
Organisation of the entire rabbit progesterone receptor mRNA and of the promoter and 5' flanking region of the gene

Micheline Misrahi, Hugues Loosfelt, Michel Atger, Cécile Mériel, Véronique Zerah, Philippe Dessen¹ and Edwin Milgrom*

INSERM U.135, Faculté de Médecine Paris-Sud, 94275 Le Kremlin-Bicêtre Cedex and ¹Laboratoire de Biochimie (UA 240 CNRS), Ecole Polytechnique, 91128 Palaiseau, France

Received February 19, 1988; Revised and Accepted May 18, 1988

Accession nos X06623, X06624

ABSTRACT

cDNA clones corresponding to the 3' and 5' non coding regions of the rabbit progesterone receptor (rPR) mRNA and genomic clones corresponding to the promoter and 5' flanking region of this gene were isolated and sequenced up to nucleotide -2761. The 3' non coding region is very long (3058-3553 nucleotides) and contains three different polyadenylation sites. Primer extension experiments and S1 mapping showed the existence of 2 transcription initiation sites 699 and 712 bp upstream from the initiator ATG. The promoter region contains two modified TATA boxes: TAGAAA at -17 and TAGA at -37bp. A CAACT sequence is present at position -100 and one consensus binding site for the transcription factor Sp1 is found at position -51. A 317 bp sequence was observed (positions -2590 to -2273) which belongs to the C family of the short interspersed repeats of the rabbit. Sequences resembling the consensus for estrogen and progesterone responsive elements are observed at several locations in the 5' flanking region. The progesterone receptor is present in tissue extracts mainly as a mixture of two molecular species (110 and 79 kDa) whose origin remains currently debated. By Northern blot analysis we have shown, using rabbit and human mRNAs, that these receptor species are not derived from separate mRNAs. Transcription-translation experiments also showed that, at least in vitro, they are not derived by use of different translation initiation sites on the same messenger RNA.

INTRODUCTION

Steroid hormone receptors are presently the subject of great interest not only because they are the intracellular mediators of a group of important regulatory molecules but also because they are among the few well characterized proteins which interact with enhancer-like sequences. We have recently cloned the cDNAs encoding the rabbit (rPR) (1) and human (hPR) (2) progesterone receptor (PR). The coding regions have been sequenced allowing the deduction of the primary structures of the proteins. It has now become possible to study the molecular mechanisms of the hormonal regulation of these proteins which are induced by estrogen and down-regulated by progesterone (3). Moreover the very well defined cell-specific expression of the progesterone receptor gene is also amenable to experimental analysis.

However, as a prerequisite to such studies it was necessary to isolate and define the entire transcription unit of the progesterone receptor and not just its coding part, and also to clone and sequence the corresponding genomic fragments especially the promoter and 5' flanking regions. Steroid hormone receptor messenger RNAs contain unusually long 3' non coding regions (1,2,4-8). These regions have been implicated in post transcriptional regulation of various receptors (9,10) growth factors and protooncogenes (11). Extensive analysis of the 5' flanking region of the gene is also necessary since hormone responsive elements have been found as far as ~ 2.7 kb upstream from the transcription start site of some steroid hormone regulated genes (12,13). Furthermore it was also necessary to examine the problem of PR heterogeneity at the messenger RNA and the protein levels since different regulation might have existed for different receptor species.

MATERIAL AND METHODS

Cloning of the cDNAs corresponding to the 3' non coding region of rPR messenger RNA: Rabbit uterine mRNAs enriched for rPR mRNA (1) were used. A cDNA library was prepared in λ gt10 according to the procedure previously described (1) with the exception that oligo dT (instead of random primers) was used for the cDNA synthesis. The library (1.2×10^6 clones) was screened with a 300 bp probe derived from a random-primed clone (λ rPR7) (1) initiating at an EcoRI site 2463 bp after the TGA stop codon.

Isolation of genomic clones: A rabbit genomic library in Charon 4A (14) was screened with a fragment of the λ rPR8 cDNA (1) clone: a 180 bp EcoRI-ScaI fragment located 323 bp upstream from the first ATG of the open reading frame for PR. Among 550 000 plaques 16 were found to be positive. A 3.2 Kb EcoRI-ScaI genomic fragment (λ rPRG1) was further characterized.

DNA sequencing was performed after subcloning into M13 vectors (15) with the use of universal or specific primers (20-25 mers synthesized using β -cyanoethyl-phosphoramidite derivatives).

Primer extension and S1 nuclease mapping: Two primers were used : a 25 mer corresponding to nucleotides 393-418 upstream from the initiator ATG (primer 1) and a 21 mer located 244 bp further upstream (primer 2). They were labeled with T4 polynucleotide kinase. Single stranded probes initiating at those primers and ending at the NsiI site 1041 nucleotides upstream from the initiator ATG were synthesized using a M13 clone derived from λ rPRG1, and purified. Hybridization of probes ($2-5 \cdot 10^5$ cpm) with 20 μ g of uterine mRNAs

enriched for receptor mRNA (1) was performed either for 5 hours at 55° (primer 1) and 60° (primer 2), in 10 μ l of 20 mM tris pH7 buffer containing 0.3 M KCl (primer extension) or for 5 h at 60° in 20 mM tris pH 7.5 buffer containing 0.3 M NaCl (S1 mapping). After extension with reverse transcriptase or digestion with 800 to 2000 u/ml of Nuclease S1 (14) the DNAs were analyzed on polyacrylamide urea gels (buffer gradient and 6 % polyacrylamide for primer 1 or 15% polyacrylamide for primer 2).

Northern blot analysis: Rabbit uterine poly(A)+ RNAs (10 μ g) enriched for receptor mRNA (1) were analyzed (14) with two nick translated probes corresponding to the 5' and 3' extremities of the coding sequence: a 345 bp Sall-Tth111I fragment starting 5 bp upstream from the initiator codon ATG and a 350 bp HindIII-EcoRI fragment ending 70 bp after the TGA stop codon. Human poly(A)+ RNAs were prepared from the T47D breast cancer cell line. The 5' probe consisted of two successive PpuMI-PpuMI fragments (180 and 170 bp), starting 20 bp downstream from the initiator codon ATG. The 3' probe was a 330 bp HindIII-BstXI fragment ending 40 bp before the TGA stop codon.

Transcription-translation experiments: PR cDNAs containing the entire coding sequence were constructed by using clones λ rPR5 and λ rPR6(1), and λ hPR1 and λ hPR4 (2), for rabbit and human receptors, respectively. In vitro transcription and translation were performed in PGEM4 vector as described by the manufacturer (Promega Biotec). Immunoprecipitation of the ³H leucine labeled translation mixtures was performed (16) with Let 126 (10 μ g) (17) antireceptor monoclonal antibody. Total translation mixtures or immunoprecipitates were analysed by SDS-polyacrylamide gel (9%) electrophoresis and fluorography.

RESULTS AND DISCUSSION

Cloning, sequencing and analysis of the entire rabbit progesterone receptor mRNA.

3'non coding region (Figure 1): In previous experiments where the coding region of PR cDNA was cloned, we had prepared a library containing random primed cDNAs (1). This was necessary due to the extremely long 3' non coding region of receptor mRNA but in turn resulted in clones lacking the 3' end of this messenger. Consequently an oligo dT primed cDNA library was prepared and screened yielding 45 positive clones. Analysis by Southern blotting (14) after EcoRI digestion of DNA showed that the clones fell into 3 categories according to the size of the 3' EcoRI-EcoRI fragment. Two phages from each group were isolated and their inserts were sequenced. The cDNAs were strictly identical and colinear. They corresponded to different sites of

```

+1      10      20      30      40      50      60      70      80      90      100     110
PCA ATGTCAAATTTATTTTCAAAGAAATTAAGTGTGGTGTATGGTCTTCGTTTGGTCAGGATATGACGCTCGAGCTCGAGTGTTTATAATATCTTCTGAAAGCCCTTACATTAATTAAC
130     140     150     160     170     180     190     200     210     220     230
ATATCATACCGGTGAATTTAGAGGAAGATTTGGAGACTCAATATTTCCGTTAAGGCAATTAAGTTTTAAAGTGTGTTCTGGCCCCCATATTTCTTTAAAGTTTAAA
250     260     270     280     290     300     310     320     330     340     350
GAGTTTAAAGTTGAAAAGTACTAAACGATGATTATGAAGTAAGCTATGGCTTACCATACTATTCATACGATTTAGTGGAGATTTTAACTTTTACATATAACAATCTCTACT
370     380     390     400     410     420     430     440     450     460     470
TTAGAGAAAGAAATCTCACATGTAATAAATAAAGCTACTATTATGTTACTCTAGATAGCTCCCTTTTCTCCTGACTGTACTTCAAAGTGAACCTTTAAATGGTATGCAAA
490     500     510     520     530     540     550     560     570     580     590
AACTTTGTCTACTAGGTGTGGTGTGTATATATATATACACACACAGATATATATGTTTCTGAAAGGAAATTAACAACGATTTCTAAGAGTTTTTAATGACAAA
610     620     630     640     650     660     670     680     690     700     710
AATATAGACAGAAAGAGTAAAAAACAACCTAAACAGATTACCATTTTTCAGACTAGACAAAACAGTCTATGTTGAAGGATTTCTTATATGGAAACCAATCTATAAGGAAAT
730     740     750     760     770     780     790     800     810     820     830
TAGTAATGAGTAGAGTTTACTGTTATGCCAACAGTATGCGAGCTTTTGTAGGGCACTTTTGGTAACCTTTATAGGAGGGAGTTGCTCTTAACAAGATATTGACTGAAACATATC
850     860     870     880     890     900     910     920     930     940     950
TGTGATCTGACTTTTACCATCTGGCATGGGAAGTTTCAATTTTCTCACTTTATGTTGGTACAATGATATCTTCTTCCCAAAATCTCTTGGTACTGAGGCTCTTTTAACT
970     980     990     1000    1010    1020    1030    1040    1050    1060    1070
CTTCTCAAAATCAAGAAGGAGGGATGGAAAGGGAGATAGGGAGGGAAGAGAGGGATATAAATCCCTTTTCCACTCCAAAGACTTAAGAGTGGTGTCCCTTCTGCTGAGTGT
1090    1100    1110    1120    1130    1140    1150    1160    1170    1180    1190
AATAGATGAATCTCAGGTTGAAAATTTTGGATGCCCTTGTATCCAGTTTTATTAGAGATGCTCAGGAAATAGTAGTGGTCTTAAAATTAACAACCAATTAACAATTAACCAAGAA
1210    1220    1230    1240    1250    1260    1270    1280    1290    1300    1310
CACAAAATCACATCTGTTGAACTGAAATATTACTTTAGAAAGCAGATAGTCTTTTCTGAGGTAATCACTGGAGATGATTTCTGAAAAGAAAGGAAAGGATATATTTATATA
1330    1340    1350    1360    1370    1380    1390    1400    1410    1420    1430
TACTATCCATACATAATAACATGGGAAGACCAAGTTCATACCATTGGTGAATTCAGTGAACCCGGTGTTCCTTCTAGATAAAGTTAGTTATACCAGTTTTAGGACGTATGAG
1450    1460    1470    1480    1490    1500    1510    1520    1530    1540    1550
TATACATACTTTTTTAATTAACCTCAACAAAGCTAATAGGATAAATTTTTTCTAAATTCGCCAAAATGGCAATATCATTTCTCAAAATAAAATACATCTCCCTCAGGACGATTA
1570    1580    1590    1600    1610    1620    1630    1640    1650    1660    1670
CTAGTTGTAATTAGAAAATAAAGTAACCTTGAATAATTTTACATTTGAAAACAATATAGAAAAAATTTGCTGAAATTTCTGAGTCTTAATGTGTTTCAATCTCCATCTCTCTT
1690    1700    1710    1720    1730    1740    1750    1760    1770    1780    1790
TGGTAATCTTCCAGTTTACCAGGGATAATAACATAAGAAATTTCTAGCAAAATCTGCTGACCCCTAAAGCACAAGTAAGTGGTTTCTAGTGATAAATATGGGCAATAAT
1810    1820    1830    1840    1850    1860    1870    1880    1890    1900    1910
ATTACATAGCTGGTTTTTAAACCAATAATCTCAACTTTCAAAAATATAGCTATGAGGTTCTAAATAGGGTGTCTATAGCCATTAAATAAGATTTGAAACATGCTCCAGA
1930    1940    1950    1960    1970    1980    1990    2000    2010    2020    2030
TAGAGCCCTTGGCTGGATTCATTAACCTATGGCCAAATAAATAGGTTATATACATTTTAAATTTTAGCTATAACAGCTCAACAAAATTTAAGTAACTTCAAGTAAAGCCCTC
2050    2060    2070    2080    2090    2100    2110    2120    2130    2140    2150
AAAGATTTAAATTAGCCGCTTTTGTAGTATTTCTGTGATCATTTTTCTGTGGATGATTTGTTGATAGTGTATATAGACTGTGACAGGCATACCTCCACCCTAGCAGCAGTT
2170    2180    2190    2200    2210    2220    2230    2240    2250    2260    2270
ATCATCTAAAAGGGATCTACTTTGCCATTTCTCTCACTGTGAACCTTTCAGTATTTCCAGTCACTGTCAGCTCCAGCTCAGCACTCAAGATTTCCAAATGAACCTTAAAGAATGC
2290    2300    2310    2320    2330    2340    2350    2360    2370    2380    2390
TATTACTGAGCTCTCTTCTCAGTGTATTAATGGGCTACTTGGCATTAACTGCAGGCACAGAAATAGTATTTTGAATATGTTTCAACTGTCAATCTTTTACTTTTGTGCTTTTACCC
2410    2420    2430    2440    2450    2460    2470    2480    2490    2500    2510
TCCACTGTAACTGTATGTCGCCAGCAGAAATCTTTACCCACCATTTTTTACCTGCGCGAAATTCGCCACCTTTTGACAAAGACAGCAAGTATCTTTCTGTAAAGTCTTCCCT
2530    2540    2550    2560    2570    2580    2590    2600    2610    2620    2630
TAGTCTTCAGACAGTGGTTGCTTCTCTCTCAGTTCACCTACCGTTTTTGGCAGATCTGAATGACTGTAGCACAACCTGAAATCTCTCAAAATCTGAGGTTGAGGTTGGAAC
2650    2660    2670    2680    2690    2700    2710    2720    2730    2740    2750
ATTGTGCTGTGGACTGAAAGGTTCAATGACACATCTCTCTGTAACCTGGTATCACCATGAAGTCTTCACTCAGCTTTGATTAGCATCATGCTATGCTGTGCTGACTTATGTTTACCACCT
2770    2780    2790    2800    2810    2820    2830    2840    2850    2860    2870
TCTCTTTTCAATCTTTTATCTTCTTTTGGTGTGAAATTTTCCAACCTTTTTTCCGGCTTTGAGAGATTTCTCCCTTTACCAACTATAAATAGGAAATGAAAGTACATCTTCTG
2890    2900    2910    2920    2930    2940    2950    2960    2970    2980    2990
TTTATTAATTAAGATACTTTATCTGTCTTAAATTTAACAATGAATGAATGGAATACTTTTTTCTCCCTAGAAAATTTCCAGACTGTGTTCTTATGTTTAAGTAGCCGAAAATGTT
3010    3020    3030    3040    3050    3060    3070    3080    3090    3100    3110
TTTGTATTGTAACCTCAGTATGAAATTTGTTATTTGATTTCAATTAACAAAATTTTCTGTTATTGGCATGAGAAAATGTCATACACTTTGTGAGCAATTTAGCTCAAGA
3130    3140    3150    3160    3170    3180    3190    3200    3210    3220    3230
AAGACTATAATTAAGACAGCTGGTAGCTCAGTAGTAAATATAAATGACAGGGTGATAAATTTGAGGCACTCAAAATATGTTGACATGACATGTTGATGGTGCAAAAGCATTATGTACATA
3250    3260    3270    3280    3290    3300    3310    3320    3330    3340    3350
TAATTTGACTAGACAAATGAGCAGCAAAATATACTAACTTTGAATAAATTAATTAACAACTTTTGTCAAGTGCCTAAATTTCTTTAATTAATCTATGTGACTGTTTGTTTACT
3370    3380    3390    3400    3410    3420    3430    3440    3450    3460    3470
AAAAGCACTTTCCCATATCTTTCCCTACCGTGTATCCAGTAAATAAGATTTATAGTCCCATATGCAATTAATTAAGGAAATATCCCTCTGGGAGTGCCACTTCCAGAAAT
3490    3500    3510    3520    3530    3540    3550    3560    3570    3580    3590
GGATTATGACATAAATGAATGATGTTTTCGAAAAATAAATAAATGAGAACAATGTTGATGCCCTTTAAAAAATAA

```

Figure 1: Nucleotide sequence of the 3' non coding region of rPR mRNA:
 The sequence starts with the TGA stop codon (which is boxed) ending the open reading frame encoding rPR. Position +1 is assigned to the next nucleotide. The 3 polyadenylation signals are underlined. The terminal poly A stretch and the two other poly A tails are indicated at positions 3058 and 3300. The 52 bp alternating purine pyrimidine sequence beginning at position 494 and the poly A rich region which follows it at nucleotide + 617 are underlined. The ATTTA sequences are squared. Three bases have been found to differ in different clones, they are noted by an asterisk.

polyadenylation localized 3058, 3300 and 3553 bp (only 2 clones from 45) after the TGA stop signal. The corresponding polyadenylation signals are a modified AACAAA sequence, then two canonical AATAAA signals, respectively. The sequence of the 3' non-coding region of the rPR mRNA showed no homology to other receptor mRNAs of the steroid-thyroid hormone family whose structures in this region have been determined (4,7,18) nor to chick progesterone receptor (8) except for 140 bp (at positions 7-146 and 69-208 after the TGA stop codon for the chicken and rabbit PR, respectively) which exhibits 62.9% homology but with many insertions. The use of different polyadenylation signals has been described for human and rat glucocorticoid receptors (GR) (4,5). Some striking features can be observed in the 3' non coding part of the rPR mRNA: A stretch of 25 nucleotides (494 bp after the TGA stop codon) showed alternating purines and pyrimidines as has been observed in the Z conformation of DNA. Such sequences are considered as possible sites of regulation of transcription (19), they may belong to repetitive elements and are often associated with a 3' polyA tail (20). Indeed a (A)₁₀ tract is present 63 bp downstream from the alternating purine pyrimidine region. The role of ATTA sequences in controlling messenger half life has been discussed (11). Several such sequences are found in rPR mRNA, especially at the beginning of the 3' non coding region and next to the polyadenylation signals.

5' non coding region (Figure 2): The 5' non coding region of rPR mRNA extends 712 bp (see below) and is thus markedly longer than that of most eukaryotic mRNAs (21). It shows no homology to the shorter (366 bp) equivalent region of the chicken PR mRNA (8). The exact length of the 5' non coding region has not been reported for other receptors (4,5,7,18, 22-24) except for the human estradiol receptor (hER) (6). The 5' non coding region of the rabbit PR is very G+C rich (63.8%). Several regions predicted to have very stable secondary structures are observed as has been described in many eukaryotic mRNAs (21). The 5' non coding region contains three short open reading frames initiating with ATG codons. Upstream ATG codons are rare in eucaryotes except in protooncogenes (21) but have been described in ER (6,25) and hap (26) mRNAs.

Cloning, sequencing and analysis of the promoter region and of the 5' flanking region of rPR gene:

A fragment of the rPR8 cDNA clone (1) (323 to 503 nucleotides upstream from the first ATG of the open reading frame) was used to screen 550 000 plaques of a rabbit genomic library. Sixteen positive clones were

detected. Preliminary analysis by Southern blotting led to the selection and sequencing of a clone containing an EcoRI-ScaI fragment extending ~ 3.5 Kbp upstream from the first ATG of the open reading frame.

```

-2750   -2740   -2730   -2720   -2710   -2700   -2690   -2680   -2670   -2660   -2650
TCITTCATCAAAGATTCAATATTCAGGAAGTTCTCTGGCCAACCATATCATCACTAACGATCAAATCTTGTGGATAAATGGTGAATGCGATATGCTGGGCACCTTATTATTATTTT
*****

-2630   -2620   -2610   -2600   -2590   -2580   -2570   -2560   -2550   -2540   -2530
CCCAGTCAAGATTTTCATCTTGGTGGAAAATAACTAGCAGGGGCCAGCGCTGTTGGGTAGCAGGTTAAGCCCGCCCTGTAGCTGGCATTCCCATTTGGCACCGGTTTCGAGCCCTGGC
*****

-2510   -2500   -2490   -2480   -2470   -2460   -2450   -2440   -2430   -2420   -2410
TGCTCCATTTTGTATAGCTCTCTGCTATGGCTGGGAAAGCCTACAAAGATGACCCTAGTCATTGGGCCCTGCAACCCATGTGGGAAGACCCGGAGGAGGCCTCTGGCTTGGATCAG
*****

-2390   -2380   -2370   -2360   -2350   -2340   -2330   -2320   -2310   -2300   -2290
TGCAGCTCCTGGCACTCGGCCAGTTGGGGAGTGAACCAGCAGATAGAGAAGCTTTCTCTCTCTCGCTCTGCTTCTCTCTCTGCAATACTCTGACTTTCAAATAAAATAAATAATT
*****

-2270   -2260   -2250   -2240   -2230   -2220   -2210   -2200   -2190   -2180   -2170
TAAAAAGAAAGAACTAGCCAGACTCTCCGCTTAGAGGGAAGATGTAAAATCAACCACAGTTAAAATAAATTGCAGATTTAGGATGAGGTCAGGTCGCCAAATTTACGGCAGCAAC
*****

-2150   -2140   -2130   -2120   -2110   -2100   -2090   -2080   -2070   -2060   -2050
CTAGAACAACAGCTGTAGTCATCTGATATTTAATAACCACTCTAATAAGACTATTTGGCCATTGGGCATAAATCTGCTCAATGCTACTTACTATATTTATAATAGTGGTAACTTTGT
*****

-2030   -2020   -2010   -2000   -1990   -1980   -1970   -1960   -1950   -1940   -1930
GTCTCACTTTTCCCAACAAGGAGAGATAAGGACTATTTTGGAGCTTTGAAGTATTAATGAGTTTTACACACATATAGGGCTTAGAAGATTACTGGCAATGCTGTCCAGATT
*****

-1910   -1900   -1890   -1880   -1870   -1860   -1850   -1840   -1830   -1820   -1810
TAATATTCTCAGACTTCTAGCCCGCCATCAAGGATACACAATCAAAAGCTTCTGTGATACCCCTGATACAGCAGAACCAAGAAATAAAAAAAGACTAGGTTCGTAAGCAT
*****

-1790   -1780   -1770   -1760   -1750   -1740   -1730   -1720   -1710   -1700   -1690
TGCTAGTGTGTAGAAAATGAACTCTGTAACCTGGTGGGAAACATAAAATGCACTCTTTGAAGAAAACATCCAGACTCTTGTACGCACACATGCAATATATCTTGGCAATA
*****

-1670   -1660   -1650   -1640   -1630   -1620   -1610   -1600   -1590   -1580   -1570
AAAATTAATGACGTTACTGATACATACATATGATAGTTGAACATAACTTCTGAAGTAAAGAAACTGTCAAAAAAAAAGACTCTATGGGATGGTAAATTTATATAAAATTTCCAAA
*****

-1550   -1540   -1530   -1520   -1510   -1500   -1490   -1480   -1470   -1460   -1450
ATATGAAATTCACAGACATATAATAGTGGTCTGTGGGAGTTGGGGGAAGTAGGGAATAACTACTAAAGGATAGGATTTCTTTTAGGATTTAATAAATTTTAAAGTGA
*****

-1430   -1420   -1410   -1400   -1390   -1380   -1370   -1360   -1350   -1340   -1330
GAATGATGATGTTTGCACATGAAACATGCTAAACCCATTGAATTTATGTTAAATGGGTAACCTGATGTACCTGAATATATCTCAATAAGACACTTGTGTTTAAAGTGTA
*****

-1310   -1300   -1290   -1280   -1270   -1260   -1250   -1240   -1230   -1220   -1210
CTGAAAAGAGGAAAAGCCGATTAACATTTAAATATGATAATAGTTTCTGTGGGTATGACTTCCATCCCACTACTTTCTGTACTCTTTTGTAAATGATCCCTCCACAGTTC
*****

-1190   -1180   -1170   -1160   -1150   -1140   -1130   -1120   -1110   -1100   -1090
CTACATTATGTTTCTTAAACCCCTGTGACCACAAATGGCTTCGCCAGATTTCTTTTCTCTCCTACTCCTGCTAACTCTTGTGACACATTTAAATAGCCAGTCAAGGAGCAT
*****

-1070   -1060   -1050   -1040   -1030   -1020   -1010   -1000   -990   -980   -970
TAATGCTGGAACAATCAGTCCCAACCCAGTCAACTAATCGTAAAGGTTTATTCTCATCCATAAGCCCCCTCCCAGGTGATCCAGGGTCTCCCTCCAGTGCTGTCTCCATC
*****

-950   -940   -930   -920   -910   -900   -890   -880   -870   -860   -850
CTGGAGCTCCTGGAGCCCTCCACTAGACCTCTATTTCTGGCTGAGAGAAATTTGGATAGCAGGGTGGTCTACAAGAGGTTTCCAAATGATCAACTCTATCTTGAAGATGATGCCA
*****

-830   -820   -810   -800   -790   -780   -770   -760   -750   -740   -730
GGACACCAATGATCTGGTGTGGCAGAAGCTGTACTCTGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT
*****

-710   -700   -690   -680   -670   -660   -650   -640   -630   -620   -610
TGTGTTCAAGAATTTAGAGGTGATCTGTTTAAATAGCTGAAAAGACAAATTTCTAAAGTAAGAATTTGCAGGAGAAAATTAATCAAGCAGTCTTTGGGGAATCTTATGATGT
*****

-590   -580   -570   -560   -550   -540   -530   -520   -510   -500   -490
GGCACAGAGTTCAAGAAATCTCTTAATATTTATAAAGAAATGGATTGTGTCTGTAAGACAGATCAAGACACAGCATAACAATTCAAGCATCTCGCATGGGTAGAAGGTTGTAGC
*****

-470   -460   -450   -440   -430   -420   -410   -400   -390   -380   -370
AACCATCTCCCATTCCAAATTTGGTTCACATCAAGGACTGTTAAAGTAAAGCACTGTCATGCGATTGGAGTGAAGAAAATTAAGAAGCACTCACTTTTCATGGACTCAAGGGAGAGA
*****

-350   -340   -330   -320   -310   -300   -290   -280   -270   -260   -250
CTCTCTTAAGATGATTAGGATCCTCTGATCGATGATTTCTAAAGACAATTTGGCCAGAGGTTCTCTCTGACAATCTCTATAAAACAGAAATCTGAAAACGGTGCAAAACCG
*****

-230   -220   -210   -200   -190   -180   -170   -160   -150   -140   -130
TGCTTTCTCACTCAACTATTTCCGATTTCCAAATTAATGGGATGCTTCTTAACTACTAATAACAAAACCTCATATGCTTACGGGGCTGCAATCCCAAGCACCTGCTTATGGGAGTTGATT
*****

-110   -100   -90   -80   -70   -60   -50   -40   -30   -20   -10
CAGACGATGGTGGGAAATCAACTTGTAGAGTTTCAGCTCTGCTGGTCCGATTTGGGATAGGGAGGGGCTTTGGGCGGGCCTCTCTAGAGGGGAGGAGCCGTTTGTAGAAAGCCATGGGGC
*****

+1       10       20       30       40       50       60       70       80       90       100       110       120
CAGTCCCGCTGTCCTACTCGGTTAAGCTTTTAGACTCTCGGTTGTGAGGATTCGTGGTGTGGGAGCCCTTCAATGAGAACAAGCTTCCACTTCCGATCGGACCTTGGCCACCA
*****

130      140      150      160      170      180      190      200      210      220      230      240
GAGGGCCCTCTTCTCCCGGAAACCGAGGAAACCGACGACGACGAGGAGCACAGTGTGGCCCGCTCTAGCGGGCGGGCAGAGGCCATCTCTGGGAAATGAGGCTGTACAGAGGGTCAGCT
*****

250      260      270      280      290      300      310      320      330      340      350      360
AGTCCGACGTGCTAGAGAGGGGCACGGGAACGCACTCAGGGGCTCGGAGCTTCCGACGGTGCACGGGTTTGAAGCTCTGGGAGAAATCGGGCGCTTAATGAATGCAAGCGCTGC
*****

370      380      390      400      410      420      430      440      450      460      470      480
CGCAGCCACTCCGGGGGACTTGTGAGTATCTCCGCTCGCCCGCTCGBGACAGAGCCAGAGGACTTCCAAATCGGAGACGAGCCCTGCAACTACTTCTCTGCTTCTCCCAT
*****

490      500      510      520      530      540      550      560      570      580      590      600
TGCCCGAGGACTGAGGACCGCAGCCCTGCCCTGACCAAGAGTGAAGGCGGCAAGGCGAGCGGACCGGACCGGCCCCCTCCGACCCAGGAGGTGGAGATCCCGGCGGTCCAGCCA
*****

610      620      630      640      650      660      670      680      690      700      710
AACCCCCACCCATTTTCTCCCTCTGCCCCTATATACCGGGACCCCTCCTCTCTCTTCCCTCTCTCCGAGACGGGGAGGAGAAAAGGGGAGTTCACGGTCCAG
*****

```

Determination of the transcription start site (Figure 3): Using probe 1 (see Material and Methods) S1 nuclease mapping revealed 2 protected bands. The same bands were observed in primer extension experiments. With the shorter probe 2 and a more proximal primer the two transcription start sites were precisely localized to the adenines present at position 699 and 712 upstream from the first ATG of the open reading frame. The most 5' transcription initiation site was defined as nucleotide +1.

Analysis of the promoter and of the 5' flanking region of the rPR gene (Figure 2): The sequence of the 2761 bp upstream from the first transcription initiation site was determined. The rPR gene has no typical TATA box but two motifs resembling that sequence are found: a TAGAAA motif is found at position -17 and another TAGA motif is present at position -37 bp. Such modified TATA boxes are often associated with multiple initiations (review in 27). A CAACT sequence is found at position -100 which may correspond to a CAAT box. Up to this point there is a relatively high G+C content in the promoter region (57%) due in particular to a high proportion of G, so that no stable secondary structure can be predicted which is in

Figure 2: Nucleotide sequence of 5' flanking, promoter and 5' non coding region of rPR gene: The 5' non coding sequence is indicated in italics. The sequence derived from cDNA clones was from +209 to +712, the sequence derived from genomic clones was from -2761 to +389. One divergence in the sequence between cDNA (on top of the line) and the gene is indicated by an asterisk. The putative initiator codon is boxed at the end of the sequence (the upstream in frame stop codon is also boxed). In the 5' non coding region other initiation and stop codons giving rise to open reading frames (indicated by dashed lines) are also squared. The two transcription initiation sites are indicated by closed triangles. The most 5' one is noted nucleotide +1. In the promoter region the two modified consensus for TATA boxes at positions -17 and -37 and the putative CAAT box at position -100 are squared. A consensus sequence for Sp1 binding site at position -51 is underlined. The direct repeats in the promoter region are shown. The modified ERE consensus sequences are squared at positions -427 and -243. In the 5' upstream region: the C Sines repeat is bracketed (positions -2590 to -2273) and the direct repeats in its flanks are indicated. ATG and TGA codons defining an open reading frame within that region are boxed. The sequence homologous to the bovine acetylcholine receptor alpha subunit mRNA is also bracketed (positions -1677 to -1344). The motifs resembling the consensus described for the palindromic ERE at positions -2477, -2462 and -2392 and at positions -1940 and -743 are boxed. The cluster of close PRE consensus sequences is indicated: each element is sublined by a thick arrow according to its direction. The other 3 tandemly associated sequences resembling the PRE consensus are indicated at positions -1998, -1972 ; -1774, -1753 ; -458 and -429. Within the cluster, two sequences matching silencer consensus element are underlined by dots at positions -1142 and -1216. Are also indicated: the sequence matching the heat shock factor binding site at position -2736 (which is underlined) and the sequences similar to the SV40 core enhancer (indicated by stared arrows at positions -2689, -2637, -2618).

contrast to the adjacent 5' non coding region. Further upstream the G+C content decreases strikingly (39% from -100 to -600). One putative binding site for the transcription factor Sp1 (28) is found at position -51. Between positions -269 to -226 a complex structure composed of two 10 bp direct repeats is observed, flanked by a pair of heptameric direct repeats. When found in promoter regions such structures have often been shown to be involved in regulatory mechanisms (29). A palindromic sequence very similar to the heat shock factor binding site consensus sequence (see 28) is present at position -2736 and a cluster of 3 sequences homologous to the SV₄₀ core enhancer (27) occurs in both orientations between positions -2689 and -2611. Sequences located at -2689 and -2618 are part of a direct 13 mer repeat. Close to this region a 317 bp long sequence (-2590 to -2273) was found belonging to the C family of rabbit short interspersed repeats (SINES) (30). This sequence is flanked by two direct repeats of 8 bp and ends with a poly A. It is intriguing that repetitive sequences of the same family are found in the 5' flanking region (in about the same position: -2.7 Kb to -3 Kb) and in the first intron of a rabbit progesterone regulated gene: the uteroglobin gene (31). Repetitive sequences are involved in the regulation of the expression of several genes, especially in the case of the silencer elements of the rat insulin and chick lysosyme genes (32).

Progesterone receptor is induced by estrogen and down regulated by progesterone (1,3,16). It was thus interesting to investigate the possible presence of sequences resembling the consensus sequence of estrogen responsive element (ERE) (33,34) and progesterone responsive element (PRE) (13). This required the sequencing of about 3000 bp of the 5' flanking region of the gene since hormone responsive elements have been described either close to the transcription initiation site or up to -2.7 kb (12, 13). Three modified consensus sequences for ERE GGTCANNNTGACC were effectively found. They are located within the C SINES repetitive sequence (nucleotides -2477 to -2380) and correspond to three closely related, imperfect palindromes. In the case of the vitellogenine gene ERE it has been shown that such repeated imperfect palindromes, when close to each other, cooperate to be functionally active (35). Other isolated sequences differing in 3 bp from the ERE consensus were also found: one of them lies within the direct repeat present in the promoter region. A cluster of nine consensus for PRE TGTTCACT (with at most 2 mismatches) in both orientations, are found between nucleotides -1446 and -1044. Such an arrangement resembles other functional PRE for positively controlled genes (36,13 and J.F. Savouret unpublished results). These

sequences are more divergent from the imperfect palindromic GRE consensus (36,37). Some other tandemly associated PRE consensus sequences are found in this sequence. Hormone responsive elements have been shown in many cases to form functional clusters (review in 36). Transfection and receptor binding studies will be necessary to define the functional significance of these sequences. It is of interest that two sequences matching the consensus ANCTCTCC described for silencer elements (32) are found within the cluster of the putative PRE at positions -1216 and -1142.

Between positions -1677 and -1344, is observed a sequence which exhibits a striking 61.3% homology with the 3' non coding region of the bovine acetylcholine receptor mRNA alpha subunit precursor (38).

Receptor messenger RNA heterogeneity and origin of the "A" and "B" forms of Progesterone receptor in rabbit and human cells:

Since progesterone receptors were observed in cellular extracts mainly as two different species of apparent molecular weight 110 kDa (subunit B) and 79 kDa (subunit A) (39,40) and since rabbit PR mRNA appeared on Northern blots as a doublet (1), we examined the possibility that two different messenger RNAs might have coded for the two protein species. In a previous study we have shown that the A form of receptor was derived from the B form by deletion of the N-terminal part (41). We thus prepared two radioactive probes corresponding to the N-terminal and C-terminal extremities of the protein and compared the Northern blots obtained with these probes. As shown in Fig. 4 exactly the same pattern was observed indicating that there was no discrete mRNA species lacking the N-terminal part of the coding region. Receptor mRNA heterogeneity is thus probably linked to variations in non coding regions, and especially the 3' non coding region where sequencing has shown variability in the polyadenylation site. Another explanation has been proposed for the existence of the B and A forms of receptor: a unique messenger RNA could have been alternatively translated using two different ATG to yield the two proteins (8,42). To test this hypothesis we transcribed and translated rPR cDNA in vitro (Fig. 4). This experiment yielded a major band of apparent molecular weight \sim 110 kDa (thus corresponding to the B form) and many minor shorter bands. The latter did not correspond to internal initiations of translation but to premature stops as shown by the fact that these abortive peptides were immunoprecipitated by Let 126 antibody. This antibody recognizes an epitope localized at the N-terminus of the protein between aminoacids 1 and 60 (41) (see below). No discrete and quantitatively

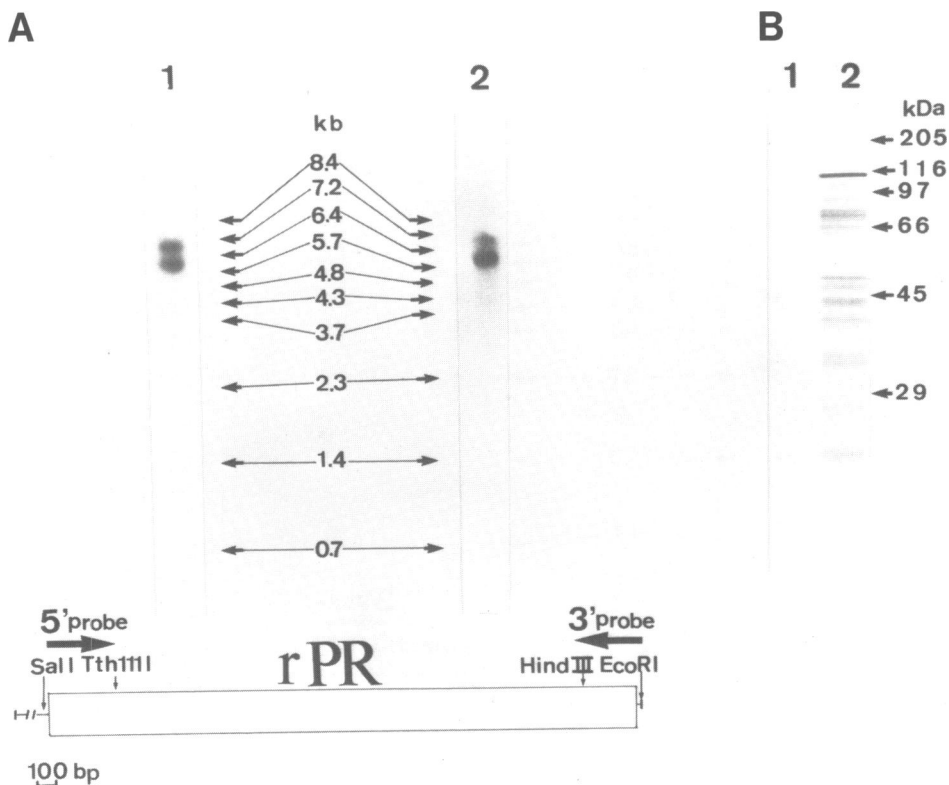


Figure 4: rPR mRNA heterogeneity and origin of the B and A forms of rabbit progesterone receptor:

Figure A: Northern blot experiments with rabbit uterine mRNAs (see Material and Methods). 1: products of endogenous translation (without addition of rPR mRNA). 2: products of translation of rPR mRNAs (obtained by transcription of rPR cDNA).

Size markers are shown by arrows. The coding region of rPR messenger RNA is shown in the lower part.

important band was observed corresponding to a protein of ~ 79 kDa. All these results are in good agreement with the data presented previously suggesting that the A form is derived by artefactual proteolysis from the intact rabbit progesterone receptor (B form) (43). It has been argued however that the rabbit might be a special case and that in other species, especially in humans, the A form could be physiologically produced (40). We thus repeated the latter two experiments with RNAs from human cells and

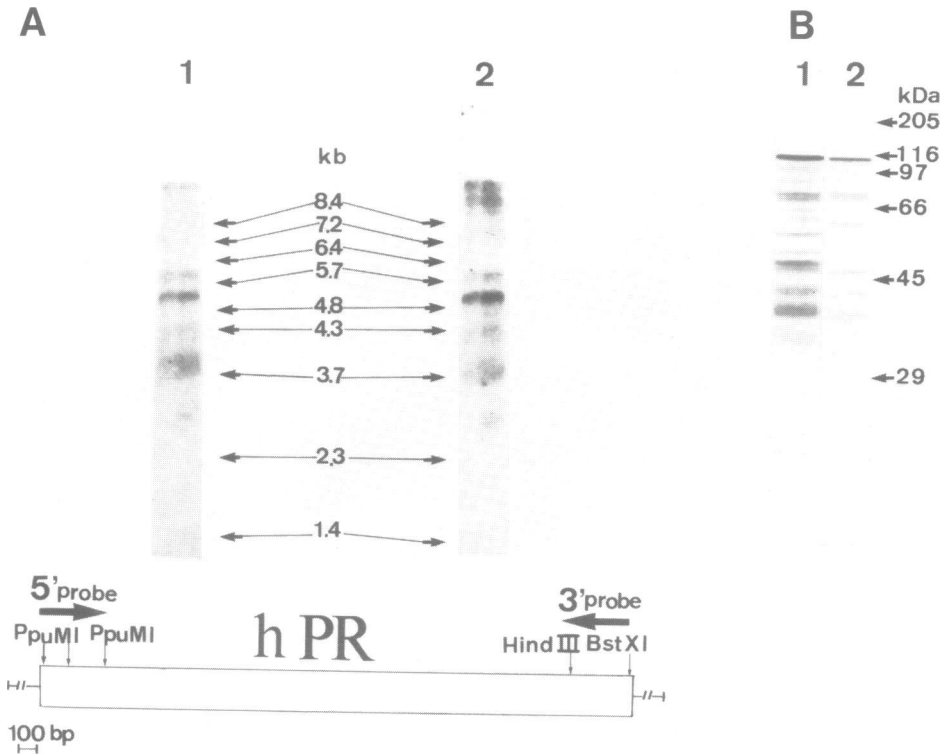


Figure 5: hPR mRNA heterogeneity and origin of the B and A forms of the human progesterone receptor:

Figure A: Northern blot experiments with human uterine mRNAs (see Material and Methods). Probes corresponding to the 5' end (lane 1) and 3' end (lane 2) of the coding region of rPR mRNA were used.

Figure B: in vitro transcription and translation of hPR cDNA (see Material and Methods). 1: product of translation of hPR mRNAs (obtained by transcription of hPR cDNA). 2: immunoprecipitation by the monoclonal antibody Let 126 of the product of translation observed in experiment 1. Size markers are shown by arrows. The coding region of hPR messenger RNA is shown in the lower part. Exactly the same patterns were observed when the hPR cDNA contained only 12 bp upstream from the initiator ATG instead of 175 bp (not shown). Let 126 immunoprecipitation of rPR messenger RNA translation products was identical to that shown for hPR (Fig. 5 B 2).

probes from human cDNA. Again no evidence was obtained for the existence of a specific messenger RNA for the A form (fig. 5). The same bands were apparent on the Northern blot using probes derived from the N-terminal and C-terminal parts of the receptor. Moreover transcription-translation experiments showed again that the first ATG was the only main initiator used during messenger translation (Fig. 5). Evidence has been obtained for the presence of a single

PR gene in various species (42,44, and M. Atger unpublished results).

The structural analysis of PR gene promoter and 5'flanking region described in this work should now be followed by functional studies in order to dissect the molecular mechanisms which direct the tissue specific expression of the receptor and its hormonal regulation. An understanding of the anomalies of receptor gene expression will be of special interest in hormone dependent cancers and in particular in breast cancer where non hormonally regulated constitutive expression or absence of expression in cancers containing estrogen receptors have been described (45).

*To whom correspondence should be addressed

REFERENCES

1. Loosfelt, H., Atger, M., Misrahi, M., Guiochon-Mantel, A., Mériel, C., Logeat, F., Benarous, R. and Milgrom, E. (1986) Proc. Natl. Acad. Sci. USA **83**, 9045-9049.
2. Misrahi, M., Atger, M., d'Auriol, L., Loosfelt, H., Mériel, C., Fridlansky, F., Guiochon-Mantel, A., Galibert, F. and Milgrom, E. (1987) Biochem. Biophys. Res. Commun. **143**, 740-748.
3. Vu Hai, M.T., Logeat, F., Warembourg, M. and Milgrom, E. (1977) Ann. N.Y. Acad. Sci. **286**, 199-209.
4. Hollenberg, S.M., Weinberger, C., Ong, E.S., Cerelli, G., Oro, A., Lebo, R., Thompson, E.B., Rosenfeld, M. G. and Evans, R.M. (1985) Nature **318**, 635-641.
5. Miesfeld, R., Rusconi, S., Godowski, P.J., Maler, B.A., Okret, S., Wikstromn A.C., Gustafsson, J.A. and Yamamoto, K.R. (1986) Cell **46**, 389-399.
6. Green, S., Walter, P., Kumar, V., Krust, A., Bornert, J.M., Argos, P. and Chambon, P. (1986) Nature **320**, 134-139.
7. Arriza, J.L., Weinberger, C., Cerelli, G., Glaser, T.M., Handelin, B.L., Housman, D.E. and Evans R.M. (1987) Science **237**, 268-275.
8. Gronemeyer, H., Turcotte, B., Quirin-Stricker, C., Bocquel, M.T., Meyer, M.E., Krozowski, Z., Jeltsch, J.M., Lerouge, T., Garnier, J.M. and Chambon, P. (1987) EMBO J. **6**, 3985-3994.
9. Owen, D. and Kuhn, L.C. (1987) EMBO J. **6**, 1287-1293.
10. Okret, S., Poellinger, L., Dong, Y. and Gustafsson, J.A. (1986) Proc. Natl. Acad. Sci. USA **83**, 5899-5903.
11. Shaw, G. and Kamen, R. (1986) Cell **46**, 659-667.
12. Jantzen, H.M., Strahle, U., Gloss, B., Steward, F., Schmid, W., Boshart, M., Miksicek, R. and Schutz, G. (1987) Cell **49**, 29-38.
13. Bailly A., Le Page, C., Rauch, M. and Milgrom, E. (1986) EMBO J. **5**, 3235-3241.
14. Maniatis, T., Fritsch, E.F. and Sambrook, J. (1982) Molecular Cloning: A Laboratory Manual (Cold Spring Harbor Laboratory, Cold Spring Harbor, NY).
15. Sanger, F., Nicklen, S. and Coulson, A.R. (1977) Proc. Natl. Acad. Sci. USA **74**, 5463-5467.
16. Loosfelt, H., Logeat, F., Vu Hai, M.T. and Milgrom, E. (1984) J. Biol. Chem. **259**, 14196-14202.
17. Perrot-Applanat, M., Groyer-Picard, M.T., Lorenzo, F., Jolivet, A., Vu Hai, M.T., Pallud, C., Spyrtos, F. and Milgrom, E. (1987) Cancer Res. **47**, 2652-2661.

18. Giguère, V., Yang, N., Segui, P. and Evans, R.M. (1988) *Nature* **331**, 91-94.
19. Azorin, F. and Rich, A. (1985) *Cell* **41**, 365-374.
20. Rogers, J. (1983) *Nature* **305**, 101-102.
21. Kozak, M. (1987) *Nucleic Acids Res.* **15**, 8125-8132.
22. Petkovich, M., Brand, N.J., Krust, A. and Chambon, P. (1987) *Nature* **330**, 444-450.
23. Weinberger, C., Thompson, C.C., Ong, E.S., Lebo, R., Gruol, D.J. and Evans, R.M. (1986) *Nature* **324**, 641-646.
24. Sap, J., Munoz, A., Damm, K., Goldberg, Y., Ghysdael, J., Leutz, A., Beug, H. and Vennstrom, B. (1986) *Nature* **324**, 635-640.
25. Krust, A. Green, S., Argos, P., Kumar, V., Walter, P., Bornert, J.M. and Chambon, P. (1986) *EMBO J.* **5**, 891-897.
26. De Thé, H., Marchio, A., Tiollais, P. and Dejean, A. (1987) *Nature* **330**, 667-670.
27. Yaniv, M. (1984) *Biol. Cell* **50**, 203-216.
28. Dynan, W.S. and Tjian, R. (1985) *Nature* **316**, 774-778.
29. Davidson, E.H., Jacobs, H.T. and Britten R.J. (1983) *Nature* **301**, 468-469.
30. Cheng, J.F., Printz, R., Callaghan, T., Shuey, D. and Hardison, R.C. (1984) *J. Mol. Biol.* **176**, 1-20.
31. Suske, G., Wenz, M. Cato, A.C.B. and Beato, M. (1983) *Nucleic Acids Res.* **11**, 2257-2270.
32. Baniahmad, A., Muller, M., Steiner, Ch. and Renkawitz, R. (1987) *EMBO J.* **6**, 2297-2303.
33. Klein-Hitpass, L., Schorpp, M., Wagner, U. and Ryffel, G.U. (1986) *Cell* **46**, 1053-1061.
34. Seiler-Tuyns, A., Walker, P., Martinez, E., Méritlat, A.M., Givel, F. and Wahli, W. (1986) *Nucleic Acids Res.* **14**, 8755-8770.
35. Martinez, E., Givel, F. and Wahli, W. (1987) *EMBO J.* **6**, 3719-3727.
36. Yamamoto, K.R. (1985) *Ann. Rev. Genet.* **19**, 209-252.
37. Strahle, U., Klock, G. and Schutz, G. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7871-7875.
38. Noda, M., Furutani, Y., Takahashi H., Toyosato, M., Tanabe, T., Shimizu, S., Kikuyotani, S., Kayano, T., Hirose, T., Inayama, S. and Numa, S. (1983) *Nature* **305**, 818-823.
39. Birnbaumer, M., Schrader, W.T. and O'Malley, B.W. (1983) *J. Biol. Chem.* **258**, 7331-7337.
40. Horwitz, K.B. (1987) *J. Steroid Biochem.* **27**, 447-457.
41. Lorenzo, F., Jolivet, A., Loosfelt, H., Vu Hai, M.T., Brailly, S., Perrot-Appianat, M. and Milgrom, E. *Eur. J. Biochem.*, in press.
42. Huckaby, C.S., Conneely, O.M., Beattie, W.G., Dobson, A.D.W., Tsai, M.J. and O'Malley B.W. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 8380-8384.
43. Logeat, F., Pamphile, R., Loosfelt, H., Jolivet, A., Fournier, A. and Milgrom, E. (1985) *Biochemistry* **24**, 1029-1035.
44. Rousseau-Merck, M.F., Misrahi, M., Loosfelt, H., Milgrom, E. and Berger, R. (1987) *Hum. Genet.* **77**, 280-282.
45. McGuire, W.L. (1980) *Recent Prog. Horm. Res.* **36**, 135-156.