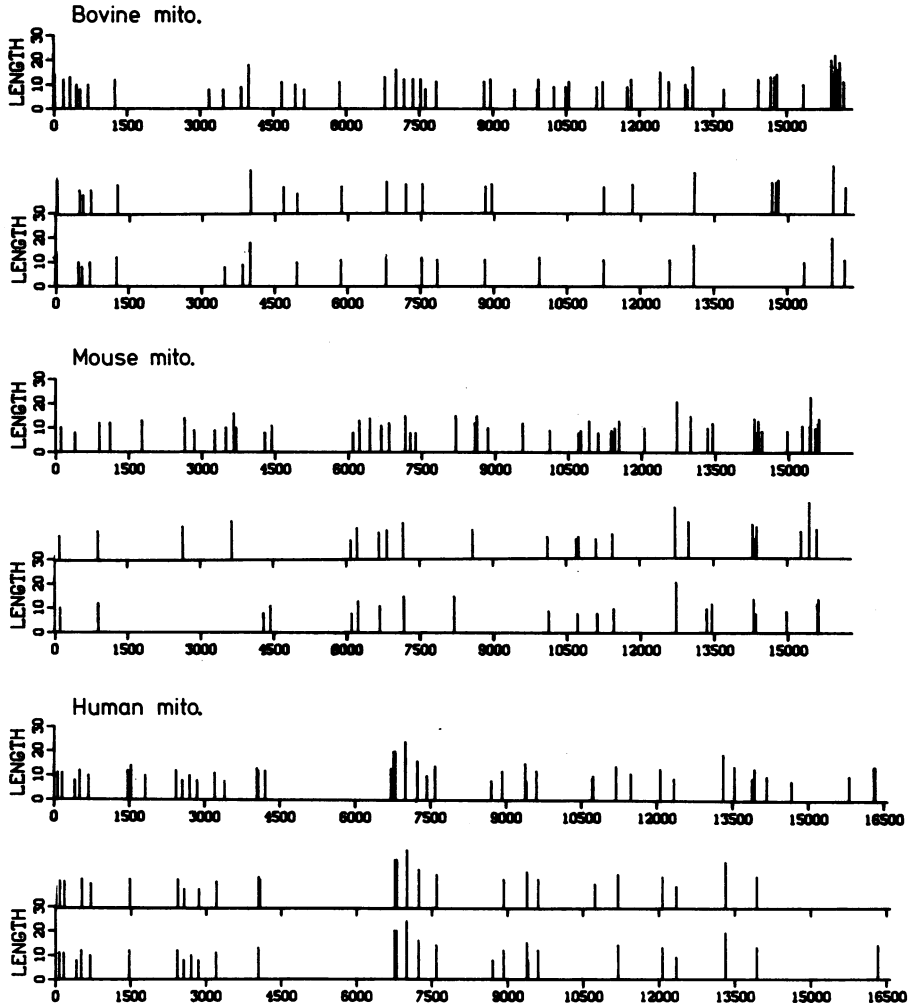


Fig. 4. Linear maps of qAFs and potential Z-DNA sites in mitochondrial genomes. Conventions as in Fig. 2.

DISCUSSION

We find that in the DNAs studied:

1) Alternating purine pyrimidine tracts are under-represented compared to the expectation for random sequences.

2) Uniform alternating fragments $[d(G-C)]_n$ and $[d(C-G)]_n$ are in general absent in circular DNAs and seem to be over-represented in DNAs of linear eucaryotic viruses.

Comparison of predicted with experimentally mapped Z-DNA sites. The studies of anti Z-DNA Ig binding to plasmid pBR322 DNA (5, 8, 19) provide evidence for 3 major and other minor immunoglobulin binding sites mapped by electron microscopy techniques (8) with a resolution of about 100 bases (major sites) and 300 bases (minor sites). There is good correspondence (with the exception of site B at position 960 ± 80 ; ref. 8) with the potential sites we have identified at positions (Figures 2 and 5): 237, 258, 1410, 1452, 2107, 2290, 2315, 2785 and 3099.

The anti Z-DNA Ig binding sites detected in the SV40 genome by filter-binding studies (18) and by immuno-electron microscopy (27) show the existence of 3 major antibody binding sites in the nucleotide sequences associated with the transcriptional enhancers within the nucleosome-free "gap" region of the papovaviral chromatin. These sites occur at positions 126, 198 and 258 and are predicted in this paper (Figures 2 and 5). Three other potential Z-DNA regions predicted by our algorithm (positions 1056, 3218 and 3575 of the SV40 genome) may correspond to minor antibody binding sites observed in the electron microscopy studies. Thus, among 10 predicted potential Z-DNA sites in SV40, 3 and possibly 6 have been experimentally detected, at least within the resolution currently available.

Mapping of anti Z-DNA Ig binding sites in ϕ X174 DNA provides further experimental evidence supporting the predictions of potential Z-DNA sites made in this paper. According to the listing (Fig. 5) and plot (Fig. 2) there are 13 potential Z-DNA sites in ϕ X174 DNA. Nine of them (positions 763, 811, 826, 1027, 2146, 3555, 4161, 4911 and 5345) correspond well with antibody binding sites identified by high resolution darkfield electron microscopy (23). Revet et al. (23) identify a site (no. 8) at position 3542 ± 62 and consider its possible relationship to the sequence starting at nucleotide 3504. By our criteria, this fragment is rejected due to its high alternating A-T content (Fig. 5). However, we note as a potential site the qAF at 3555 (which meets both definitions) and which is within the resolution limits of the site identified by e.m.

Studies of the PM2 bacteriophage genome also provides evidence for the correspondence of anti Z-DNA Ig binding sites to tracts of purine-pyrimidine repetitions (24, 25). The immunoelectron microscopy mapping of anti-Z-DNA Ig binding sites in the purine-pyrimidine rich region of this phage DNA (26) shows the existence of Z-DNA within a protein coding region. There are 13 potential Z-DNA sites predicted by our algorithm (Fig. 5). Ten of them (positions 129, 212, 345, 483, 528, 699, 812, 1194, 1205, 1494) correspond well with antibody binding sites identified by immunoelectron microscopy (26).

The correspondence between experimentally mapped Z-DNA sites in supercoiled circular DNAs and those predicted by the criteria we have defined is satisfactory but not perfect [some experimental positions we cannot account for and others we identify have not (yet) been observed]. As additional data emerge, the specific values of the empirical parameters (a, b, c in definitions 1, 2) will require adjustment. In any event, we expect that they will depend on superhelix density and, to a degree, each other. For example, alternating fragments exclusively composed of G and C are under-represented (Table 1, below) but where they do occur (23, 26) the Z conformation may be expressed even for lengths smaller than the value 8 used in this work. In addition, we do not address the means for defining a *hierarchy* in Z-forming potential, for which the experimental data provide some indications. It is obvious that the ultimate but as yet unattainable goal will be to replace the empirical criteria employed here with rigorous thermodynamically defined relationships.

The under-representation of potential Z-DNA It has already been suggested that Z-DNA could play a role in the control of transcription (22). In the circular DNA molecules, such processes would be coupled to changes in the free energy of supercoiling. (Since the B to Z transition lowers the negative superhelix density, one Z-forming tract may affect the potential of another; 23). Thus, it would seem reasonable that in such genomes the number of sites allowing a B to Z transition would be limited and highly regulated. Furthermore, the genomes examined in this work are almost fully transcribed. For these various reasons, the observed under-representation of potential Z-DNA forming sequences is not unexpected. In this connection, it is noteworthy that $[d(G-C)]_n$ and $[d(C-G)]_n$ tracts are avoided in circular DNA genomes, whereas these are the sequences which undergo the B-Z transition most readily *in vitro*. One can envisage positive as well as negative selection processes accounting for this phenomenon. Clearly, the intervention of proteins with specificity for different helical conformations as well as other factors determining higher order structure of DNA *in vivo* will determine which of the sites we and others have identified actually undergo the B \rightarrow Z transition and if so, whether functional roles are involved.

ACKNOWLEDGEMENTS

We thank Dr. E. Trifonov for discussions concerning the definition of potential Z-DNA sequences, and Drs. G. Hamm and K. Stüber for discussion and help in exploitation of EMBL Nucleotide Sequence Library. We are indebted to Dr. J.H. van de Sande for providing manuscripts prior to publication. Ms. Melanie Harvey is acknowledged for typing the manuscript.

*Present address: Stanford University, Department of Chemistry, Stanford, CA 94305, USA

+Present address: University of California-Berkeley, Naval Biosciences Laboratory, Oakland, CA 94625, USA

REFERENCES

1. Pohl, F.M. and Jovin, T.M., (1972) *J. Mol. Biol.* 67, 375-379.
2. Wang, A.H-J., Quigley, G.J., Kolpak, F.J., van der Marel, G., van Boom, J.H., and Rich, A., (1981) *Science* 211, 171-176.
3. Drew, H.R. and Dickerson, R.E., (1981) *J. Mol. Biol.* 151, 535.
4. Jovin, T.M., McIntosh, L.P., Arndt-Jovin, D.J., Zarling, D.A., Robert-Nicoud, M., van de Sande, J.H., Jorgensen, K.F. and Eckstein, F., (1983) *J. Biomol. Struct. Dynam.* 1, 21-57.
5. Nordheim, A., Lafer, E.M., Peck, L.J., Wang, J.C., Stollar, B.D., and Rich, A., (1982) *Cell* 31, 309-318.
6. Singleton, C.K., Klysik, J., Stirdivant, S.M., and Wells, R.D., (1982) *Nature* 299, 312-316.
7. Pohl, F.M., Thomae, R., and DiCapua, E., (1982) *Nature*, 300, 545-546.
8. DiCapua, E., Stasiak, A., Koller, T., Brahm, S., Thomae, R., and Pohl, F.M., (1983) *EMBO J.* 2, 1531-1535.
9. Wang, J.C., Peck, L.J., and Becherer, K., (1983) *Cold Spring Harbor Symp. Quant. Biol.* 47, 85-92.
10. Peck, L.J. and Wang, J.C., *Proc. Natl. Acad. Sci. USA* (1983), 80, 6206-6210.
11. Zarling, D.A., Arndt-Jovin, D.J., McIntosh, L.P., Robert-Nicoud, M., and Jovin, T.M., (1984a) *J. Biomol. Struct. Dynam.* 1, 1081-1107.
12. Zarling, D.A., Arndt-Jovin, D.J., Robert-Nicoud, M., McIntosh, L.P., Thomae, R., and Jovin, T.M., (1984b) *J. Mol. Biol.*, 176, 369-415.
13. Nordheim, A., Pardue, M.L., Lafer, E.M., Moller, A., Stollar, B.D., and Rich, A., (1981) *Nature*, 294, 417-422.
14. Lemeunier, F., Derbin, C., Malfroy, B., Leng, M., and Taillandier, E. (1982) *Exp. Cell Res.* 141, 508-513.
15. Arndt-Jovin, D.J., Robert-Nicoud, M., Zarling, D.A., Greider, C., Weimer, E., and Jovin, T.M., (1983) *Proc. Natl. Acad. Sci. USA* 80, 4344-4348.
16. Robert-Nicoud, M., Arndt-Jovin, D.J., Zarling, D.A., Jovin, T.M., (1984) *EMBO J.* 3, 721-731.
17. Russel, W.C., Precious, B., Martin, S.R. and Bayley, P.M., (1983) *EMBO J.* 2, 1647-1653.
18. Nordheim, A. and Rich, A., (1983) *Nature* 303, 674-679.
19. Azorin, F., Nordheim, A., Rich, A., (1983) *EMBO J.* 2, 649-655.
20. Quadrifoglio, F., Mannzini, G., Yathindra, N., and Crea, A., (1983) *Nucleic Acids: The Vectors of Life.* Pullman, B. and Jortner, J. (eds.) pp. 61-74. D.Reidel, Dordrecht, Holland.
21. Rich, A., Nordheim, A., and Azorin, F., (1983) *J. Biomol. Struct. Dynam.* 1, 1-19.
22. Revet, B., Zarling, D.A., Jovin, T.M. and Delain, E., (1984) *The EMBO J.* in press.
23. Rich, A., (1983) *Cold Spring Harbor Symp. Quant. Biol.* 47, 1-13.
24. Stockton, J.F., Miller, F.D., Jorgenson, K.F., Zarling, D.A., Morgan, A.R., Rattner, J.B. and van de Sande, J.H., (1983) *EMBO J.* 2, 2123-2128.
25. Miller, F.D., Jorgenson, K.F., Winkfein, R.J., van de Sande, J.H., Zarling, D.A., Stockton, J. and Rattner, J.B., (1983) *J. Biomol. Struct. Dynam.* 1; 611-620.
26. Miller, F.D., Winkfein, R.J., Rattner, J.B., and van de Sande, J.H., (1984) *Bioscience Reports*, in press.
27. Hagen, F.K., Zarling, D.A., and Jovin, T.M., (1985) *EMBO J.*, in press.