# Codon usage in bacteria: correlation with gene expressivity

M.Gouy and C.Gautier

Equipe Evolution moléculaire, Laboratoire de Biométrie, Université Lyon I, 43 Bd du 11 Nov 1918, 69622 Villeurbanne, France

ABSTRACT

The nucleic acid sequence bank now contains over 600 protein coding genes of which 107 are from prokaryotic organisms. Codon frequencies in each new prokaryotic gene are given. Analysis of genetic code usage in the 83 sequenced genes of the Escherichia coli genome (chromosome, transposons and plasmids) is presented, taking into account new data on gene expressivity and regulation as well as iso-tRNA specificity and cellular concentration. The codon composition of each gene is summarized using two indexes: one is based on the differential usage of iso-tRNA species during gene translation, the other on choice between Cytosine and Uracile for third base. A strong relationship between codon composition and mRNA expressivity is confirmed, even for genes transcribed in the same operon. The influence of codon use on peptide elongation rate and protein yield is discussed. Finally, the evolutionary aspect of codon selection in mRNA sequences is studied.

## INTRODUCTION

Bacteria, and especially Escherichia coli, furnish the best documented example of the effect of natural selection on codon usage. The resulting bias in genetic code usage has two main components: correlation with tRNA availability in the cell and non random choices between pyrimidine ending codons. We present here two simple indexes to quantify these components and characterize all 107 bacterial sequences in the Lyon sequence bank ACNUC.

The E. coli sample (83 sequences) can now be considered as representative of the whole genome since it contains sequences from the chromosome as well as from plasmids and transposons. These sequences code for structural proteins (membrane and ribosomal proteins), for enzymes (amino acid biosynthesis operons) and for regulatory proteins (repressors, lexA and recA genes). Also, copious data on primary structure, codon recognition and relative abundance of E. coli tRNA species are available. Finally, regulating mechanisms and cellular concentrations of translation products are known for many genes in our sample, the latter, however, often on a qualitative basis. These data permit a more precise approach to the relationship between codon usage and translation efficiency and a somewhat new evolutionary point of view.

All genes studied here can be identified with Table 1. Sequence sources appear in the list of sequence references at the end of the paper. Table 2 shows the codon composition of 85 new genes not already presented in our last two compilations (1,2). As previously, codon compositions are expressed per 1000 and both initiator and stop codons are excluded from the data. E. coli ilvG gene (symbol ECOILVG) is given again because this nucleotide sequence has been redetermined and another reading frame has been proposed (3,4). Genes for elongation factor G (ECOEFG) and ribosomal protein S7 (ECOS7) are also re-included because more of their codons have been sequenced.

## CODON USAGE AND tRNA AVAILABILITY

The first relationship between codon occurrences in mRNAs and cellular tRNA abundancies described in the literature is known as the tRNA adaptation theory (5,6). This phenomenon has been observed in many specialized cells mainly involved in the translation of very few mRNA species. This is not the case of E. coli cell. However, the high use of codons corresponding to major tRNAs was remarked by Post et al. upon sequencing ribosomal protein (r-protein) genes (7,8). Since then, various workers have observed the same tendency in many other highly expressed genes, i.e. genes coding for abundant proteins (see the r-proteins, elongation factors and membrane proteins items in the list of sequence references). Ikemura quantified the relative frequencies of tRNA species in E. coli and presented an elegant demonstration of the correlation between codon use in messengers and tRNA cellular quantities (9,10). The linear correlation coefficient computed on several genes discriminates between high and low expressivity genes with exceptions that seem mainly due to mathe-matical artifacts (10). Following the work of one of us (11,12) we construct here a measure of the influence of tRNA availability based on a model of pro-tein synthesis dynamics

At the beginning of each polypeptide elongation cycle, a codon is found at the ribosomal A-site. The ternary complexes (aminoacyl-tRNAs bound with elongation factor Tu and GTP) diffusing in the cytoplasm interact with the codon and the ribosome at the A-site. Most often the codon does not belong to the tRNA recognition spectrum and therefore the aminoacyl-tRNA dissociates from the ribosome. When the specificity condition is fulfilled, the elongation cycle starts: transpeptidation and translocation occur. Thus, each codon can be characterized by the average number of codon-tRNA interactions at the A-site during one elongation cycle. The relative concentration of the codon-cognate tRNA is equivalent to its probability of colliding with the A-site

Table 1: mRNA sequence portfolio

| Symbol | Species and gene | No Codons |
|---|---|---|
| ECOLPP | E. coli outer membrane lipoprotein | 77 |
| ECOOMPA | E. coli ompA gene | 345 |
| ECOS2 | E. coli ribosomal protein S2 | 240 |
| ECOS4 | E. coli ribos. prot. S4 3' partial seq. | 60 |
| ECOS7 | E. coli ribos. prot. S7 5' & 3' partial seq. | 112 |
| ECOS10 | E. coli ribosomal protein S10 | 102 |
| ECOS11 | E. coli ribos. prot. S11 5' partial seq. | 55 |
| ECOS12 | E. coli ribosomal protein S12 | 123 |
| ECOS13 | E. coli ribos. prot. S13 partial sequence | 35 |
| ECOS17 | E. coli ribos. prot. S17 3' partial seq. | 24 |
| ECOS20 | E. coli ribosomal protein S20 UUG initiator | 86 |
| ECOL1 | E. coli ribosomal protein L1 | 233 |
| ECOL3 | E. coli ribos. prot. L3 5' partial seq. | 79 |
| ECO712 | E. coli ribosomal proteins L7/L12 | 120 |
| ECOL10 | E. coli ribosomal protein L10 | 164 |
| ECOL11 | E. coli ribosomal protein L11 | 141 |
| ECOL14 | E. coli ribos. prot. L14 5' partial seq. | 19 |
| ECOTUFA | E. coli elong. factor Tu (tufA) GUG initiator | 393 |
| ECOTUFB | E. coli elong. factor Tu (tufB) | 393 |
| ECOEFG | E. coli elongation factor G 5' & 3' partial seq. | 114 |
| ECOEFTS | E. coli elongation factor Ts (tsf) | 282 |
| ECORECA | E. coli recA | 352 |
| ECORPOA | E. coli RNA polymerase α-subunit 5' partial seq. | 158 |
| ECORPOB | E. coli RNA polymerase β-subunit | 1341 |
| ECORPOC | E. coli RNA polym. β'-subunit 5' partial seq. | 510 |
| ECORPOD | E. coli RNA polymerase σ-subunit | 612 |
| ECOALRS | E. coli alanyl-tRNA-synthetase | 875 |
| ECOUNC1 | E. coli hypoth. gene uncl GUG initiator | 129 |
| ECOUNC2 | E. coli a-subunit of ATP-synthase (unc2) | 270 |
| ECOUNC3 | E. coli c-subunit of ATP-synthase (unc3) | 78 |
| ECOUNC4 | E. coli b-subunit of ATP-synthase (unc4) | 155 |
| ECOUNC5 | E. coli δ-subunit of ATP-synthase (unc5) | 176 |
| ECOUNCA | E. coli α-subunit of ATP-synthase (uncA) | 512 |
| ECOUNCG | E. coli γ-subunit of ATP-synthase (uncG) | 286 |
| ECOUNCD | E. coli β-subunit of ATP-synthase (uncD) | 459 |
| ECOUNCC | E. coli ε-subunit of ATP-synthase (uncC) | 132 |
| ECOLACI | E. coli lacI lac operon repressor | 359 |
| ECOLACZ | E. coli lacZ β-galactosidase 3' partial seq. | 24 |
| ECOLACY | E. coli lacY lactose permease | 416 |
| ECOLACA | E. coli lacA transacetylase 5' partial seq. | 25 |
| ECOTRPL | E. coli trpL attenuation of trp operon | 13 |
| ECOTRPE | E. coli trpE anthranilate synthetase component I | 519 |
| ECOTRPD | E. coli trpD | 530 |
| ECOTRPC | E. coli trpC | 451 |
| ECOTRPB | E. coli trpB | 397 |
| ECOTRPA | E. coli trpA | 267 |
| ECOTRPR | E. coli trpR trp & aroH operons repressor | 107 |
| ECOILVL | E. coli ilvL attenuation of ilv operon | 31 |
| ECOILVG | E. coli ilvG (mutant Valine resistant) | 520 |
| ECOILVE | E. coli ilvE 5' partial sequence | 80 |
| ECOTHRL | E. coli thrL attenuation of thr operon | 20 |
| ECOTHRA | E. coli thrA | 819 |
| ECOPHEL | E. coli pheL attenuation of phe operon | 14 |

| Symbol | Species and gene | No Codons |
|--------|------------------|-----------|
| ECOHISL | E. coli hisL attenuation of his operon | 15 |
| ECOAMPC | E. coli ampC cephalosporinidase | 376 |
| ECOARAC | E. coli araC ara operon regulation | 291 |
| ECOAROH | E. coli aroH 5' & 3' partial sequence | 152 |
| ECOASNA | E. coli asnA asparagine synthetase | 329 |
| ECOFOL | E. coli dihydrofolate reductase | 158 |
| ECOLEXA | E. coli lexA | 201 |
| ECOLTA | E. coli heat-labile toxin, A-subunit partial seq. | 124 |
| ECOLTB | E. coli heat-labile toxin, B-subunit | 123 |
| ECONDH | E. coli ndh NADH dehydrogenase | 433 |
| ECOPHOA | E. coli phoA alkaline phosphatase | 76 |
| ECOR1EN | E. coli restriction endonuclease R1 | 276 |
| ECOR1ME | E. coli restriction methylase R1 | 325 |
| ECOTNAA | E. coli tnaA tryptophanase | 470 |
| ECOUVRB | E. coli uvrB | 28 |
| ECOTN3L | E. coli transposon TN3 β-lactamase | 285 |
| ECOTN3T | E. coli transposon TN3 tnpA transposase | 1014 |
| ECOTN3R | E. coli transposon TN3 repressor | 184 |
| ECOTN9 | E. coli transposon TN9 chloramphenicol resistance | 218 |
| ECO903K | E. coli transposon TN903 kanamycin resistance | 270 |
| ECPFOL | plasmid R388 dihydrofolate reductase | 77 |
| CLOIMMU | plasmid CLO-DF13 immunity gene | 84 |
| CLOH | plasmid CLO-DF13 gene H | 48 |
| R1PORI1 | plasmid R1 hypoth. gene ori1 10500 d protein | 85 |
| R1PORI2 | plasmid R1 hypoth. gene ori2 7000 d protein | 60 |
| SM1RA1 | plasmid SM1 repA1 gene GUG initiator | 284 |
| SM1RA2 | plasmid SM1 repA2 gene | 102 |
| SAULRER | plasmid PE194 erythromycin resist. regulator | 18 |
| SAURERY | plasmid PE194 erythromycin resistance | 243 |
| SAUBLAC | Staphylococcus aureus β-lactamase | 46 |
| STYTRPL | Salmonella typhimurium trpL attenuation trp operon | 13 |
| STYTRPE | Salmonella typhimurium trpE | 519 |
| STYTRPD | Salmonella typhimurium trpD 5' partial seq. | 193 |
| STYTRPB | Salmonella typhimurium trpB | 396 |
| STYTRPA | Salmonella typhimurium trpA | 267 |
| STYILVL | Salmonella typhimurium ilvL attenuation ilv operon | 31 |
| STYHISL | Salmonella typhimurium hisL attenuation his operon | 15 |
| STYHISJ | Salmonella typhimurium hisJ | 259 |
| STYLEUL | Salmonella typhimurium leuL attenuation leu operon | 27 |
| STYARGT | Salmonella typhimurium argT | 259 |
| SALHIN | Salmonella hin protein | 189 |
| SMALPP | Serratia marcescens outer membrane lipoprotein | 76 |
| SMATRPE | Serratia marcescens trpE 3' partial sequence | 26 |
| SMATRPG | Serratia marcescens trpG | 192 |
| SDYTRPE | Shigella dysenteriae trpE 3' partial sequence | 26 |
| SDYTRPD | Shigella dysenteriae trpD 5' partial sequence | 193 |
| EAMLPP | Erwinia amylovora outer membrane lipoprotein | 77 |
| KAETRPB | Klebsiella aerogenes trpB 3' partial sequence | 19 |
| KAETRPA | Klebsiella aerogenes trpA | 268 |
| KPNNIFH | Klebsiella pneumoniae nifH nitrogenase 2 | 292 |
| ANANIFH | Anabaena nifH nitrogenase reductase | 298 |
| RHINIFH | Rhisobium melitoti nifH nitrogenase | 296 |
| BLIPEN | Bacillus licheniformis penicillinase | 306 |
| HALRHO | Halobacterium halobium bacteriorhodopsin | 261 |

## SEQUENCE REFERENCES (Table 1 order)

ECOLPP    Nakamura, K., Inouye, M.: Cell, 18, 1109-1117(1979)

ECOOMPA    Beck, E., Bremer, E.: Nucleic Acids Res., 8, 3011-3027(1980)

Movva, N.R., Nakamura, K., Inouye, M.: Proc. Natl. Acad. Sci. USA, 77, 3845-3949(1980)

ECOS2    An, G., Bendiak, D.S., Mamelak, L.A., Friesen, J.D.: Nucleic Acid Res., 9,4162-4172(1981)

ECOS4    Post, L.E., Nomura, M.: J. Biol. Chem., 254, 10604-10606(1979)

ECOS7    Post, L.E., Nomura, M.: J. Biol. Chem. 255, 4660-4666(1980)

ECOS10    Olins, P.O., Nomura, M.: Cell, 26, 205-211(1981)

ECOS11    Post, L.E., Arfsten, A.E., Davis, G.R., Nomura, M.: J. Biol. Chem., 255, 4653-4659(1980)

ECOS12    same reference as ECOS7

ECOS13    same reference as ECOS11

ECOS17    Post, L.E., Arfsten, A.E., Reusser, F., Nomura, M.: Cell, 15, 215-229(1978)

ECOS20    Mackie, G.A.:J. Biol. Chem., 256, 8177-8182(1981)

ECOL1    Post, L.E., Strycharz, G.D., Nomura, M., Lewis, H., Dennis, P.P.: Proc. Natl. Acad. Sci. USA, 76, 1697-1701(1979)

ECOL3    same reference as ECOS10

ECO712    same reference as ECOL1

ECOL10    same reference as ECOL1

ECOL11    same reference as ECOL1

ECOL14    same reference as ECOS17

ECOTUFA    Yokota, T., Sugisaki, H., Takanami, M., Kaziro, Y.: Gene, 12, 25-31(1980)

ECOTUFB    An, G., Friesen, J.D.: Gene, 12, 33-39(1980)

ECOEFG    Post, L.E., Nomura, M.:J. Biol. Chem., 255, 4660-4666(1980)

Yokota, T., Sugisaki, H., Takanami, M., Kaziro, Y.: Gene, 12, 25-31(1980)

ECOEFTS    same reference as ECOS2

ECORECA    Sancar, A., Stachelek, C., Konigsberg, W., Rupp, W.D.:Proc. Natl. Acad. Sci. USA,77,2611-2615(1980)

ECORPOA    same reference as ECOS4

ECORPOB    Ovchinnikov, Y.A., Monastyrskaya, G.S., Gubanov, V.V., Guryev, S.O., Chertov, O.Y., Modyanov, N.N., Grinkevich, V.A., Makarova, I.A., Marchenko,T.V., Polovnikova, I.N., Lipkin, V.M., Sverdlov,E.D.: Eur. J. Biochem., 116, 621-629(1981)

ECORPOC    same reference as ECORPOA

Squires, C., Krainer, A., Barry, G., Shen, W.F., Squires, C.L.: Nucleic Acids Res., 9, 6827-6840 (1981)

ECORPOD    Burton, Z., Burgess, R.R., Lin, J., Moore, D., Holder, S., Gross, C.A.: Nucleic Acids Res.,9, 2889-2903(1981)

ECOALRS    Putney, S.D., Royal, N.J., De Vegvar, H.N., Herlihy, W.C., Biemann, K., Schimmel, P.:Science, 213, 1497-1500(1981)

ECOUNC1   Gay, N.J., Walker, J.E.: Nucleic Acids Res., 9,
        3919-3926(1981)
ECOUNC2   same reference as ECOUNC1
ECOUNC3   same reference as ECOUNC1
ECOUNC4   same reference as ECOUNC1
ECOUNC5   same reference as ECOUNC1
ECOUNCA   Gay, N.J., Walker, J.E.: Nucleic Acids Res., 9,
        2187-2194(1981)
ECOUNCG   Saraste, M., Gay, N.J., Eberle, A., Runswick, M.J.,
        Walker, J.E.: Nucleic Acids Res., 9, 5287-5296(1981)
ECOUNCD   same reference as ECOUNCG
ECOUNCC   same reference as ECOUNCG
ECOLACI   Farabaugh, P.: Nature,274,765-769(1978)
ECOLACZ   Buchel, D.E., Gronenborn, B., Muller-Hill, B.:Nature,
        283,541-545(1980)
        Maizels, N.:Proc. Natl. Acad. Sci. USA, 70,
        3585-3592(1973)
ECOLACY   Buchel, D.E., Gronenborn, B., Muller-Hill, B.:
        Nature, 283, 541-545(1980)
ECOLACA   same reference as ECOLACY
ECOTRPL   Platt, T.: Cell, 24, 10-23(1981)
ECOTRPE   Nichols, B.P., Van Cleemput, M., Yanofsky, C.: J.
        Mol. Biol., 146,45-54(1981)
ECOTRPD   Nichols, B.P., Miozzari, G.F., Van Cleemput, M.,
        Bennett, G.N., Yanofsky, C.: J. Mol. Biol., 142,
        503-517(1980)
        Horowitz, H., Christie, G.E., Platt, T.: J. Mol.
        Biol., 156, 245-256(1982)
ECOTRPC   Christie, G.E., Platt, T.: J. Mol. Biol., 142,519-530,
        (1980)
ECOTRPB   Crawford, I.P., Nichols, B.P., Yanofsky, C.: J.
        Mol. Biol., 142, 489-502(1980)
ECOTRPA   Nichols, B.P., Yanofsky, C.: Proc. Natl. Acad. Sci.
        USA, 76, 5244-5248(1979)
ECOTRPR   Gunsalus, R.P. Yanofsky, C.: Proc. Natl. Acad. USA,
        77, 7117-7121(1980)
ECOILVL   Nargang, F.E., Subrahmanyam, C.S., Umbarger, H.E.:
        Proc. Natl. Acad. Sci. USA,77,1823-1827(1980)
        Lawther, R.P., Calhoun, D.H., Adams, C.W.,Hauser,C.A.,
        Gray, J., Hatfield, G.W.: Proc. Natl. Acad. Sci.
        USA, 78, 922-925(1981)
ECOILVG   Lawther, R.P., Calhoun, D.H., Adams, C.W., Hauser,
        C.A., Gray, J., Hatfield, G.W.: Proc. Natl. Acad.
        Sci. USA, 78, 922-925(1981)
ECOILVE   Lawther, R.P., Nichols, B., Zurawski, G., Hatfield,
        G.W.: Nucleic Acids Res., 7, 2289-2301(1979)
ECOTHRL   Gardner, J.F.: Proc. Natl. Acad. Sci. USA, 76,
        1706-1710(1979)
ECOTHRA   Katinka, M., Cossart, P., Sibilli, L., Saint-Girons,
        I., Chalvignac, M.A., Le Bras, G., Cohen, G.N.,
        Yaniv, M.:Proc. Natl. Acad. Sci. USA,77,
        5730-5733(1980)
ECOPHEL   Zurawski, G., Brown, K., Killingly, D., Yanofsky,
        C.: Proc. Natl. Acad. Sci. USA,75,4271-4275(1978)

ECOHISL    Di Nocera, P.P., Blasi, F., Di Lauro, R., Frunzio, R.
           Bruni, C.B.: Proc. Nat. Acad. Sci. USA, 75,
           4276-4280(1878)
ECOAMPC    Jaurin, B., Grundstrom, T.: Proc. Natl. Acad. Sci.
           USA, 78, 4897-4901(1981)
ECOARAC    Miyada, C.G., Horwitz, A.H., Cass, L.G., Timko, J.,
           Wilcox, G.: Nucleic Acids Res.,8,5267-5274(1980)
           Wallace, R.G., Lee, N., Fowler, A.V.: Gene, 12,
           179-190(1980)
           Stoner, C.M., Schleif, R.:J. Mol. Biol., 154,
           649-652(1982)
ECOAROH    Zurawski, G., Gunsalus, R.P., Brown, K.D., Yanofsky,
           C.:J. Mol. Biol.,145,47-73(1981)
ECOASNA    Nakamura, M., Yamada, M., Hirota, Y., Sugimoto, K.,
           Oka, A., Takanami, M.:Nucleic Acids Res., 9,
           4669-4676(1981)
ECOFOL     Smith, D.R., Calvo, J.M.: Nucleic Acids Res., 8,
           2255-2274(1980)
ECOLEXA    Miki, T., Ebina, Y., Kishi, F., Nakazawa, A.:
           Nucleic Acids Res., 9,529-543(1981)
           Horii, T., Ogawa, T., Ogawa, H.: Cell,23,689-697(1981)
           Markham, B.E., Little, J.W., Mount, D.W.: Nucleic
           Acids Res., 9, 4149-4161(1981)
ECOLTA     Spicer, E.K., Kavanaugh, W.M., Dallas, W.S., Falkow,
           S., Konigsberg, W.H., Schafer, D.E.: Proc. Natl.
           Acad. Sci. USA, 78, 50-54(1981)
ECOLTB     Dallas, W.S., Falkow, S.: Nature, 288,499-501(1980)
ECONDH     Young, I.G., Rogers, B.L., Campbell, H.D.,
           Jaworowski, A., Shaw, D.C.: Eur. J. Biochem., 116,
           165-170(1981)
ECOPHOA    Kikuchi, Y., Yoda, K., Yamasaki, M., Tamura, G.:
           Nucleic Acids Res., 9, 5671-5678(1981)
ECOR1EN    Greene, P.J., Gupta, M., Boyer, H.W., Brown, W.E.,
           Rosenberg, J.M.: J. Biol. Chem. , 256, 2143-2153
           (1981)
           Newman, A.K., Rubin, R.A., Kim, S.H., Modrich, P.:
           J. Biol. Chem., 256, 2131-2139(1981)
ECOR1ME    same reference as ECOR1EN
ECOTNAA    Deeley, M.C., Yanofsky, C.:J. Bact., 147,787-796(1981)
ECOUVRB    Van Den Berg, E., Zwetsloot, J., Noordermeer, I.,
           Pannekoek, H., Dekker, B., Dijkema, R., Van Ormondt,
           H.: Nucleic Acids Res.,9,5623-5643(1981)
ECOTN3L    Heffron, F., Mc Carthy, B.J., Ohtsubo, H., Ohtsubo,
           E.: Cell, 18, 1153-1163(1979)
           Sutcliffe, G.: Proc. Natl. Acad. Sci. USA, 75,
           3737-3741(1978)
ECOTN3R    Heffron, F., Mc Carthy, B.J., Ohtsubo, H., Ohtsubo,
           E.:Cell,18, 1153-1163(1979)
           Chou, J., Lemaux, P.G., Casabadan, M.J., Cohen, S.N.:
           Nature, 282, 801-806(1979)
ECOTN3T    same reference as ECOTN3R
ECOTN9     Alton, N.K., Vapnek, D.: Nature, 282, 864-869(1979)
           Marcoli, R., Iida, S., Bickle, T.A.: FEBS Lett., 110,
           11-14(1980)

ECO903K    Oka, A., Sugisaki, H., Takanami, M.:J. Mol. Biol.,
     147, 217-226(1981)

     Heidekamp, F., Baas, P.D., Van Boom, J.H., Veeneman,
     G.H., Zipursky, S.L., Jansz, H.S.: Nucleic Acids
     Res., 9, 3335-3354(1981)

ECPFOL    Zolg, J.W., Hanggi, U.J.: Nucleic Acids Res.,9,
     698-709(1981)

     Swift, G., Maccarthy, B.J., Heffron, F.: Mol. Gen.
     Genet., 181,441-447(1981)

CLOIMMU    Elzen, P.J.M., Gaastra, W., Spelt, C.E., De Graaf,
     F.K., Veltkamp,E., Nijkamp, H.J.J.: Nucleic Acids
     Res.,8,4349-4363(1980)

CLOH    Stuitje, A.R., Spelt, C.E., Veltkamp, E., Nijkamp,
     H.J.J.:Nature,290,264-267(1981)

R1PORI1    Stougaard, P., Molin, S., Nordstrom, K., Hansen, F.G.:
     Mol. Gen. Genet.,181,116-122(1981)

R1PORI2    same reference as R1PORI1

     Stougaard, P., Molin, S., Nordstrom, K.: Proc. Natl.
     Acad. Sci. USA,78,6008-6012(1981)

SM1RA1    Rosen,J., Ryder, T., Inokuchi, H., Ohtsubo, H.,
     Ohtsubo,E.: Molec. Gen. Genet.,179,527-537(1980)

     Rosen, J., Ryder, T., Ohtsubo,H., Ohtsubo, E.:Nature,
     290,794-797(1981)

SM1RA2    same reference as SM1RA1

SAULRER    Horinouchi, S., Weisblum, B.: Proc. Natl. Acad. Sci.
     USA,77,7079-7083(1980)

SAURERY    same reference as SAULRER

SAUBLAC    Mc Laughlin, J.R., Murray, C.L., Rabinowitz, J.C.: J.
     Biol. Chem.,256,11283-11291(1981)

STYTRPL    Yanofsky, C., Van Cleemput, M.: J. Mol. Biol.,154,
     235-246(1982)

STYTRPE    same reference as STYTRPL

STYTRPD    Nichols, B.P., Miozzari, G.F., Van Cleemput, M.,
     Bennett, G.N., Yanofsky, C.: J. Mol. Biol., 142,
     503-517(1980)

STYTRPB    same reference as ECOTRPB

STYTRPA    same reference as ECOTRPA

STYILVL    Taillon, M.P., Gotto, D.A., Lawther, R.P.:Nucleic
     Acids Res.,9,3419-3432(1981)

STYHISL    Barnes, W.M.:Proc. Natl. Acad. Sci. USA,75,4281-4285
     (1978)

STYHISJ    Higgins, C.F., Ames, G.F.L.: Proc. Natl. Acad. Sci.
     USA,78,6038-6042(1981)

STYLEUL    Gemmill, R.M., Wessler, S.R., Keller, E.B., Calvo,
     J.M.: Proc. Natl. Acad. Sci. USA,76,4941-4945(1979)

STYARGT    same reference as STYHISJ

SALHIN    Zieg, J., Simon, M.: Proc. Natl. Acad. Sci. USA,77,
     4196-4200(1980)

SMALPP    Nakamura, K., Inouye, M.: Proc. Natl. Acad. Sci. USA,
     77,1369-1373(1980)

SMATRPE    same reference as STYTRPD

SMATRPG    same reference as STYTRPD

SDYTRPE    same reference as STYTRPD

SDYTRPD    same reference as STYTRPD

EAMLPP    Yamagata, H., Nakamura, K., Inouye, M.: J. Biol.
     Chem.,256,2194-2198(1981)

KAETRPB    Nichols, B.P., Blumenberg, M., Yanofsky, C.: Nucleic
           Acids Res., 9,1743-1755(1981)
KAETRPA    same reference as KAETRPB
KPNNIFH    Sundaresan, V., Ausubel, F.M.: J. Biol. Chem., 256,
           2808-2812(1981)
ANANIFH    Mevarech, M., Rice, D., Haselkorn, R.: Proc. Natl.
           Acad. Sci. USA,77,6476-6480(1980)
RHINIFH    Torok, I., Kondorosi, A.: Nucleic Acids Res.,9,
           5711-5723(1981)
BLIPEN     Neugebauer, K., Sprengel, R., Schaller, H.: Nucleic
           Acids Res.,9,2577-2588(1981)
HALRHO     Dunn, R., Mc Coy, J., Simsek, M., Majumdar, A.,
           Chang, S.H., Rajbhandary, U.L., Khorana, H.G.:
           Proc. Natl. Acad. Sci. USA,78,6744-6748(1981)

codon. Hence, if this probability is f, the mean number of codon-tRNA interactions necessary for the elongation cycle to occur is 1/f. The rarest tRNA species in E. coli is $tRNA_2^{Ile}$ decoding codon AUA (0.3% of the total tRNA population, (9) and legend to fig. 1); the most abundant species is $tRNA_3^{Gly}$ decoding GGPy codons (6.5% of the tRNA population). Consequently, 1/0.003 = 333 non specific codon-tRNA interactions are necessary, on the average, to translate the AUA codon, whereas an average of 1/0.065 = 15 such interactions occur when translating GGU or GGC codons. For a given gene, we weight the mean number of non specific codon-tRNA interactions for each codon by the relative frequency of the codon in the sequence. The value thus obtained is the average number of tRNA discriminations per elongation cycle (P1 index). This index is scaled on the abcissa of fig. 1.

All highly expressed genes appear to the left of the figure. These are genes coding for major membrane proteins, for the 3 elongation factors Tu, Ts and G, for r-proteins, and the recA gene, which is not constitutively highly expressed as are the preceding messengers. These genes are therefore highly optimized for a small number of tRNA discriminations. We calculate that the best theoretical messenger for the L1 r-protein (ECOL1) would require 22.3 mean discriminations instead of the observed value of 25.5. Translation of a codon whose cognate tRNA is rare in the cell requires more non specific codon-tRNA interactions at the A-site than do codons decoded by major tRNAs. Consequently, translation of the latter codons might be faster and/or less error-prone than translation of other synonymous codons. A quantitative assessment of this phenomenon is not yet possible since the duration of a non specific codon-tRNA interaction is unknown (11). Therefore, given the available data, P1 is the best possible index for quantifying the effect of codon use on translation rate.

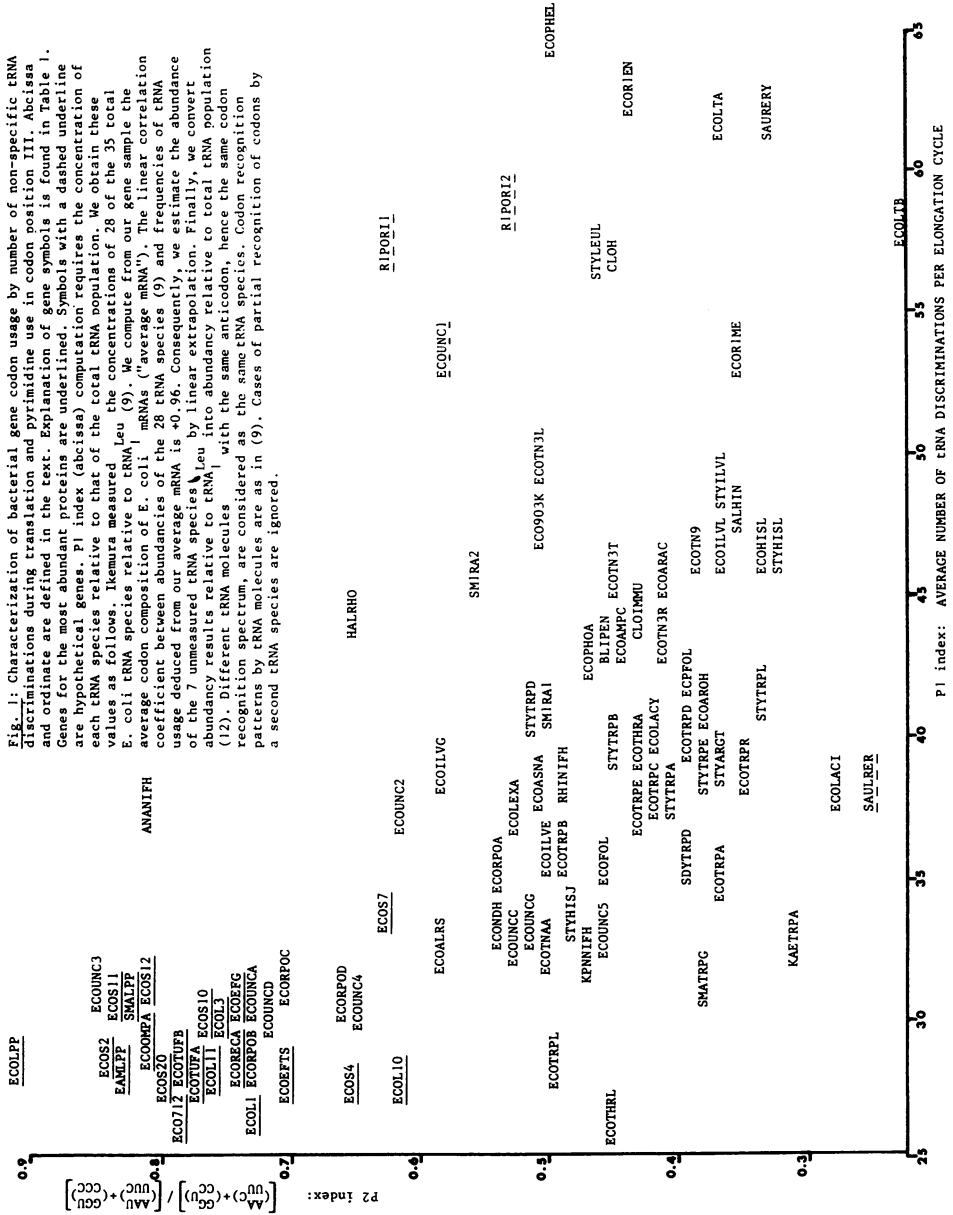| | | ECO S2 | ECO S4 | ECO S7 | ECO S10 | ECO S17 | ECO S20 | ECO L3 | ECO L14 | ECO TUFA | ECO TUFB | ECO EFG | ECO EFTS | ECO RPOA | ECO RPOB | ECO RPOC | ECO RPOD | ECO ALRS | ECO UNC1 | ECO UNC2 | ECO UNC3 | ECO UNC4 | ECO UNC5 | ECO UNCA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Fig. 1: Characterization of bacterial gene codon usage by number of non-specific tRNA discriminations during translation and pyrimidine use in codon position III. Abcissa and ordinate are defined in the text. Explanation of gene symbols is found in Table I. Genes for the most abundant proteins are underlined. Symbols with a dashed underline are hypothetical genes. PI index (abcissa) computation requires the concentration of each tRNA species relative to that of the total tRNA population. We obtain these values as follows. Ikemura measured the concentrations of 28 of the 35 total E. coli tRNA species relative to $tRNA_{Leu}^1$ (9). We compute from our gene sample the average codon composition of E. coli mRNAs ("average mRNA"). The linear correlation coefficient between abundancies of the 28 tRNA species (9) and frequencies of tRNA usage deduced from our average mRNA is +0.96. Consequently, we estimate the abundance of the 7 unmeasured tRNA species relative to $tRNA_{Leu}^1$ by linear extrapolation. Finally, we convert abundancy results relative to $tRNA_{Leu}^1$ into abundancy relative to total tRNA population (12). Different tRNA molecules with the same anticodon, hence the same codon recognition spectrum, are considered as the same tRNA species. Codon recognition patterns by tRNA molecules are as in (9). Cases of partial recognition of codons by a second tRNA species are ignored.

PI index: AVERAGE NUMBER OF tRNA DISCRIMINATIONS PER ELONGATION CYCLE

The results in fig. I suggest that translation of highly expressed mRNAs is faster than translation of others. The complex relationship between translation rate and mRNA protein yield has been analyzed by Bergmann and Lodish (13). Briefly, these authors show that the main factor determining protein yield is initiation rate, but that elongation and initiation rates can interact and modulate protein yield.

Plasmid and transposon genes generally have a PI index value greater than 43, indicating that the cellular tRNA population is poorly adapted to their codon usage. Conversely, few chromosomal genes are located in the same region of fig. I. E. coli mRNAs falling there are very weakly expressed: toxin genes (ECOLTA and ECOLTB), the Eco RI restriction enzyme system (ECORIEN endonuclease and ECORIME methylase), two leader genes of amino acid biosynthesis operons (ECOPHEL and STYLEUL) and an hypothetical regulatory gene in unc operon (ECOUNCI).

Genes from enterobacteriaceae other than E. coli (Salmonella, Serratia, Shigella, Klebsiella and Erwinia (14)) closely follow E. coli codon usage as seen with genes sequenced in both species, namely lpp and trp operon genes. Genes from non enteric bacteria (Anabaena, Bacillus, Rhizobium, Staphylococcus and the archaebacterium Halobacterium) have been included for completeness, but the PI index, based on the E. coli tRNA population, may be meaningless for them.

## CHOICE BETWEEN PYRIMIDINES IN CODON POSITION III

Grosjean et al. first noted in the MS2 phage genome a bias in the choice between C and U bases in codon position III (15). They found that nucleotides in the degenerate position consistently yield a codon-anticodon binding energy of intermediate strength. For example, AAC appears more often than AAU, and GCU more often than GCC. This bias is independent of amino acid composition (whatever the bases $X_1$ and $X_2$, codons $X_1X_2C$ and $X_1X_2U$ code for the same amino acid) and of tRNA availability ($X_1X_2C$ and $X_1X_2U$ are mainly decoded by the same tRNA (9)). Our group showed by statistical analysis of 29 bacterial genes that this bias is not present in all sequences, as already recognized by Fiers and Grosjean in the lac I gene (16), but strongly depends on the expressivity level of the gene (2). Ikemura later stated the same conclusion (9,10). If the first two bases of a codon are both A or U, C in third position will give a codon-anticodon binding energy nearer the mean than would U. Likewise, if the first two bases are both C or G, the "right choice" for third base is U, for C would give a strong binding energy.

Consequently, we characterize each messenger by the frequency of "right

choices" between the pyrimidines among codons beginning with AA, AU, UA, UU, CC, CG, GC or GG. This frequency, called P2 index, is scaled on the ordinate of fig. 1. As with P1 index, a discrimination between messengers for abundant and rare proteins occurs. The highly expressed genes appear in the upper left part of the figure.

In order to summarize the optimal choice between synonymous codons, we compute the mean frequency of each codon after pooling all highly expressed genes, that is genes with underlined symbols in fig. 1. The results, HIGH. EXPR. column of Table 2, clearly show the two tendencies in codon usage discussed (the codon recognition patterns of E. coli tRNA species are given in (9)).

The interpretation of pyrimidine selection given by Grosjean et al. involves the energy level of the codon-anticodon interaction (15) and is connected to translational fidelity (29). However, the nature of the effect of U/C choice on mRNA translation remains unclear. This effect cannot be restricted to fidelity because we show here that pyrimidine choice is related to gene expressivity.

AN INTRA-OPERON RELATIONSHIP BETWEEN GENE EXPRESSIVITY AND CODON USAGE

To examine more precisely the relationships between codon usage and gene expressivity, we reviewed the available quantitative data (Table 3). Although these data are not numerous, a general agreement with fig. 1 results: clearly, all genes coding for the most abundant proteins are clustered to the top left of fig. 1. The nine unc (or atp) operon genes and the four RNA-polymerase subunit genes of E. coli merit more comment, however.

The completely sequenced unc operon comprises 8 structural genes coding for the 8 polypeptides that form the ATP-synthase complex and a hypothetical control gene unc1 (17-19). All these genes are controlled by the same promoter. Products of genes uncA, uncD, unc5, uncG and uncC form the $F_1$-part of the ATP-synthase complex with the stoichiometry 3:3:1:1:1 (18-21). The unc2, unc3 and unc4 gene products, respectively the a, c and b polypeptides, form together the $F_0$-part of the same complex (17). The stoichiometry of the $F_0$-part is un-clear, but the c peptide appears in more than 5 copies and the a and b peptides in 1 or 2 copies in the complex (21). The positions in fig. 1 of these 8 struc-tural genes correlate with the gene product copy numbers in the ATP-synthase complex. Genes with high copy numbers (unc3, uncA and uncD) have a higher P2 index value and generally a lower P1 value than genes with low copy numbers (unc5, uncG and uncC). As stressed by Fillingame, the coordinate expression of the 8 unc operon structural genes probably involves unidentified post-transcrip-tional mechanisms (21).

Three of the four RNA-polymerase genes (rpoB, C and D) are located in
fig. 1 with the highly expressed genes while rpoA is not. Table 3 shows that
these genes correspond to proteins of medium abundance and that rpoA is more
highly expressed than the other three genes because its product, the $\alpha$-subunit
appears in 2 copies in the holoenzyme (22). Note however that this is  cons-
titutive expressivity level; it is therefore difficult to compare with the lac
and trp genes whose levels are for maximum derepression. The rpoA gene is co-
transcribed with other r-protein genes in this order: rpsM (S13), rpsK (S11),
rpsD (S4), rpoA and rplQ (L17) (23). No regulatory feature in the operon pri-
mary structure is apparent between the first 3 r-protein genes and rpoA,
although the expressivities of these genes are markedly different (24). The
locations in fig. 1 of rpsK (ECOS11),rpsD (ECOS4) and rpoA genes show that the
same correlation between codon usage and gene expressivity holds for this ope-
ron as for the unc operon. However, only part of these last 3 genes has been
sequenced, hence we cannot be sure of their overall coding strategy. On the
other hand, rpoD is part of a single-gene operon (25) while rpoB and C are co-
transcribed with rplJ (L10) and rplL (L7/L12, symbol ECO712) but a transcrip-
tion attenuator site, which has been shown to function actively in vivo (26),
is present in the nucleotide sequence between rplL and rpoB, C genes (27).

Table 3: Cellular contents of various E. coli gene products.

| Symbol, gene and product | No. molecules per genome | Ref. |
|---|---|---|
| ECOLPP  lpp (outer membrane lipoprotein) | 330,000 | 31[§] |
| ECOTUF- tufA,B (elongation factor Tu) | 89,000[&] | 11[§], 32 |
| ECORECA recA | 38,000[&] | 33 |
| ECOOMPA ompA (outer membrane protein) | 36,800 | 34 |
| ECO712  rplL (L7 and L12 r-proteins) | 25,000 | 34 |
| other r-proteins | 9,200 | 11[§] |
| ECOEFTS tsf (elongation factor Ts) | 9,200 | 11[§] |
| ECOEFG  fusA (elongation factor G) | 9,200 | 11[§] |
| ECORPOA rpoA ($\alpha$-subunit of RNA-polymerase) | 4,000[&] | 34 |
| ECOLACY lacY (lactose permease) | 3,300[&] | 35 |
| ECOTRP- trpA-E (tryptophan biosynthesis) | 3,100[&] | 36 |
| ECORPOB rpoB ($\beta$-subunit of RNA-polymerase) | 1,400 | 34 |
| ECORPOC rpoC ($\beta'$-subunit of RNA-polymerase) | 1,400 | 34,37 |
| ECOALRS (aminoacyl-tRNA-synthetases) | 500-1,300 | 34 |
| ECOARAC araC (ara operon regulation) | 100 | 38 |
| ECOTRPR trpR (trp and aroH operons repressor) | 30 | 39 |
| ECOLACI lacI (lac operon repressor) | 10 | 40 |

Data have been computed for a cell growth rate of 1.5 doublings per hour
when possible.
[&] variable with growth conditions, value for maximum derepression.
[§] see references therein.

In addition to the RNA-polymerase genes and unc operon, other cases exist where several genes in one operon are sequenced. Three operons have been completely sequenced in E. coli. 1) trp operon: trpL, trpE, trpD, trpC, trpB and trpA genes (also partly sequenced in Salmonella typhimurium, Serratia marcescens, Shigella dysenteriae and Klebsiella aerogenes); 2) rpl operon: rplK (L11) and rplA (L1) genes; 3) elongation factor Ts operon: rpsB (S2) and tsf (EF-Ts) genes. Examples of partly sequenced operons are rpsL (S2), rpsG (S7), fusA (EF-G) and tufA (EF-Tu) genes; rpsJ (S10) and rplC (L3) genes and finally genes for the Eco RI restriction and modification enzymes (symbols ECORIEN and ECORIME). In all these cases, both the expressivity level and codon composition of genes from each operon (fig. 1) are similar. We cannot interpret Fig. 1 location of ilv operon genes (ilvL, G and E) since data on their expressivity are lacking. Moreover, ilvG is a peculiar cryptic gene (4).

The leader genes for amino acid biosynthesis operons (hisL, ilvL, pheL, thrL and trpL in E. coli and S. typhimurium) are exceptions to this observation, being widely scattered throughout fig. 1. We believe they are not adequately represented by this method because peculiar constraints are exerted on their structure (28).

This analysis of genes in the same operon argues for a relationship between P2 index (U/C choice in degenerate position) and the relative level of gene expressivity. Thus, codon usage may be involved in expressivity tuning of genes in the same operon according to cell needs. Nevertheless, direct experimentation is needed to confirm this hypothesis, and other mechanisms may exist to regulate translation rates. One alternate hypothesis is that transcriptional efficiency is involved. However, differences between expressivity of genes in the same transcription unit correspond to variation of P2. This indicates that transcriptional efficiency is not involved.

DARWINIAN EVOLUTION AND mRNA TRANSLATION

Whatever the process, codon usage optimization implies the existence of strong evolutionary pressures. This point of view has been largely developed by Ikemura (9,10). His main idea is a balance between selection for optimal codons and the mutational process. When a gene does not have a high translation rate, selective pressure is weak and the mutational process may blur its effect. For highly expressed genes the pressure is stronger and optimization clearly occurs.

We first want to point out that codon usage bias is the clearest example known of overall optimization of gene sequences. Selective pressure does not

act on one particularly important codon, but on the relative frequency of each codon. This is an argument in favour of Darwinian evolution as a process of cumulative small variations. Although we have no means to measure the phenotypic effect of one silent nucleotide substitution, it would be surprising if the fitness were greatly affected. Direct competition experiments between strains varying in codon usage would be of great evolutionary interest.

Our second remark derives from comparisons among P2 index values. Table 2 in (2) indicates that weakly expressed genes are not free of constraints on C/U choice in codon position III but that the rule is opposite that for highly expressed genes. Ikemura's evolutionary scheme, wherein each gene independently evolves toward an optimum, does seem valid for tRNA considerations, but may not be true for C/U choice. Our study indicates that the selection target also includes harmonization of expressivity and C/U choice, particularly when the genes are in the same operon. This would explain the positions of unc genes in fig. 1. Although we do not find data on ATP-synthase (unc operon product) cell content, this enzyme is most probably not as abundant as r-proteins or elongation factors. Hence, the high use of optimal codons in unc3, uncA and uncD genes has to be explained by their expressivity relative to that of remaining genes in the unc operon, rather than by their absolute expressivity.

Finally, the relationship between codon usage and gene expressivity demonstrated in E. coli may be valid for eukaryotic organisms as well. Indeed, a similar correlation between codon usage, tRNA population and gene expressivity has been recently described in the eukaryotic, undifferentiated yeast cell (30). However, the yeast catalog of "optimal" codons is different from that in E. coli probably because anticodon sequences, codon recognition patterns and relative concentrations of iso-tRNA species differ widely between these two organisms.

REFERENCES

1. Grantham, R., Gautier, C. and Gouy, M. (1980) Nucleic Acids Res. 8, 1893-1912.
2. Grantham, R., Gautier, C., Gouy, M., Jacobzone, M. and Mercier, R. (1981) Nucleic Acids Res. 9, r43-r74.
3. Lawther, R.P., Nichols, B., Zurawski, G. and Hatfield, G.W. (1979) Nucleic Acids Res. 7, 2289-2301.
4. Lawther, R.P., Calhoun, D.H., Adams, C.W., Hauser, C.A., Gray, J. and Hatfield, G.W. (1981) Proc. Natl. Acad. Sci. USA 78, 922-925.
5. Garel, J.-P., Mandel P., Chavancy, G. and Daillie, J. (1970) FEBS lett. 7, 327-329.
6. Chavancy, G. and Garel, J.-P. (1981) Biochimie 63, 187-195.
7. Post, L.E., Strycharz, G.D., Nomura, M., Lewis, H. and Dennis, P.P. (1979) Proc. Natl. Acad. Sci. USA 76, 1697-1701.
8. Post, L.E. and Nomura, M. (1980) J. Biol. Chem. 255, 4660-4666.

9. Ikemura, T. (1981) J. Mol. Biol. 146, 1-21.
10. Ikemura, T. (1981) J. Mol. Biol. 151, 389-409.
11. Gouy, M. and Grantham, R. (1980) FEBS lett. 115, 151-155.
12. Gouy, M. (1981) Thèse de troisième cycle, Université Lyon I.
13. Bergmann, J. E. and Lodish, H.F. (1979) J. Biol. Chem. 254, 11927-11937.
14. Sanderson, K.E. (1976) Ann. Rev. Microbiol. 30, 327-349.
15. Grosjean, H., Sankoff, D., Min Jou, W., Fiers, W. and Cedergren, R.J. (1978) J. Mol. Evol. 12, 113-119.
16. Fiers, W. and Grosjean, H. (1979) Nature 277, 328.
17. Gay, N.J. and Walker, J.E. (1981) Nucleic Acids Res. 9, 3919-3926.
18. Gay, N.J. and Walker, J.E. (1981) Nucleic Acids Res. 9, 2187-2194.
19. Saraste, M., Gay, N.J., Eberle, A., Runswick, M.J. and Walker, J.E. (1981) Nucleic Acids Res. 9, 5287-5296.
20. Futai, M. and Kanazawa, H. (1980) Curr. Topics Bioenerg. 10, 181-215.
21. Fillingame, R.H. (1981) Curr. Topics Bioenerg. 11, 35-106.
22. Chamberlin, M.J. (1974) Ann. Rev. Biochem. 43, 721-775.
23. Lindahl, L., Post, L., Zengel, J., Gilbert, S.F., Strycharz, W.A. and Nomura, M. (1977) J. Biol. Chem. 252, 7365-7383.
24. Post, L.E. and Nomura, M. (1979) J. Biol. Chem. 254, 10604-10606.
25. Burton, Z., Burgess, R.R., Lin, J., Moore, D., Holder, S. and Gross, C.A. (1981) Nucleic Acids Res. 9, 2889-2903.
26. Little, R., Fiil, N.P. and Dennis, P.P. (1981) J. Bact. 147, 25-35.
27. Post, L.E., Strycharz, G.D., Nomura, M., Lewis, H. and Dennis, P.P. (1979) Proc. Natl. Acad. Sci. USA 76, 1697-1701.
28. Keller, E.B. and Calvo, J.M. (1979) Proc. Natl. Acad. Sci. USA 76, 6186-6190.
29. Grosjean, H.J., de Hénau, S. and Crothers, D.M. (1978) Proc. Natl. Acad. Sci. USA 75, 610-614.
30. Bennetzen, J.L. and Hall, B.D. (1982) J. Biol. Chem. 257, 3026-3031.
31. di Rienzo, J.M., Nakamura, K. and Inouye, M. (1978) Ann. Rev. Biochem. 47, 481-532.
32. Furano, A.V. (1975) Proc. Natl. Acad. Sci. USA 72, 4780-4784.
33. Gudas, L.J. and Pardee, A.B. (1976) J. Mol. Biol. 101, 459-477.
34. Pedersen, S., Bloch, P.L., Reeh, S. and Neidhardt, F.C. (1978) Cell 14, 179-190.
35. Zabin, I. and Fowler, A.W. (1980) in The Operon, Miller, J.H. and Reznikoff, W.S. Eds., Cold Spring Harbor Laboratory.
36. Morse, D.E., Mosteller, R.D. and Yanofsky, C. (1969) Cold Spring Harbor Symp. Quant. Biol. 34, 725-739.
37. Dennis, P.P. (1977) J. Mol. Biol. 115, 603-625.
38. Casabadan, M.J. (1976) J. Mol. Biol. 104, 557-566.
39. Rose, J.K. and Yanofsky, C. (1974) Proc. Natl. Acad. Sci. USA 71, 3134-3138.
40. Gilbert, W. and Müller-Hill, B. (1966) Proc. Natl. Acad. Sci. USA 56, 1891-1898.