
Compilation and analysis of sequences upstream from the translational start site in eukaryotic mRNAs

Marilyn Kozak

Department of Biological Sciences, University of Pittsburgh, Pittsburgh, PA 15260, USA

Received 15 November 1983; Accepted 21 November 1983

ABSTRACT

5'-Noncoding sequences have been tabulated for 211 messenger RNAs from higher eukaryotic cells. The 5'-proximal AUG triplet serves as the initiator codon in 95% of the mRNAs examined. The most conspicuous conserved feature is the presence of a purine (most often A) three nucleotides upstream from the AUG initiator codon; only 6 of the mRNAs in the survey have a pyrimidine in that position. There is a predominance of C in positions -1, -2, -4 and -5, just upstream from the initiator codon. The sequence CCG^ACCAUG(G) thus emerges as a consensus sequence for eukaryotic initiation sites. The extent to which the ribosome binding site in a given mRNA matches the -1 to -5 consensus sequence varies: more than half of the mRNAs in the tabulation have 3 or 4 nucleotides in common with the CCACC consensus, but only ten mRNAs conform perfectly.

INTRODUCTION

Two years ago I prepared a compilation of the then-available 5'-noncoding sequences of eukaryotic mRNAs (Curr. Topics Microbiol. Immunol. 93, 81-123, 1981). Apart from a bevy of globin and histone mRNAs, only 32 other cellular mRNA sequences were known at that time. In contrast, there are 166 cellular mRNAs in the present compilation, not counting the globins and histones. I have excluded the mRNAs of lower eukaryotes and viruses, only to keep the survey manageable. One of my objectives was to determine whether certain patterns noted in the earlier compilation would be evident with this larger, more diversified set of sequences.

A few points about the selection and presentation of the sequences require explanation. In cases where numerous representatives of a gene family have been sequenced, I have omitted many and chosen those in which the leader sequences show the most divergence. There are exceptions, however. It seemed useful to include certain pairs of mRNAs in which the leader sequences show extensive homology *except near the AUG initiator codon* (e.g. human versus rat preproinsulin). The opposite pattern is also provocative; i.e., sequence conservation *only* near the AUG codon, as in human versus rat immunoglobulin E. Upon inspecting the completed compilation, only two families of mRNAs appeared to be excessively

represented: histones and globins. The 5'-noncoding sequences of histone mRNAs are sufficiently varied that they pose little danger of distorting the search for homology among ribosome binding sites. This is less true of the globin sequences, and I have controlled for this as described later in the text.

Nucleotide sequences determined by analyzing genomic DNA have been included only when there is sufficient supplementary data to identify introns that lie upstream from the AUG codon (or to verify their absence) and to estimate the location of the 5'-end of the mRNA. The 5'-end of each mRNA in the table was identified according to one of the following criteria:

- (a) Direct sequence analysis of the purified mRNA.
- (b) Primer-extension and/or mapping with a single-strand specific nuclease, such as S1. With these techniques there is often a 2- to 4-nucleotide ambiguity in pinpointing the cap site.
- (c) Termination of the longest cDNA clone. When the cDNA clone is known to stop significantly short of the 5'-end of the mRNA, the sequence in the table is preceded by an ellipsis (...).
- (d) Sequence homology with the corresponding gene from a closely-related species in which the 5'-terminus of the mRNA has been mapped.
- (e) Presence in the genomic DNA sequence of an appropriately-positioned TATA box, 25- to 30-nucleotides upstream from the presumptive cap site. In the absence of other supporting data this criterion is rather weak.

The following criteria, identified by code letters in the rightmost column of the table, were used to identify the AUG initiator codon in each message:

- (a) The nucleotide sequence corresponds to the known N-terminal amino acid sequence of the *primary* translation product. In some cases amino acid and nucleotide sequence data were derived from two different but related organisms.
- (b) The N-terminal amino acid sequence has been determined only for the *mature* protein, which is known (or presumed) to derive from a precursor that carries an N-terminal extension (the "signal peptide") of 15 to 30 amino acids. The indicated AUG triplet is the only candidate initiation site compatible with the synthesis of such a precursor.
- (c) The nucleotide sequence has a single open reading frame which either corresponds in size to the known molecular weight of the encoded protein, or includes peptides that are known to be present in the mature protein.
- (d) The indicated AUG triplet occurs at the beginning of the longest open reading frame, but the exact size of the primary translation product is not available for comparison. This criterion is rather weak.
- (e) The initiation site was deduced from sequence homology with the corresponding gene from a closely-related species in which the start site has been defined.
- (f) Under conditions that allow formation of initiation complexes *in vitro*, the indicated AUG triplet was protected by ribosomes against nuclease digestion.

In 13 of the mRNAs in the table the functional initiator codon has not been definitively identified; the structure of the encoded protein is compatible with initiation at either of two nearby AUG triplets. In such cases I have *predicted* which AUG is most likely to be the (major) initiation site. Those entries are marked with an asterisk in the rightmost column. The AUG initiator codon was predicted based on position (i.e., proximity to the 5'-end of the mRNA) and conformity to

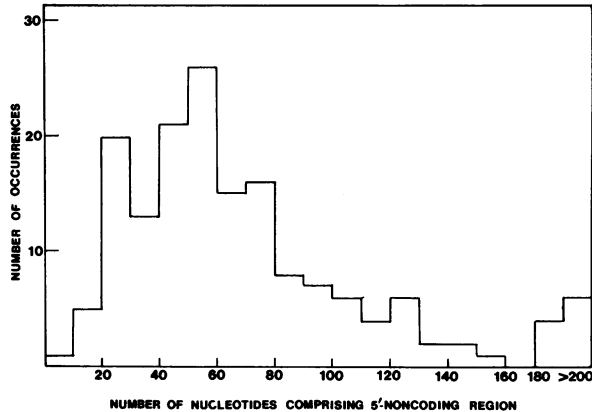


Figure 1. Length distribution of the 5'-noncoding portion of eukaryotic mRNAs. To avoid weighting the distribution by the large number of globin mRNAs in the sequence table, I scored the globin mRNAs only once in this tally.

the consensus sequence CCACCAUG. In a later section of the text I will explain in greater detail how this was done.

DISCUSSION

A few generalizations emerge from inspection of the sequences tabulated herein.

(i) The length of the 5'-noncoding region varies widely--from 3 to 572 nucleotides. However, 70% of the leader sequences are clustered in the 20- to 80-nucleotide range, as shown in Figure 1. The unusually long leader sequences occur on unusually interesting mRNAs (epidermal growth factor, oncogenes, heat shock proteins), inviting speculation that the structure of the 5'-noncoding region participates in the regulated expression of those genes.

(ii) Translation begins at the 5'-proximal AUG triplet in 95% of the mRNAs tabulated herein. There are only ten mRNAs listed in which one or more AUG triplets occur upstream from the recognized initiation site. The number of "non-functional" upstream AUG codons in each of those messages is shown in parentheses at the right edge of the table. [The upstream AUG's are called "nonfunctional" because there is as yet no evidence that ribosomes recognize those sites, but theory predicts that ribosomes should initiate (inefficiently) at the upstream AUG triplets as well as at the AUG codon that heads the long open reading frame.] The number of mRNAs with upstream AUG triplets would increase to 15 if my predictions are correct about which AUG is the major site of initiation in entries 80, 141, 146, 179 and 205. I have dealt elsewhere with the question of how ribosomes get past

ENTRY NO. 1	MESSANGER RNA	SEQUENCES FROM THE 5'-NONCODING PORTION OF EUKARYOTIC CELLULAR mRNAs ² SITE ² AUG ⁴	CAP
001	Acetylcholine receptor pre- α -subunit (torpedo)	[186]...GTTATTAGAAGTGGCAGATTGCTTGAAGGCCAAATTAATGAAAGCTGAAGAAATGATTCCTG	c b
002	Acetylcholine receptor pre- γ -subunit (torpedo)	[119]...CCCTCACCACAGGACTCACTCACAAAGCTGAGCTAGCACACTCACTGACGACCACTGACTG	c b (1)
003	Acetylcholine receptor pre- δ -subunit (torpedo)	[456]...GACTGCAAATGCTATACTGCAACAGTAAATTCGACACCTGAGCTCTTTCAAAAATGSGGAAAC	c b (8)
004	α_1 -Acid preglycoprotein (rat)	...CTCTTCTGGGCGGGTGCTCTGAGTGTCTTGGGCATGGCGCTG	*
005	α -Actin, skeletal muscle (chicken)	[73]...GAGGCGTCTGCCAGGGCCGAGCGGCTGAGTGTGAC	b a
006	α -Actin, skeletal muscle (human)	[103]...GGCCCGAGCCGAGAGTAGCAG/TTGTAGCTACCCGCCAGAAACTAGACACAAATGTGGCAG	c a
007	α -Actin, skeletal muscle (rat)	[71]...CTTCCTACCCCTGGCACCAGGGCCAGAGTGTGAGTGTGAC	b a
008	β -Actin, cytoplasmic (rat)	[80]...ACAACCTCCTTCGAGCTCCTCCGTCGCGGTCACACCCGCCACCAG/TTCCGCCATGGATGAC	b a
009	Actin, gene 79B (Drosophila)	[~150]...CTGTTGTACCCCTTGTTACCCCTTGTTACCGCCCGCACCAAACTAACCAAAACAATGTGTGAC	b a
010	Actin, gene 88F (Drosophila)	[188]...TGGAGCTAACCGTGTGCACCTCCATCTCCCTCCAGATAAAACAACCTGCCAAGATGTGTGAC	b a (3)
011	Alcohol dehydrogenase adult form (Drosophila)	[123]...TACTTAATTGATCAA/ATCGAAAGAGCCCTGCTAAAGCAAAAAGAAGTCAACCAATGTGGTTT	b a
012	pre- α -Amylase (barley)	[96]...AAGAAAAGGAGTGTCTGTACTGTAAAGTGAAGACAATGACAGTAGCGCGCCCATGGGGGAAG	c c
013	pre- α -Amylase, pancreatic (mouse)	[17] GACAACCTTCAAAGCAAAAATGAAGTTC	a, b a
014	pre- α -Amylase salivary (mouse)	[95]...GCACATG/AAATAAATTAGTTGTAGAAAGAACTACTGCCAACAGCATAGCAAAAATGAAATTC	b a (1)
015	preAngiotensinogen (rat)	[205] CCAAAATG/ [61]...GCTTGTCTGGGCTGGAGCTAAAGGACACAGAAAGCAAGTCCACAGATCCGATGACTCCC	(1) c b
016	preAntithrombin III (human)	[~67] ...GATCACATATCTCCACTTCCCGCCCTGTTGGAAAGATTAGCGGCCCATGTATTCC	b b, c
017	preProcalcitonin (rat)	[68]...CATCAGGACCCGGCAGTCTCAGCTCCAAGTCTCGCTCACCCAG/GGAGGGCATCATGGGCTTT	c c
018	Calmodulin (chicken)	[91]...CGCACCGTGCGCGGTTAGTCTCCACCAGCCCGCAGCCGGGCGGAGCCACCATGGCTGAT	c a
019	pre ϵ -Casein (mouse)	...TAGGAAGCAAGGACTCATTAAATCATGAAGTTC	c e

020	pre α -Casein (rat)	[61]...GATCATCTCCAGCTTCTCTCACCCCTACTCTTCAAGATCTTGAAGCATGAAACTT	c a
021	pre β -Casein (rat)	[52] ATCCCTGAGCTTTCATCTTCTCTTGTCCGCTAAAGGACTTGAACAGCCATGAAGGTC	b a
022	pre γ -Casein (rat)	[56]ATCATCATCTACCTATTCCTCTGCTTCCACTTGGGAAGCAAG ⁵ ATCAAGTAACCATGAAGTTC	b, e a
023	preChorion proteins (silkworm)	...CTGAATATCCAGCATCATGTTTACC	c a, c
024	-gene 292a, A-family	ATCATTCTAGATTCAGCAAGACTGTGTGATCAATATGCTACT	b b
025	-gene 18b, A-family	GTCACTTCGAAATTTAAATATCTCAATTCAGTCCAAACATGCTCACC	b b
026	-gene 401a, B-family	ATCATTCTCAGCTTTTGAATTTCAAAGTAAACATGAACACT	b b
027	-gene 10a, B-family	ATCATTATTGAGTTTCTCTCATACAAACAAAATGGCAGCC	b b
028	Chorionic gonadotropin pre- β -subunit (human)	...AGACAAGGCGAGGGGAGCCACCAAGGATGGAGATG	c a
029	preproChymosin (bovine)	AGCAGCGGCTGGACCCAGATCCAAGATGAGGTGT	c b
030	preproCollagen- $\alpha 2$, type I (chicken)	[133]...CCAGTCCGGGGGCTCTGCAACACAAGGAGTCTGATGTCTAGCAAGTAGACATGCTCAGC	b c (2)
031	preComplement, C3 (mouse)	[56]GCCTCTGCCACCCCTGCCCTTACCCCTTCACTTCTTCCACCTTTTTCCTCCTACTATGGGACCA	b b, c
032	preConalbumin (chicken)	[76]...CTGTGACCAACACCGCTGCCCTCTCAACACCAGCTGCCCTGCCCAACATGAAGCTC	a a
033	preproCorticotropin releasing factor (ovine)	[>127]...AGCGCTGGCCCTCGCTCACCTGCAGAAGCACCTCAGAAG/CGCCCCCTAAACATGCGACTG	c d
034	preCuticle protein I (Drosoph)	ATCAGTCAAAGTTCGTTCTCGACACAGCAAGTCAAGCAATATGTTCAAG	b b
035	preCuticle protein III "	ATCAGTCTTAGAAGATTTCTAGTCGGACAATCCACCCAAATCAAAATGTTCAAG	b b
036	Cytochrome P-450 (rat)	ACTGAAGTCTACCGTGGTTACACCCAGGACCATGGAGCCC	b a
037	Dihydrofolate reductase (mouse)	[115]...TTGACGGCAATCCTAGCTGAAAGGCTGGTAGGATTTTATCCCGGCTGCCATCATGGTTCGA	b a
038	preproElastase I (rat)	[21] GTGGTACTACTCTCTCCACAACATGCTGCGC	b b, c
039	preproElastase II (rat)	ACAGACATCCACGSAACACACCATGATCAGG	b, c c
040	preproEnkephalin (human)	[129]...CTGAACCCGGCTTTTCCAAATGGCTGCTCCATCCGAACAGGGTCAAC/TCATGGGGCGG	b c
041	preproEpidermal growth factor (mouse)	[353]...TCAGAGGCTCTCGAGAGGTGCAG ⁵ AGGACCTGGAAAGGCGAGCTAAATAAAAGATGCCCTGG	b, c d
042	Fatty acid binding protein (rat liver)	[39] CTGTTGGTGGCAGCTGGGAAAGGAAACCTCATTGGCCACCATGAACCTC	c a
043	pre α -Fetoprotein (human)	[42] ATGTGCTTCCACCACCTGCCAATAACAAAATAACTAGCAACCATGAAGTGG	b c
044	pre α -Fetoprotein (mouse)	[44] ACATCCCACTTCCAGCACTGCCGGTGAAGGAAC ⁵ AGCAGCCATGAAGTGG	b, c a
045	preFibrinogen, A α (human)	[>60]...TCCTTTCTTTCAGCTGGAGTGTCTCAGGAGCCAGCCGCCACCTTTAGAAAAGATGTTTTCC	c *

046	prefibrinogen, γ (rat)	[53]	AGAGGTCACAGTCTG66CTGTAAGGGGCTGGAGACACCG6TCAACCAGACACTATGAATTGG	b, e
047	prefibroin (silkworm)	[24]	ATCAGTTGGTTCOAACCTCTCAAGATGAGAGTCA	a, b, c
048	preprogastrin (porcine)	[61]	...ACTGAGGACCAAG66COAACAGCAGCACACTGTCCTCCAGCTGTCGAGTCAAGATGCAAGCGA	c (1)
α-Globin family:				
049	-chicken, embryonic (π^1)	[55]	ACAACTGCTCTGGGTTCTCACTGAAGGAGCCTGAGCCAGCACTCTCTCTGACAAATGGCACTG	d, e
050	-duck, major adult (α^A)	[36]	ACCCGTGCTGGGGCTGCCAAAGCCGGAGCTGCAACCAATGGTGCTG	b
051	-duck, minor adult (α^B)	[42]	ACAGAAACCGTCTAGTTGCGAGCTGCCAGCGCGCTCCGCAATGCTGACC	a
052	-human, adult	[37]	ACTCTTCTGTTCCACAGACTCAGAGAAAGCCCACTGGTGCTG	a
053	-human, embryonic (ζ)	[55]	ACCAAGCCAGCTCTGAGCAGCCCAACTCCAGTGCAGCTGCCACCCTGCCCAATGCTCTG	b
054	-mouse, adult	[32]	ACTTCTGATTTCTGACAGACTCAGGAAGAAACCAATGGTGCTC	a
055	-rabbit, adult	[36]	ACACTTCTGTTCCAGTCCGACTCCGACTGAGAAGAAACCACTGGTGCTG	a
056	-Xenopus, major adult		...TGCACACACAAACAGCAACCAATGCTTCTT	c
β-Globin family:				
057	-chicken, adult	[77]	...GAGCCAGACCTCTCCGTAACGACACACAGCTACCTCCAAACCGCGCCATGGTGGAC	b
058	-duck, adult	[~60]	...AGCCGAGACCTCTCCGTAACGACACACAGCTACCTCCGCGCACACCAATGGTGGAC	c
059	-chicken, embryonic (ρ)	[45]	AGCTCTGAGTGTCCACAGCGCCAGCCAAACCCCGTCCCACTGGTGAC	c
060	-goat, adult (β^A)	[52]	ACACTTGTCTTGACACAAACCGTGTCTCACTAGCAGCTCACAAACAGACACCAATGGTGAC	d, e
061	-human, adult	[50]	ACACTTGTCTTGACACAAACCGTGTCTCACTAGCAGCTCACAAACAGACACCAATGGTGAC	a
062	-human, fetal (γ)	[53]	ACACTGCTCTGGAAAGCTGAGGTTATCAATAAGCTCTTAGTCCAGACGCCATGGTGCAT	c, d
063	-human, embryonic (ϵ)	[53]	ATATCTGCTTCGACACAGCTGCAATCACTAGCAAGCTCTCAGGCCCTGGCATCATGGTGAT	d
064	-mouse, major adult	[52]	ACATTTGCTTCTGACATAGTTGTTGACTCACAAACCCAGAAACAGACATCATGGTGAC	a
065	-mouse, embryonic	[52]	ACTTGGCTTCTGACACTCTGTGATCAGCAGCAACTCCAGACTTGCATCATGGTGAAC	a
066	-rabbit, adult	[53]	ACACTGCTTCTTGACAAACTGTTGTTTACTTTGCAATCCCGCAACAGACAGAAATGGTGAT	a
067	-rabbit, embryonic (83)	[62]	...TCTGAGACATCTGAGACTATCAGCAAGCTCAGCAGCTCTAGCCAGACATCATGGTTGAT	b
068	-Xenopus, major adult	[46]	ACTTGTCTTTTTCGAAAGCTCAGAATAAAGCTCAACTTTGGCCATGGGTTTG	b, c
069	-Xenopus, larval		...TAGCAGCTACTCCCTTACAGCCACCAATGGTGAC	c
070	pre α_2 Globulin (rat)	[~68]	...CCATCAGCAGAGAGATTGTCGACACAGAGGCAATTCATTCCTTACCACCAATGAAGCTG	c
071	preproglucagon I (anglerfish)	[~58]	...AGGAACATAACAGCACTATTTGAGGGAGAAAAAGATAAATACGGTTGTAACACATGAACCGC	c
072	preproglucagon II (")		...GAAAGCTCAAAACAATGCAAGT	c
073	preproglucagon (hamster)	[103]	...TGCACCTGCTCACCTGCTCTCCGCTCAGTCACAGCAGCAGGAGCAAAAAAATGAAGAAC	c
074	preGluc proteins (Drosophila)	[13]	TTCAAAAGTCAAGATGGCGTTG	b
075	-Sgs-4	[29]	ATCAGTTTTGTGGAGAATTAAGTAAAAAACAATGAAGCTG	c, d
076	-Sgs-3	[33]	ATCTGGTAAAGTAGTCTCAATCTAAGATAGAACCAATGAACCTG	b, d
077	-Sgs-7	[33]	ATCTGGTAAAGTTAAITTTGTTAAAGCAACCACTGAAGCTG	b, d
078	Glycoprotein hormones (human)	[100]	...CAGTCAACCGCCCTGAACACACTCTCTGCAAAAAAGCCCGCAGAGAAAG/GAGCGCCATGGATTAC	d, e
	pre- α -subunit			a

079	Glycoprotein hormones (mouse)[~100]...AATCACTGCCAGACACATCCCTCAAAAAGTCAGAGCTTCAGAAAGAGCTATGGATTAC	b
	pre- α -subunit	c
080	pregrowth hormone (bovine) [61]...AGGACCCAGTTCACCAGAGCACTCAGGTCCTGTGGACAGCTCACCAGCTATGATGGCTGCA/d,e	*
081	pregrowth hormone (human) [60]...AAGGCCCACTCCCGAACCACTCAGGTCCTGTGGAGCTCAGCTAGCTGCAATGGCTACA b,e	b
082	pregrowth hormone (rat) [60]...TCGAGCCAGATTCAAAACACTCAGGTCCTGTGGACAGATCACTGAGTGGCGATGGCTGCA b,e a,b	b
083	Heat shock (Drosophila):	
	-70K protein [247]...GAGAACTCTGAATACATTTTCAACAAGTCGTTACCGAGGAAGAATCACACAAATGGCTGTCT b,e	c
084	-22K protein [252]...ATTGGCATAAGAAGCTTTATTTGAAAAACCCAAAGTTACCTTATCAACTCAAAATGGCTGTCT b,e	c
085	-23K protein [112]...AGAAGTTTATCTTTGAAAGGAAATCATCTTGAAGCAATAAAAACAACAAAATGGCAAAAT b,e	c
086	-26K protein [181]...AGAAAAATATTTCAACTTCGCAAGGAACTAACTTAAGGAAAAAGTAAAAATGTGGCTA b,e	c
087	-27K protein [121]...AAAAATTCCTTGTCTAGACAGGTTGTGAATAAGAGAAAAAAAATCAAAAATGTCAAT b,e	c
	Histocompatibility (MHC) antigens:	
088	-Class I: mouse (pre)H-2K ^d [25]	AAGTCGCTAAATGCCGACCCAGTGGATGGCACCC c,e
089	-Class II: mouse (pre)E _A 34k [48]	GCCTCAGTCTGGGATCGCTTCTGAAACCAACCCAAAGAAAGAAATGGCCACA b
090	- " mouse (pre)I-A _g 29K	...GGCATTTACCTGGCTTAGAGATGGCTGT -
091	- " human (pre)DR _g 29K	...CTCCTCTGGCCCTGGTCTCTCTTCTCCAGCATGGTGTT c
	Histones:	
092	-chicken H1, embryonic [38]	AGTGGTCCCCCGATCTGTGGAAAGAGTCGTCACCATGTGGAG b
093	- " H2A.1, adult [146]...GTTGCGCGTGGTTCCTGGCTGTGGTCTCTAGTGTTCAGTCGCTGGATGTGGGG e	a
094	- " H2A.F, embryonic [70]...GAGCGGTTGAGCGGGAATTCGGGACCGGGCTCGGGCGGCGGCCATGGCAGT c	c
095	- " H2B [36]	ACTCGCTGCGCGGAGGGATCTGGAGAGTTGACATGGCTGTGAG e
096	- " H4, embryonic [26]	GGTGGAGCACTCAGGCTCTGGCATGTCTGGC b
097	- " H5, erythrocyte [109]...CTTTTTAAGCTCCCTAACCCAGTGCCTCGGGTGAAGCGGGGCCATGACGGAG c,e	a
098	-human H3 [27]	ACTTTTGGTTGACAGCAGTTCTGCGAAATGGCGGT e
099	- " H4 [36]	ACAAAGCGTTGGGTGAGAGCCCTCTTGTCTGTCTGATGTCTGGC e
100	-sea urchin, early, H1 [38]	ATTGGTTTGAACCTCCGACGCCGTATATCAAGATGGCTAG b
101	(S. purpuratus) H2A [70]...TTGTAGTCTGAACCTGGCTTCGATTTTCAAACTATCAAAACATCATGTCTGGC b	a
102	H2B [78]...CTTGACACTGTTCGGTCTTTACAGCAAAAACCTCAATTCATCATGGCTCA b	a
103	H3 [55]	ATTCATCCGCTCACTGTTTGAAGTACTGAATTAACCTGTCGCAAGCAACTATGGCAGC b
104	H4 [67]...ATATCCAGTCTACAGGATCACAGCACTGCTCAACTATCAATCATCATGTCCAGT a,b	a
105	-sea urchin, early H2A [60] ...CGTGTGTTCACAAGCTGGTCTGCTCAGCTGTGTTAAACCAACCACTCATGTCTGGA b	a
106	(P. miliaris, h2 ²) H2B [76]...TCGTGACCTCGCATGGTCTCTAAGACATCAGAAAACTTCATCTCACCATG b	a
107	H3 [54]	ATTCATCTGTCACCCCTGTTTGAACACTGATCAACTCTCCCAATCAACCATGGCAGC b
108	H4 [60]...AGCGAAACGTCAGTGGTCAGCATGCGCAATGAACTCTCTCAACTCCATAATGTCCAGC d	a
109	-sea urchin, late H3 [~27]	ATCAGTTGACTATCGAAAAATCAACATGGCTGTGGA b,e
110	(L. pictus) H4 [~23]	ATCTCAAAAGAAACTTCAAAATGGCTGTGGA b,e
111	-Xenopus H1 [28]	AGTTGGCTGAGTAAATTTCAAAAGATGACAAA e
112	" H2A [47]	AGTCTACAACATCTTGTGATTTGATTTGTAGCACAGTAATCATGTCTGGA e
113	" H2B [35]	ACAGTTTTGTAGGCTGAGAGAAAGCAGCAATTTATGGCTGAA e
114	" H3 [80?] ...CGGGTTACCCGGTTCACAGCTTAGGCTTTCTTAACTGATACATATGGCCGT e	a
115	" H4 [28]	ATATTGTTTCAAGAGCTCAAGAAAGATGTCTGGA e

148	apoLipoprotein II, VLD, chick [77]...GAAAGCAGGACAG/GTCTCTTGGTGAAGGGCTGAACCTGGTACCAACAACAAACCATGGTGCAA	b	a	
149	apoLipoprotein A-I (human)		a	
150	apoLipoprotein E (rat)	...TCCCCACGGCCCTTCAGGATGAAAAGCT	c	
151	preLysozyme (chicken)	...ACTGGCCAATCACAACTGGGAAGATGAAAGGCT	c	
152	preProMelittin (honeybee)	AGTCCCGCTGTGTGTAGCACACTGGCAACATGAGGTCT	b	
153	Metallothionein-II (human)	...AGCGAATTAAACAGATTAAACAGGAAGGAAGGACGATCGGAGAATAATCATGAAAATTC	c	
154	Metallothionein-I (mouse)	[69]...CCAGCGAACCGCGTGAACCTGTCCCGACTAGCGCCCTTTCAGCTCGCCATGGAATCCC	b	
155	pre-β2-Microglobulin (mouse) [52]...ATTTTCAGTGGCTCTACTCGGGCTTACCGTAGCTCCAGCTTACCAGATCTCGGAAATGGACCCC	b	a	
156	Myoglobin (seal)	[70]...CAGGACACCGAGTCAGCCGGGACTTGTCTTCTTGTCTTCTCCAGACTGCACCATGCGGCTCG	c	
157	Myosin, skel. I chain (chick)	...CCTCTCAGCTAATCCCTCCGGCCCGTCCGCTACTTTTTCCAACCTCAATCATGTCTCTC	c	
158	prepro-β-Nerve growth factor [74]...GCTGGCCTTATATTGGATCTCCGGGACGCTTTTTGGAAACTCCTAGTGAACATGCTGTGC	c	d	
159	protoOncogenes: -c-fos (human)	[154]...CCCACCTGTCTCCGCCCTCGGCCCTCGGCCGCTTTCCTAACCGCCACGATGATGTTC	e	*
160	-c-myc (human)	[572] ...AGAGGCTGGATTTTTTCGGGTAGTGGAAAACCAG/CAGCCCTCCCGGACGATGCCCCCTC	c	d
161	-c-Ki-ras2 (human)	[>181]...GGGGCCAGAGGCTCAGCGGCTCCCAAGTGGGGAGAGAG/GCCTGCTGAAAAATGACTGAA	c	c
162	-c-src (chicken)	...TGGCGTACCACCTGTGGCCAGGCGGTAGCTGGGACGTGCAG/CCACACCACATGGGGAGC	-	c
163	preproPiomelanocortin (bovine)	[128]...GGAGCCGCCCGAGGCAGCTTCCCGGTGACAG/AGCCTCAGCCCTGCGTGGGAAGATGCCGAGA	b	a
164	preproPiomelanocortin (human)	[107]...CCCGCCCTCAGAGAGCAGCCTCCCGAGACAG/AGCCTCAGCCTGCTGGGAAGATGCCGAGA	d,e	a
165	Ovalbumin (chicken)	[64]...AAAGCTGATTTGCCCTTTAGCAGTCAAGCTCGAAAG/ACAACCTCAGAGTTCACCATGGGCTCC	a,b	a,f
166	preOvomucoid (chicken)	[53] ATCTCAGGAGCAGACAGCAGCGCTCCAGAGCGGGCAGTACCTCACCATGGCCATG	b,c	a
167	preproParathyroid hormone (bovine)	[100]...TCAGCTGCTAATACATTTGAAAGAAGATTGTATCTCTAAGACGTGTG/TTAAATATGATGTCT	b	a
168	preproParathyroid hormone (human)	[90]...AGCTACTAACATACCTCAAGGAAGATCTTGTTCAGACATTTGATG/TGAAGATGATACCT	d,e	a (1)
169	prePepsinogen (human)	[54] CTGCACCTTCTCCCGTTCCTCCCTTCTCCCTCGAGTTGGGACCCGGGAAGAACCACTGAAGTGG	b,e	b
170	pre(?)Phaseolin	[77]...CTACTACTACTATAATACCCAACCACTATAATTCATACTACTCTACTACTGATGAGA	c	*
171	Phosphoglycerate kinase, human [80]...GGCTCCCTCGTTGACCGGAATCACGACCTCTCTCCCGAGCTGATTTTCCAAAAATGTGGCTT	c	a	
172	prePlacental Lactogen, human	...CTGTGGACAGCTCACCTAGCGGCAATGGCTGCA	c	a

173	prePlasminogen activator (human)	[84]...	GCAGGAAGAGGAGGCAAGCCGTGAATTTAAGGGACGCTGTGAAGCAATCATG&ATG&A	c	*
174	preProlactin (rat)	[51]	AGTGGTCTCTTAGGACTCTTGGGGGAAGTGTGGTCCAGTGGTCAACACCATTGAAACAGC	c, e	a
175	preProlactin (bovine)	[67]...	GCCATAGGACGAGAGCTTCTGGTGAAGTGTCTTTGTAATCATCACCACCATTG&ACAGC	c	b
176	Protamine (trout)	[14]	ATCCATCAATCACAATG&CC&AGA	b	a
177	Pyruvate kinase (chicken)	[80]...	GCTTTGGGCACGGGGGGGAAGCAGCAGCAGGAGACACCGA&ACTCCAGTA&CC&ATG&TCGAAG	c	c
178	preproRelaxin (rat)	[~60]	...CTGAACGGCCAGGAGACACCGCCAGGAGCAGCC&CC&GGAATG&TC&CAGC	c	b
179	preproRenin (mouse)	[55]	CTGGGCTACACAGCTCTTAGAAGCCCTTGGCTGAAC&CAGATG&ACAGGAGGAGGATG&CCTC	c	*
180	preRibonuclease, panc.	(rat) [~75]...	CAATTTGCTCG&AATCAAAGCTTAGGCTCCTCAGACGACGAA&CC&ACTATG&GGTCTG	c	b
181	Ribosomal protein S19, Xenopus	[>46]	...AATTCCTTAATCCTTTCTTGTCC&CGT&AGAGATAG&CC&G&CAAGATG&AATGAC	c	c
182	preRibulose biphosphate carboxylase (soybean)	[45]	ATCTGGCAGCAGAAAA&CA&AGTAGTTGAGAACTAAGA&A&A&A&A&ATG&C&TTCC	b	a
183	Seminal vesicle pre-secre-tory protein IV (rat)	[22]	AGTCAAGAGCTTTTTCTGG&CA&AGTGAAGTCT	b	b
184	pre(?)Sericin (silkwmoth)	[54]	ATAGTGGTCTTATCATCGGGTCTTAAGGATCA&AGCGATCCA&AGACCG&CA&ACATG&CGTTTC	b	c
185	Serum preproalbumin, chick	[41]	AGCATTTTTG&A&ATAATTTAGCC&C&ACATCAATCTG&C&AGCCAATG&AAGTGG	b	a
186	Serum preproalbumin, human		...GCTTTTCTCTTCTGTC&A&CCCC&C&AC&GGCTTTTGG&C&ACAATG&AAGTGG	c	a
187	preproSomatostatin (human)	[112]...	TGCTCGGAGCAG&GG&ATATCTG&CG&CATC&AGT&GCC&CC&CG&GGTGA&AG&GATG&CC&ACTC	b	c
188	preproSomatostatin-II (anglerfish)	[59]...	CA&A&CC&C&AG&CA&AAC&C&AGTAG&A&AAC&C&AG&CA&AG&CA&C&AG&C&AG&C&AG&C&ATG&C&AGTGT	c	c
189	preproSomatostatin-22 (catfish)		...CG&CC&C&AG&CTCG&AAA&TCTTCC&C&AG&CTACC&A&A&A&AGATG&TC>CT	c	c
190	preproSomatostatin-14 "	[114]...	GCCCTCCCTCC&AC&CA&ATTTTTCC&AC&CA&A&ATCC&AGCTTTATTTCTTTTTG&A&AGATG&CC&CTCC	c	c
191	preproSomatostatin-I (human)	[105]...	CG&GC&CTAG&AGTTTG&AC&CG&C&ACTC&C&AG&CTG&GGCTTTGG&GG&CG&CC&G&AGATG&CTG&TCC	c	c
192	preproSomatostatin (rat)	[81]...	CTGGGCTAG&ACTG&ACC&C&AC&GG&CTC&A&AG&CTGG&CTGTG&TAG&GC&AG&GG&G&AGATG&CTG&TCC	c	c
193	pre-Steroid binding protein Cl, rat prostate	[44]	TG&A&AG&AGTTTCA&TTTGTCC&C&C&ATTTGCTA&AGTAG&A&A&A&A&CTG&A&A&ATG&AG&C&C	b	b
194	pre-Steroid binding protein C2, rat prostate	[42]	TGG&AG&AGT&TCC&ATTTG&CTC&AGTCT&A&A&A&A&G&A&A&A&A&CTG&AG&C&C&A&ATG&AG&GG&CTG	b	b
195	pre-Steroid binding protein C3, rat prostate	[55]	G&AGTTC&CTG&ATTTCTG&TCTGG&A&C&A&G&C&A&C&A&C&C&C&C&AG&G&G&A&CTGC&CTC&A&A&C&ATG&A&AG&CTG	b	b
196	preproThaumatin	[31]	AA&AG&CG&C&AG&C&CTCA&ATTTGG&C&ATC&AT&C&A&T&CA&ATG&C&C&CG&CC	b	c
197	Thyrotropin, pre-β-subunit (rat)	[89]...	CG&AGT&G&AG&A&A&A&A&A&A&ATTTCTG&CTT&C&AGT&G&A&AG&AG&CTGG&GGTGTTC&A&A&AG&C&ATG&AGT&G&CT	c	b
198	preproTrypsinogen (rat)	[12]	CC&TTCTG&CC&C&C&C&C&C&ATG&AGT&G&C&A	b	b

199	α -Tubulin (rat)	[100]...AACACCTCTCTCTCGCCCTCCGCCATCCACC GGGGAGCGGGGAAAGCAGCAACCATG/CGTGA	b	a
200	β -Tubulin (chicken)	[87]...AGAGCGGGAGGTGACGGAGGGGAGCGGGGACCGGCAGACACCGGCATCATCGGTGAG	c	a
201	β -Tubulin (human)	[159]...TTTTCTTGGCCCATACATACCTTGGGGAGCAAAAAAATAAATTTTAAACCATGAGGGAA	b	a
202	pre β -terroglobin (rabbit)	[47] AGATCACCGGATCCAGAGCCGCCAGAGCCCTTCCCATCTGCGCACCATGAAAGCTC	b, c	a
203	preproVasoactive intestinal polypeptide (human)	GGGGAGCAGCACTGGGCGAGGACACAGAAATGGACACC	b	c
204	preproVasopressin neuro-physin II (bovine)	[49] GCACAGTCTACAGAGCAGCAGCTGGCACGCTGTGCCACGCGTGCAGGATGCCCGAC	b	b
205	preproVasopressin neuro-physin (rat)	AGCAGAGCAGAGCTGCAGCAGTATGCTCGCCATGATGCTCAAC	b	*
206	Vitellogenin II (chick)	[13] ATTCACCTTCGCTATGAGGGGG	b	d
207	Vitellogenin I (Drosophila)	[258]...ACTCACTCAGTGTGAAGTCGCATCCGAGACCAAAATCCCAAATCCGAAACCATGAACCCC	b	*
208	Vitellogenin II (Drosophila)	[251] ATGCAGTACAATTTGGTACGGTCTGAAAAAGTGCAGTGGAAAGCCACAAATGAATCCT	b	*
209	pre β Whey acidic protein (mouse)	[26] ATCAGTCACTTGGCTGACCCGGTACCATGCGTTGC	c	c
210	preZein, 19K (maize)	[57]...CACATATTATTGAGACCAACTAGCAATATAGAAAGCACAATATTGTACCAATAATGGCAGCC	c	b
211	preZein, 22K (maize)	[>67]...TCAGCATTCAAAAACACACACCAAGCGAAGCCACTAGCAAC ² ACCTTACAC ³ CAATGGCTACC	c	c

Footnotes:

¹The number assigned here to each mRNA is used again to identify the corresponding references in the bibliography.

²The table shows the sequence of the plus strand of DNA, from which the sequence of mRNA can be derived by substituting U for T. All ATG triplets are shown in italics. The sequences are aligned by using the ATG triplet that is known (see footnote 4) or predicted (see text) to be the functional initiator codon. The positions of introns are indicated by a diagonal line. For mRNAs in which the 5'-noncoding sequence exceeds 56 nucleotides, I have shown only the portion nearest the ATG triplet. The number in brackets preceding the sequence indicates the (approximate) full length of the 5'-noncoding sequence, not counting the m/G cap or the ATG triplet. There is likely to be a 2- to 4-nucleotide uncertainty when S1 nuclease was used to map the cap site. In cases where cDNA clones were analyzed, the 5'-noncoding sequence may be a little longer than indicated. If the missing portion of the leader is suspected to be more than a few nucleotides, no figure is given for the overall length of the 5'-noncoding region.

³The criteria (a-e) used to identify the 5'-terminus of each mRNA are summarized in the text.

⁴The criteria (a-f) used to identify the ATG initiator codon are summarized in the text. An asterisk in this column means that the ATG triplet used to align the mRNA is predicted, but is not known to be the functional initiator codon. I have also marked in this column those mRNAs that have ATG triplets upstream from the functional initiation site. The number of upstream ATG's is indicated in parentheses. Entry 119 (marked +) is, to my knowledge, the only cellular mRNA in which two functional initiator codons have been identified: the ATGs in positions 4-6, and 19-21.

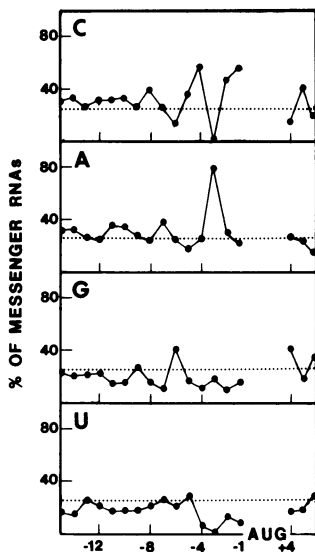


Figure 2. Frequency distribution of each nucleotide around the functional initiator codon in 198 mRNAs listed in the table. The calculations presented here do not include the 13 mRNAs in which the translational start site was predicted but has not been verified. The nucleotide immediately preceding the AUG codon is numbered -1; nucleotides +4 to +6 represent the start of the protein coding sequence. The dotted line across each panel indicates the 25% value that would be expected on a random basis. To ensure that the results were not distorted by the inclusion of too many closely-related globin mRNA sequences, I recalculated the frequency of occurrence of each nucleotide in positions -1 through -8, omitting all of the globin sequences. Although the absolute values changed somewhat (e.g. G in position -6 dropped from 40% to 36%), the rank order of nucleotides in each position remained unchanged.

the upstream AUG triplet(s) in such messages (Kozak, *Microbiol. Rev.*, 47, 1-45, 1983; Kozak, manuscript submitted). The main point to note here is that such mRNAs are rare. The "first-AUG-rule" holds for 93% to 95% of the entries in the table.

(iii) The sequences in the table have been searched manually for signs of a conserved motif that might uniquely identify AUG initiator codons. The most conspicuous conserved feature is presence of a purine (most often A) in position -3; i.e., three nucleotides upstream from the initiator codon. As illustrated in Figure 2, 79% of the mRNAs that were counted have A in that position, 18% have G, and only 3% (a total of 6 messages) have a pyrimidine in position -3. The strong preference for a purine in position -3 is peculiar to AUG triplets that serve as initiator codons. *Pyrimidines* are favored in the -3 position preceding AUG triplets that lie *upstream* from the initiation site, in those rare mRNAs that have upstream AUGs (Kozak, *Nuc. Acids Res.* 9, 5233-5252, 1981); and the nucleotide frequency in position -3 is almost perfectly random around AUG triplets that code for methionine at *internal* positions in polypeptide chains (Kozak, 1983, *op. cit.*). Although no other position is as highly conserved as the purine in position -3, the distribution of nucleotides is decidedly nonrandom in every position from -1 through -6, and perhaps beyond. The predominance of C in positions -1, -2, -4 and -5 was evident in an earlier survey (Kozak, 1981, *op. cit.*) and is confirmed here. The preference for G in position +4, noted in the previous survey,

is less evident here. From the data in Figure 2, the sequence CC^A_GCCAUG(G) emerges as a consensus sequence for eukaryotic initiation sites. The extent to which a given message matches the -1 to -5 consensus sequence varies considerably: only 10 mRNAs in the table conform perfectly to the CCACC sequence; in more than half of the mRNAs, 3 or 4 of the nucleotides directly preceding the AUG codon match the consensus sequence; about 10% of the mRNAs have a purine in position -3 but otherwise differ entirely from the -1 to -5 consensus. The 6 mRNAs in the tabulation that lack a purine three nucleotides upstream from the initiator codon do not seem to compensate by conforming closely to the other four consensus positions. Recent site-directed mutagenesis experiments (Kozak, manuscript submitted) have confirmed the importance of the purine in position -3, but there is as yet no evidence that cytosine in positions -1, -2, -4 and -5 contributes to recognition of eukaryotic initiation sites.

Obviously the (semi)conserved sequence revealed by a survey such as this need not correspond to the most favorable context for initiation, since the table includes mRNAs that vary in translational efficiency. Nonetheless, reference to the consensus sequence, *especially the highly conserved -3 position*, can be of help when searching a new mRNA sequence to locate the translational initiation site. It is important to avoid two errors when using this approach:

(a) If inspection of the sequence near the 5'-end of the mRNA were to reveal two AUG triplets that conform approximately equally to the consensus sequence, it would be incorrect to conclude that either AUG is equally likely to be the initiator codon. Because 40S ribosomal subunits most likely scan the 5'-end of the mRNA in a linear fashion (Kozak, Cell 34, 971-978, 1983), *the 5'-proximal AUG triplet is the first to be "inspected."* If the sequence preceding the first AUG triplet conforms closely to the consensus, especially if an A occurs in position -3, *the search ends there.* There are two exceptions to this rule. The first involves a small number of mRNAs in which the reading frame following the first ANNAUG sequence is short, terminating upstream from a second AUG codon *to which ribosomes seem to gain access by reinitiating!* The second exception consists of a single example: the mRNA derived from influenza B virus genome segment 6 allows ribosomes to initiate efficiently at the first and the second AUG codons, although the first AUG triplet occurs in a "good" context (ANNAUGA) and is not followed by a terminator codon (Shaw et al., Proc. Natl. Acad. Sci. USA 80, 4879-4883, 1983). I have no explanation for this at present.

(b) An AUG triplet that deviates from the consensus in the crucial -3 position *can nevertheless serve as the initiator codon.* This is evidenced by a few mRNAs in the table (entries 40, 98, 129, 133, 134, 196) and also by experimental manipu-

lation of the sequence flanking the initiator codon (Sherman et al., Cell 20, 215-222, 1980; Kozak, manuscript submitted). As a consequence of initiating at a "weak" AUG codon, however, those rare messenger RNAs are predicted to have two special properties: translation should be inefficient; and ribosomes should initiate not only at the first (weak) AUG but also at the next AUG that lies downstream. Such mRNAs should therefore have the potential to direct synthesis of two proteins. This has been shown to occur with a few viral mRNAs (Kozak, 1983, op.cit.) but it has yet to be demonstrated for cellular mRNAs.

The -1 to -5 consensus sequence detected in this survey differs from previously-suggested eukaryotic consensus sequences (Hagenbüchle et al., Cell 13, 551-563, 1978; Baralle and Brownlee, Nature 274, 84-87, 1978; Stiles et al., Cell 25, 277-284, 1981) in both its high frequency of occurrence and its constant position relative to the AUG initiator codon. None of the previously-suggested consensus sequences met those criteria. Until further experiments are carried out, it is premature to speculate about the mechanism by which flanking nucleotides might modulate recognition of the AUG initiator codon; but the temptation is irresistible. Sargan et al. (FEBS Lett., 147, 133-136, 1982) have noted an intriguing complementarity between the sequence CCACC in mRNA and the sequence GGUGG at the base of the 3'-terminal hairpin structure in 18S ribosomal RNA. The possibility of base pairing between mRNA and rRNA thus seems worth exploring. An alternative rationalization for the conserved sequence preceding the initiator codon is that it might base-pair with a complementary sequence just downstream from the AUG codon. The resulting hairpin could help to identify the initiation site. Although some mRNAs (see entries 18, 66, 151) have the potential to form a stable hairpin structure centered about the AUG codon, this is by no means universal. Moreover, comparison of closely-related sequences does not reveal compensatory changes that would preserve the potential hairpin structure.

Bibliography. The numbers used here to identify each mRNA correspond to those in column 1 of the table. Bibliographic data are presented in condensed form: first author, year, journal title, volume, first page. Personal communications are indicated by the letters *pc* after an individual's name. My thanks are extended to those individuals.

- | | |
|----------------------------------|---------------------------------------|
| 001. Noda(1982)Nature 299,793. | 011. Benyajati(1983)Cell 33,125. |
| 002. Noda(1983)Nature 302, 528. | 012. Rogers(1983)JBC 258, 8169. |
| 003. Noda(1983)Nature 301, 251. | 013. Hagenbüchle(1980)Cell 21, 179. |
| 004. Ricca(1981)JBC 256, 11199. | 014. Hagenbüchle(1981)Nature 289,643. |
| 005. Fornwald(1982)NAR 10, 3861. | 015. Ohkubo(1983)PNAS 80, 2196. |
| 006. Hanauer(1983)NAR 11, 3503. | 016. Bock(1982)NAR 10,8113. |
| 007. Zakut(1982)Nature 298, 857. | 017. Amara(1982)Nature 298,240. |
| 008. Nudel(1983)NAR 11, 1759. | 018. Putkey(1983)JBC 258, 11864. |
| 009. Sanchez(1983)JMB 163, 533. | 019. Hennighausen(1982)EJB 126, 569. |
| 010. Sanchez, op.cit. | 020. Hobbs(1982)NAR 10, 8079. |

021. Blackburn(1982)NAR 10, 2295.
022. Hobbs(1982)NAR 10, 8079;
Yu-Lee(1983)JBC 258, 10794.
023. Rodakis(1982)PNAS 79,3551.
024. Jones(1980)Cell 22, 855.
025. Jones, op.cit.
026. Jones, op.cit.
027. Jones, op.cit.
028. Fiddes(1980)Nature 286, 684.
029. Harris(1982)NAR 10,2177.
030. Vogeli(1981)PNAS 78, 5334.
031. Wiebauer(1982)PNAS 79,7077.
032. Jeltsch(1982)EJB 122,291.
033. Furutani(1983)Nature 301,537.
034. Snyder(1982)Cell 29,1027.
035. Snyder, op.cit.
036. Mizukami(1983)PNAS 80,3958.
037. Nunberg(1980)Cell 19,355.
R. Schimke, *pc*.
038. MacDonald(1982)Biochem. 21,1453.
039. MacDonald, op.cit.
040. Noda(1982)Nature 297,431.
Comb(1982)Nature 295,663.
041. Scott(1983)Science 221,236.
Gray(1983)Nature 303, 722.
042. Gordon(1983)JBC 258,3356.
043. Morinaga(1983)PNAS 80,4604.
044. Law(1981)Nature 291,201.
Eiferman(1981)Nature 294,713.
045. Kant(1983)PNAS 80,3953.
046. Crabtree(1982)Cell 31,159.
047. Tsujimoto(1979)Cell 18,591.
048. Yoo(1982)PNAS 79,1049.
049. Engel(1983)PNAS 80,1392.
050. Erbil(1982)Gene 20,211.
051. Erbil(1983)EMBO 2,1339.
052. Baralle(1977)Cell 12,1085.
Wilson(1980)JBC 255,2807.
053. Proudfoot(1982)Cell 31,553.
054. Baralle(1978)Nature 274,84.
055. Baralle(1977)Nature 267,279.
Heindell(1978)Cell 15,43.
056. Kay(1983)NAR 11,1537.
057. Dolan(1983)JBC 258,3983.
058. Hampe(1981)Gene 14,11.
059. Roninson(1981)PNAS 78,4782.
060. Haynes(1980)PNAS 77,7127.
061. Baralle(1977)Cell 12,1085.
062. Slightom(1980)Cell 21,627.
063. Baralle(1980)Cell 21,621.
064. Konkel(1978)Cell 15,1125.
Baralle(1978)Nature 274,84.
065. Hansen(1982)JBC 257,1048.
066. Baralle(1977)Cell 10,549.
Efstratiadis(1977)Cell 10,571.
067. Hardison(1981)JBC 256,11780.
068. Patient(1983)JBC 258,8521.
069. Banville(1983)JBC 258,7924.
070. Laperche(1983)Cell 32,453.
071. Lund(1982)PNAS 79,345.
072. Lund(1983)JBC 258,3280.
073. Bell(1983)Nature 302,716.
074. Muskavitch(1982)Cell 29,1041.
075. Garfinkel(1983)JMB 168,765.
076. Garfinkel, op.cit.
077. Garfinkel, op.cit.
078. Fiddes(1979)Nature 281,351.
Fiddes(1981)JMAG 1,3.
079. Chin(1981)PNAS 78,5329.
080. Miller(1980)JBC 255,7521.
Woychik(1982)NAR 10,7197.
081. DeNoto(1981)NAR 9,3719.
082. Page(1981)NAR 9,2087.
083. Török(1980)NAR 8,3105.
Ingolia(1980)Cell 21,669.
084-087. Ingolia(1981)NAR 9,1627.
Southgate(1983)JMB 165,35.
088. Lalanne(1983)NAR 11,1567.
Kvist(1983)EMBO 2,245.
089. Mathis(1983)Cell 32,745.
090. Malissen(1983)Science 221,750.
091. Long(1983)EMBO 2,389.
092. Sugarman(1983)JBC 258,9005.
093. D'Andrea(1981)NAR 9,3119.
094. Harvey(1983)PNAS 80,2819.
095. Grandy(1982)JBC 257,8577.
096. Sugarman(1983)JBC 258,9005.
097. Ruiz-Vazquez(1982)NAR 10,2093.
Krieg(1983)NAR 11,619.
098. Clark(1981)NAR 9,1583.
099. Heintz(1981)Cell 24,661.
100-104. Sures(1980)PNAS 77,1265.
105-108. Busslinger(1980)NAR 8,957.
Hentschel(1980)Nature 285,147.
109-110. Childs(1982)Cell 31,383.
111. Turner(1983)NAR 11,4093.
112-113. Moorman(1982)FEBS Lett.144,235.
114-115. Moorman(1981)FEBS Lett.136,45.
116. Jolly(1983)PNAS 80,477.
117. Konecki(1982)NAR 10,6763.
118. Selsing(1981)Cell 25,47.
119. Kelley(1982)Cell 29,681.
120. Bernard(1978)Cell 15,1133.
Picard(1983)PNAS 80,417.
121. Early(1980)Cell 19,981.
122. Kataoka(1982)JBC 257,277.
123. Kenten(1982)PNAS 79,6661.
124. Hellman(1982)NAR 10,6041.
125. Perler(1980)Cell 20,555.
126. Bell(1980)Nature 284,26.
127. Lomedico(1979)Cell 18,545.
Cordell(1979)Cell 18,533.
128. Hobart(1980)Science 210,1360.
129. Chan(1981)JBC 256,7595.

130. Sorokin(1982)Gene 20,367.
131. Goeddel(1980)Nature 287,411.
Lawn(1981)PNAS 78,5435.
132. Nagata(1980)Nature 287,401.
- 133-134. Goeddel(1981)Nature 290,20.
Lawn(1981)Science 212,1159.
135. Shaw(1983)NAR 11,555.
136. Houghton(1980)NAR 8,1913.
Ohno(1981)PNAS 78,5305.
137. Higashi(1983)JBC 258,9522.
138. Derynck(1982)NAR 10,3605.
Gray(1982)Nature 298,859.
139. Devos(1983)NAR 11,4307.
Taniguchi(1983)Nature 302,305.
140. Mason(1983)Nature 303,300.
141. Swift(1982)PNAS 79,7263.
142. Steinert(1983)Nature 302,794.
143. Powell(1983)NAR 11,5327.
144. Nawa(1983)PNAS 80,90.
145. Hall(1982)NAR 10,3503.
146. Hoffman(1982)NAR 10,7819.
147. Brisson(1982)PNAS 79,4055.
148. van het Schip(1983)NAR 11,2529.
149. Cheung(1983)NAR 11,3703.
150. McLean(1983)JBC 258,8993.
151. Grez(1981)Cell 25,743.
Jung(1980)PNAS 77,5759.
152. Vlasak(1983)EJB 135,123.
153. Karin(1982)Nature 299,797.
154. Glanville(1981)Nature 292,267.
155. Daniel(1983)EMBO 2,1061.
156. Blanchetot(1983)Nature 301,732.
157. Nabeshima(1982)NAR 10,6099.
158. Ullrich(1983)Nature 303,821.
159. VanStraaten(1983)PNAS 80,3183.
160. Watt(1983)Nature 303,725;
PNAS 80,6307.
161. Capon(1983)Nature 304,507.
162. Takeya(1983)Cell 32,881.
163. Nakanishi(1981)EJB 115,429.
164. Whitfeld(1982)DNA 1,133.
Cochet(1982)Nature 297,335.
165. McReynolds(1978)Nature 273,723.
166. Catterall(1980)J.Cell Biol.87,480.
167. Kronenberg(1979)PNAS 76,4981.
Weaver(1982)Mol.Cell.Endocrin.
28,411.
168. Hendy(1981)PNAS 78,7365.
169. Vasicek(1983)PNAS 80,2127.
170. Sogawa(1983)JBC 258,5306.
171. Slightom(1983)PNAS 80,1897.
172. Michelson(1983)PNAS 80,472.
173. Barrera-Saldaña(1983)JBC 258,3787.
174. Pennica(1983)Nature 301,214.
175. Cooke(1980)JBC 255,6502.
Cooke(1982)Nature 297,603.
176. Sasavage(1982)JBC 257,678.
177. Gregory(1982)NAR 10,7581.
178. Lonberg(1983)PNAS 80,3661.
179. Hudson(1981)Nature 291,127.
180. Panthier(1982)Nature 298,90.
181. MacDonald(1982)JBC 257,14582.
182. Amaldi(1982)Gene 17,311.
183. Berry-Lowe(1982)JMAG 1,483.
184. Kandala(1983)NAR 11,3169.
185. Okamoto(1982)JBC 257,15192.
186. Haché(1983)JBC 258,4556.
187. Dugaiczuk(1982)PNAS 79,71.
188. Gubler(1983)PNAS 80,4311.
189. Hobart(1980)Nature 288,137.
190. Magazin(1982)PNAS 79,5152.
191. Minth(1982)JBC 257,10372.
192. Shen(1982)PNAS 79,4575.
193. Funckes(1983)JBC 258,8781.
194. Parker(1982)NAR 10,5121.
195. Parker, op.cit.
196. Hurst(1983)EMBO 2,769.
197. Edens(1982)Gene 18,1.
198. Gurr(1983)PNAS 80,2122.
199. MacDonald(1982)JBC 257,9724.
200. Lemischka(1982)Nature 300,330.
201. Valenzuela(1981)Nature 289,650.
202. Lee(1983)Cell 33,477.
203. Suske(1983)NAR 11,2257.
Chandra(1981)DNA 1,19.
204. Itoh(1983)Nature 304,547.
205. Land(1982)Nature 295,299.
206. Schmale(1983)EMBO 2,763.
207. Geiser(1983)JBC 258,9024.
208. Hung(1981)NAR 9,6407.
209. Hung(1983)JMB 164,481.
210. Hennighausen(1982)NAR 10,2677.
A. Sippel, *pc*.
211. Pedersen(1982)Cell 29,1015.
212. Marks(1982)JBC 257,9976.