CHEMICAL BASIS OF BIOLOGICAL SPECIFICITY

Alan R. Fersht

Department of Chemistry, Imperial College of Science and Technology, London SW7 2AY, U.K.

Abstract - The specificity of the interaction of biopolymers with one another and with small molecules in binding and catalysis is based primarily on the fit of complementary structures. When this is inadequate to generate the high fidelity required in the replication and synthesis of DNA and proteins, specificity is enhanced by editing or proofreading mechanisms. These add kinetic control of product formation. The magnitudes of the interaction energies of groups in proteins with each other and of nucleotide pairing may be measured from simple kinetic experiments with enzymes. The energies thus estimated are far higher than those expected from experiments on simple models.

INTRODUCTION

Physical-organic chemistry has made not only an enormous contribution to our understanding of the chemical basis of enzyme catalysis but has also provided much of the framework of our knowledge of molecular interactions in aqueous solution. Work in this area has involved simple physical studies on the interaction of small solute molecules, direct and indirect measurements on the interactions of biopolymers with each other and with small molecules and, in an expanding area, the direct computation or computer simulation of these interactions using structural data from X-ray crystallography. I am interested in the interplay of catalysis and molecular interactions in providing specificity in biological systems. The method of investigation chosen is that of direct measurements on the enzymes and proteins involved in these processes. Some of the enzyme systems studied seem at first sight to be rather too complex and intractable for the extraction of information on simple processes. However, I hope to show during this lecture that, by restricting measurements in certain complex systems to the direct observation of specific and identifiable processes and by the application of basic rules and ideas of physical-organic chemistry that are used to analyse simple systems, fundamental information may be gathered. Some of the data so obtained show how model systems severely underestimate the strength of molecular interactions.

There are two components of biological specificity: thermodynamic, binding through simple molecular interactions; kinetic, the control of product formation through rates. In simple enzymic systems, catalysis and binding are intimately linked (1,2,3) and it is in some ways meaningless to separate their individual contributions (see below eq 5). There are, however, specific enzymic processes that use additional kinetic pathways to boost specificity by providing genuine kinetic control of reaction products. The chemical basis of biological specificity may accordingly be divided into the two following components.

1. Complementarity - thermodynamic control
In 1946, Pauling published a paper of seminal importance to enzymology stating explicitly that the structure of the active site of an enzyme should be complementary to the structure of the activated complex of the substrate, rather than to the substrate itself (4). This was, however, just one example of biological specificity that he quoted in an article which emphasized the importance of *shape* and *complementarity* in biological reactions. The idea of specificity in biology resulting from the interaction of complementary structures stretches back to Emil Fischer with the "lock-and-key" hypothesis for enzyme specificity and to Paul Ehrlich for the general application of this concept to the physiological activity of substances. Ehrlich's description of specificity in such terms as "the only substances that can be anchored at a particular part of an organism are those that fit into the molecule of the receptor as a piece of mosaic fits into a certain pattern" (paraphrased from ref (4)) was eventually visualised when the molecular structure of haemoglobin was solved by Perutz and colleagues (5). The subunit interfaces were seen to be close-packed, with the residues from each subunit interlocking as the pieces of Ehrlich's mosaic. Complementarity in the interaction of proteins with small molecules was seen by X-ray crystallographers in their studies on enzymes and their complexes with substrates and inhibitors. In particular, studies of the structure of the first enzyme to be solved, lysozyme, showed that its binding

with its substrate fitted precisely Pauling's (and Haldane's (6)) hypothesis (7).
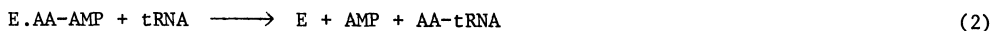
Complementarity is the basis of biological specificity in the interaction of macromolecules with one another in such processes as the assembly of larger units (e.g. oligomers) or for purposes of recognition during metabolism and antibody-antigen interactions. The specificity in the binding of small molecules by proteins is also based on complementarity. In certain cases, as demonstrated for the binding of allosteric effectors to allosteric proteins or enzymes, and postulated for processes such as the binding of physiologically active compounds to receptors, the small molecule is complementary in structure to an altered conformation of the protein. The binding thus distorts the protein and so transmits the information that binding has occurred.

The other great area dominated by complementary interactions is that of the replication of biological macromolecules. The genetic information is stored and duplicated in the structure of DNA which contains two complementary chains of paired nucleotides. Complementarity of base pairing is the essential feature in the replication of DNA and the transcription and translation of the encoded information into protein structure. However, complementary interactions by themselves are not adequate to maintain a sufficiently high accuracy in the replication of DNA and the synthesis of proteins and so a further mechanism has evolved, *editing* or *proofreading* (for a general review see ref. 8).

## 2. Editing or Proofreading - kinetic control

The overall error rate in the replication of DNA in bacteria is only one mistake per $10^{8-10}$ nucleotides polymerised (9) and the overall error rate of protein synthesis is only one mistake per $10^{3-4}$ amino acid residues incorporated (10-13). The specificity in base pairing is inadequate to account for the accuracy of DNA replication: the phenomena of keto-enol and amino-imino tautomerism scramble the specificity so that $G:T_{enol}$ base pairs, for example, are expected to occur at a frequency of $10^{-4}-10^{-5}$ as are certain other types of non-complementary base pairs (14). In protein synthesis, as pointed out initially by Pauling (15), certain amino acids are so similar to each other in structure that non-complementary pairing of amino acids with the selection of sites will occur at high levels. For example, glycine must bind to a cavity constructed to fit alanine since there are no forces of steric exclusion to preclude the binding of the smaller substrate. The binding of glycine will be weaker than that of alanine by just the binding energy of the additional methylene group of alanine to the alanine-binding site. This binding energy is smaller than that necessary to ensure a discrimination factor of $10^{3-4}$, i.e. 16-22 kJ/mol. Because of this, special biochemical pathways have evolved to reduce the error rate. For example, in the case of bacterial DNA replication, the DNA polymerase has a hydrolytic editing site. This is a $3' \rightarrow 5'$ exonuclease activity that is specific for the excision of mismatched base pairs from the chain of DNA that is being elongated in the $5' \rightarrow 3'$ direction (16). Similarly, certain aminoacyl-tRNA synthetases, the enzymes responsible for the selection of amino acids during protein synthesis, possess an esterolytic site that is specific for the hydrolysis of mischarged tRNAs (8).

The basic principles of editing are well-illustrated by the reactions of the aminoacyl-tRNA synthetases. The process of charging tRNA is a two-step reaction consisting of (1) activation of the amino acid, (2) transfer to the tRNA. The first reaction is frequently

$$E + AA + ATP \rightleftharpoons E.AA-AMP + PP_i \qquad (1)$$

$$E.AA-AMP + tRNA \longrightarrow E + AMP + AA-tRNA \qquad (2)$$

imprecise, with 'wrong' amino acids being activated. The overall reaction, however, is always very precise for the aminoacylation of tRNA with naturally occurring amino acids. The paradox of how the overall reaction is more accurate than one of the individual components was solved when it was discovered that the addition of tRNA to the complex of an aminoacyl-tRNA synthetase (that for isoleucine) and a wrong aminoacyl adenylate (Val-AMP) led to its quantitative hydrolysis rather than the formation of the mischarged tRNA (17). It has subsequently been shown that this editing reaction results from the aminoacyl-tRNA synthetase having an esterolytic site that specifically and rapidly deacylates the products of reaction (2) of the wrong aminoacyl adenylate (18,19,20). Editing has also been shown to occur in some cases at the level of aminoacyl adenylate (21).

The basic requirements of an editing mechanism are: $i$, the existence of a thermodynamically unstable intermediate or product on a reaction pathway; $ii$, a branch point on the pathway at which the unstable intermediate may be either channelled to synthesis or to destruction. The unstable 'intermediate' in DNA replication is the phosphodiester bond and in amino-acylation of tRNA the aminoacyl adenylate or the aminoacyl ester bond. The channelling to synthesis or destruction is governed by the molecular interactions of the intermediate or product with the synthetic and hydrolytic sites of the enzymes. I shall now present measurements of some of the important interaction energies responsible for binding specificity.

### MEASUREMENT OF THE INTERACTION ENERGIES OF AMINO ACID SIDE CHAINS AND PROTEINS

The energetics of folding of proteins are dominated by the interaction of amino acid residues of the chain with each other either by the packing of their side chains or by hydrogen bonding of the backbone NH and carbonyl groups. The interactions include a major contribution from the so-called hydrophobic bond (22). This bonding is also important in the binding of substrates to enzymes. The magnitudes of the interactions have been traditionally measured from the equilibrium constants for the partition of molecules between aqueous and organic phases. It is clearly more relevant to proteins to measure these energies from the interactions of proteins themselves. The aminoacyl-tRNA synthetases are particularly useful in this context since they have been subject to evolutionary pressure to maximise the binding energy of their 'right' substrate relative to that of their wrong substrates. The energetics of such interactions may be derived from the application of transition state theory to the kinetic quantities derived from the Michaelis-Menten equation (eq 3). The important kinetic constant in the analysis is the compound quantity $k_{cat}/K_M$

$$v = k_{cat}[E_o][S]/(K_M + [S])$$ (3)

because it both determines specificity and is sensitive to the contributions from binding energies (1,3). When two (or more) substrates A and B compete for the active site of an enzyme, their relative rate of reaction is given by:

$$v_A/v_B = (k_{cat}/K_M)_A[A]/(k_{cat}/K_M)_B[B]$$ (4)

Application of transition state theory shows that the value of $k_{cat}/K_M$ is related to the chemical activation energy of the reaction ($\Delta G^{\ddagger}$) and the binding energy of the enzyme and the substrate in the transition state ($\Delta G_b$) by:

$$RT\ln(k_{cat}/K_M) = \ln(kT/h) - \Delta G^{\ddagger} - \Delta G_b$$ (5)

(where $k$ is the Bolzmann's constant, $h$, Planck's constant and $T$ the absolute temperature). Therefore for two substrates A and B reacting with the enzyme:

$$RT\ln[(k_{cat}/K_M)_A/(k_{cat}/K_M)_B] = (\Delta G^{\ddagger} + \Delta G_b)_B - (\Delta G^{\ddagger} + \Delta G_b)_A.$$ (6)

where there is competition of a correct substrate (A) with a smaller wrong substrate (B) that lacks a binding group R, eq 6 reduces to eq 7 when the chemical activation energies are identical and $\Delta G_R$ is the binding energy of the group R relative to its absence in B. It is seen in eq 5 that both binding and chemical activation energies contribute to the overall

$$RT\ln[(k_{cat}/K_M)_A/(k_{cat}/K_M)_B] = - \Delta G_R$$ (7)

activation energy of ($k_{cat}/K_M$) and hence the difficulty mentioned above of separating kinetic and binding contributions to specificity. However, it is seen in eq 7 that the binding energy of a group R on a substrate may be calculated from the relative values of $k_{cat}/K_M$ of substrates containing and lacking that group. Further, it may be inferred from eqs 4 and 7 that any evolutionary pressure to maximise the specificity for A versus B will lead to a maximisation of the binding energy of the group R. Thus, the comparison of the kinetics of activation of alanine and glycine by the alanyl-tRNA synthetase, isoleucine and valine by the isoleucyl-tRNA synthetase and valine and α-aminobutyrate by the valyl-tRNA synthetase will give measurements of the binding energy of the additional methylene group of the larger substrates. In all three cases, this is calculated to be 13-14 kJ mol$^{-1}$ (23-25). The constancy of the effect of the extra methylene group of the larger substrate, independent of its distance from the seat of reaction, is consistent with changes in $k_{cat}/K_M$ being caused by the binding energy and not inductive effects on the chemical activation energy, as predicted from structure-reactivity considerations (26). The values for the binding energies are far higher than predicted from experiments with model compounds where figures of 2-4 kJ mol$^{-1}$ are typically found. Values for other interaction energies are listed in Table 1. These are also all much higher than predicted from simple partition experiments. The measurements with the aminoacyl-tRNA synthetases provide the upper limits of these interaction energies and represent the optimal values for specificity and catalysis. These approach the 'intrinsic' binding energies (27). Where tight binding is not required, lower values of the binding energy are manifested. Also listed in Table 1 are some values of repulsion energies, e.g. the unfavourable energy of cramming the additional methylene group of isoleucine into a valine-binding site, the additional -OH group of serine into an alanine-binding site.

The binding energies listed in Table 1 indicate where the problems are most acute in the selection of amino acids and how great specificity may be attained in editing by juggling the sets of values. As predicted by Pauling (15), the major difficulty is in rejecting incorrect substrates that are slightly smaller than the correct. On the other hand, the repulsion energies are sufficiently high to reject larger substrates at an adequate level.

These principles are embodied in the 'double-sieve' model for editing (28):  the activation site of the aminoacyl-tRNA synthetase rejects at a tolerable level (with one exception (23)) substrates larger than the correct whilst it accepts, to varying degrees, isosteric and smaller substrates; conversely, the editing site rejects the products of the correct substrate by steric hindrance or by specific chemical effects whilst the products of smaller or isosteric products are accepted and edited.  For example, although the energy available for discrimination against valine by a site tailored to accept isoleucine is only 13 kJ mol$^{-1}$, the repulsion energy for excluding isoleucine from a hydrolytic site tailored to accept valine derivatives is possibly greater than 30 kJ mol$^{-1}$.
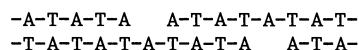
TABLE 1.   Incremental group binding energies

| Binding cavity in protein | | Unfavourable Gibbs energy kJ mol$^{-1}$ |
|---|---|---|
| constructed for: | occupied by: | |
| $CH_3-$ | $H-$ | 13 |
| $-S-$ | $H-$ | 21 |
| $HO-$ | $H-$ | 29 |
| $CH_3-$ | $HO-$ | 15 |
| $H-$ | $CH_3-$ | > 30 |
| $H-$ | $HO-$ | 15 |

Determined from relative Gibbs energies of transfer of substrates from $H_2O$ to aminoacyl-tRNA synthetases at 25 $^{\circ}$C (23,26).

MEASUREMENT OF BASE PAIRING ENERGIES DURING DNA REPLICATION

The accuracy of replication of DNA relies on complementary base pairing.  DNA polymerases have just one active site (for synthesis) to bind all four combinations of correct base pairs.  This recognises the overall shape of the pair whilst delegating specificity to the hydrogen bonding (29).  Unfortunately, the four nucleotides can form, albeit at a low frequency, complementary base pairs other than the Watson-Crick varieties that are necessary to conserve the genetic information (14).  We have been able to measure the relative frequencies of correct Watson-Crick pairings and the incorrect variants by a simple kinetic method using DNA polymerase I from *Escherichia coli* and a synthetic template, the self-complementary alternating co-polymer *poly*d(A-T) (30).  The template forms long duplexes in solution as the individual strands pair and overlap thus:

```
-A-T-A-T-A    A-T-A-T-A-T-A-T-
-T-A-T-A-T-A-T-A-T-A    A-T-A-
```

On addition of the polymerase and just one complementary deoxynucleoside triphosphate, e.g. dATP, all the 3'-termini will be 'anchored' as A.  When a non-complementary deoxynucleoside triphosphate, such as dGTP, is added to this system, the polymerase will add dGMP to the 3'-terminus and then catalyse its *removal* by the 3'→5'-exonuclease activity. The enzyme is thus a dGTPase under        these conditions.  The mispairing involved is G:A since all the 3'-termini are A and hence the next complementary base is A.  By measuring $k_{cat}/K_M$ for the dGTPase activity and comparing it with $k_{cat}/K_M$ for the incorporation of dTMP on replacing the dGTP with dTTP, the relative mispairing frequency of G:A for T:A may be calculated (= ratio of values of $k_{cat}/K_M$ for dGTPase and incorporation).  Permutation of all the combinations of dATP, dTTP, dGTP and dCTP gives the relative mispairing frequencies listed in Table 2.

TABLE 2.   Mispairing frequencies during DNA replication

| Incorrect pair | Correct pair | Ratio of correct to incorrect frequencies |
|---|---|---|
| G:T | A:T | $1.2 \times 10^4$ |
| G:A | T:A | $9 \times 10^4$ |
| C:T | A:T | $5.5 \times 10^6$ |
| C:A | T:A | $4.5 \times 10^4$ |

The most frequent mispair is that of G and T, occurring at a frequency of 1/12000 relative to the Watson-Crick pair, A:T. The C:A mispair occurs only some three times less frequently. Topal and Fresco (14) have argued that these arise via complementary pairing of the unfavoured enol or imino tautomeric forms of one of the purine or pyrimidine components with the major tautomer of the other pyrimidine or purine isomer. The resultant purine:pyrimidine mispairs have the same overall shape as a correct pair and so are accommodated by the active site of the DNA polymerase and by the DNA double helix. Model-building studies can also fit purine:purine pairs into the double helix by invoking the *syn* isomers of the purine (14). These mispairs are in fact found at frequencies of about $10^{-5}$. Two pyrimidines cannot pair in the double helix but the DNA polymerase does catalyse their association at low frequency, as seen in Table 2. In this case, it is presumably just the binding of the nucleotide to the enzyme that accounts for the apparent pairing. Whatever the mechanism for the mispairings, they occur at too high a frequency for the faithful replication of DNA. Without editing to reduce this error level, the spontaneous rate of mutation would be unacceptably high.

### MEASUREMENT OF ACCURACY OF DNA REPLICATION: 'KINETIC GENETICS'

The accuracy of DNA replication is difficult to measure *in vitro* by chemical methods because the fidelity is so high. We have developed a technique for measuring the error rate by adapting the traditional methods of genetics to a controlled genetic procedure (31,32). The classical method of measuring the fidelity of DNA replication *in vivo* is by measuring the rate of reversion of a point mutation in a bacterium or bacteriophage. This may be performed *in vitro* by replicating the DNA from mutants of a bacteriophage ($\phi$X174) using the same enzyme system that replicates the DNA *in vivo* (33). The reversion back to wild type may be induced by biasing the concentrations of deoxynucleoside triphosphates (dNTPs) according to eq 4. For example, to synthesize TGG from the template codon ATC, the reaction must be biased to favour G:T rather than A:T pairing by raising the concentration of dGTP and lowering dATP. The reaction products may be analysed using a bioassay: the synthetic DNA prepared *in vitro* is added to *E. coli* cells that have been treated to take up DNA; the cells then produce bacteriophage from the DNA. The proportion of revertant and mutant phages so formed may then be titrated on indicator bacteria by counting plaques. By this long-winded procedure, it has proved possible to produce respectable kinetic plots of the number of revertants produced against the ratios of dNTPs in the reaction mixture and to derive rate laws for mutation (Table 3). The results appear valid: when the concentration of dNTPs that are found *in vivo* are inserted into the rate equations to calculate the mutation frequency *in vivo*, there is satisfactory agreement with genetic experiments.

TABLE 3. Fidelity of DNA replication *in vitro*[a]

| Mutation | Base pairing | Rate law for production |
|----------|--------------|-------------------------|
| TAG → GAG | G:A | $1 \times 10^{-6} [dGTP]/[dTTP]$ |
| TAG → AAG | A:A | $3 \times 10^{-7} [dATP]/[dTTP]$ |
| TAG → CAG | C:A | $4 \times 10^{-7} [dCTP]/[dCTP]$ |
| TAG → TGG | G:T | $1 \times 10^{-2} [dGTP]^2/[dATP]$ |
| TAG → TCG | C:T | $1 \times 10^{-4} [dCTP][dGTP]/[dATP]$ |

[a]Replication of $\phi$X174 am16 RF DNA by DNA polymerase III holoenzyme (32).

The frequencies of misincorporation are lower than those of mispairing listed in Table 2 because editing by the $3' \rightarrow 5'$ exonuclease activity removes most of the errors. The existence of editing necessitates the modification of the rate law of eq 4 for product formation when two (or more) products compete for the active site of an enzyme. Since permanent misincorporation requires that elongation of the growing chain occur and the rate of elongation depends to some extent on the concentrations of the next nucleotide in the sequence to be incorporated ($dNTP_{fol}$), the misincorporation frequency $\nu$ is of the general form:

$$\nu = \frac{\beta [dNTP]_{inc}}{[dNTP]_{cor}} \cdot \frac{[dNTP]_{fol}}{(K + [NTP]_{fol})} \tag{8}$$

(where subscripts inc = incorrect, cor = correct, and $K$ is a compound constant composed of various rate constants as is $\beta$ (31)). Thus, the rate of spontaneous mutation by misincorporation during DNA replication follows simple kinetic rate laws.

### COST OF EDITING: THE COST-SELECTIVITY EQUATION

Editing or proofreading of incorrect intermediates implies that some of the correct inter-
mediates are also excised (34). The fraction of correct substrate hydrolysed is termed the
cost, $C$. A simple steady-state kinetic analysis based on partitioning of intermediates
shows that there is a simple relationship between the overall accuracy (the selectivity, $S$)
of an enzyme that edits and its discrimination between substrates in initial selection and
editing (8,30). If the initial discrimination between the right and wrong substrate is $f$
(i.e. the ratio of values $k_{cat}/K_M$ for the right:wrong) and the discrimination of editing
is $f'$ (the more rapid editing of the incorrect products) then:

$$S = f(1 + (f' - 1)C) \tag{9}$$

(Eq 9 is the cost-selectivity equation). For DNA replication, $f$ is equal to the reciprocal
of the mispairing frequency (as in Table 2) and $S$ is equal to the reciprocal of the
misincorporation frequency (as in Table 3). Experimentally, we find the cost for DNA
replication relatively high, $\sim 10\%$ (30), whilst the cost of amino acid selection is low (35).
The discrimination factor $f'$ has been calculated from the measurements of $S$, $C$ and $f$ using
eq 9 (30). In all cases, $f'$ is less than $f$ showing that the largest contribution to
specificity is the accuracy of the initial base pairing.

### REFERENCES

1. A.R. Fersht, Proc. R. Soc. Lond. B 187, 397-407 (1974).
2. W.P. Jencks, Adv. Enzymol. 43, 219-410 (1975).
3. A.R. Fersht, Enzyme Structure and Mechanism, W.H. Freeman, Oxford, U.K. and San
   Francisco, CA (1977).
4. L. Pauling, Chem. Eng. News 24, 1375-1377 (1946).
5. M.F. Perutz, H. Muirhead, J.M. Cox and L.C.G. Goaman, Nature 219, 131-139 (1968).
6. J.B.S. Haldane, Enzymes, Longmans Green and Co., p. 182 (1930).
7. C.C.F. Blake, L.N. Johnson, G.A. Mair, A.C.T. North, D.C. Phillips and V.R. Sarma,
   Proc. R. Soc. B167, 378-388 (1967).
8. A.R. Fersht, Proc. R. Soc. Lond. B212, 351-379 (1981).
9. E.C. Cox, Ann. Rev. Genet. 10, 135-56 (1976).
10. R.B. Loftfield, Biochem. J. 89, 82-92 (1963).
11. R.B. Loftfield and D. Vanderjagt, Biochem. J. 128, 1353-1356 (1972).
12. P. Edelmann and J. Gallant, Cell 10, 131-137 (1977).
13. J. Parker, T.C. Johnston and P.T. Borgia, Molec. Gen. Genet. 180, 275-281 (1980).
14. M.D. Topal and J.R. Fresco, Nature 263, 285-289 (1976).
15. L. Pauling, Festschr. Prof. Dr. Arthur Stoll Siebzigsten Geburstag, 597-602 (1958).
16. D. Brutlag and A. Kornberg, J. Biol. Chem. 247, 241-248 (1972).
17. A.T. Norris and P. Berg, Proc. Natl. Acad. Sci. 52, 330-337 (1964).
18. M. Yarus, Proc. Natl. Acad. Sci. 69, 1915-1919 (1972).
19. E.W. Eldred and P. Schimmel, J. Biol. Chem. 247, 2961-2964 (1972).
20. A.R. Fersht and M. Kaethner, Biochemistry 15, 3342-3346 (1976).
21. J. Jakubowski and A.R. Fersht, Nucl. Acids. Res. 9, 3105-3117 (1981).
22. W. Kauzmann, Adv. Protein Chem. 14, 1-63 (1957).
23. W.-C. Tsui and A.R. Fersht, Nucl. Acids. Res. 9, 4627-4637 (1981).
24. R.B. Loftfield and E.A. Eigner, Biochim. Biophys. Acta 130, 426-448 (1966).
25. A.R. Fersht and C. Dingwall, Biochemistry 18, 1238-1245 (1979).
26. A.R. Fersht, J.S. Shindler and W.-C. Tsui, Biochemistry 19, 5520-5524 (1980).
27. W.P. Jencks, Proc. Natl. Acad. Sci. USA 78, 4046-4050 (1981).
28. A.R. Fersht and C. Dingwall, Biochemistry 18, 2627-2631 (1979).
29. A. Kornberg, DNA Replication, W.H. Freeman and Co. San Francisco Cal. (1980).
30. A.R. Fersht, J.W. Knill-Jones and W.-C. Tsui, J. Mol. Biol. 156, 37-51 (1982).
31. A.R. Fersht, Proc. Natl. Acad. Sci. USA 76, 4946-4950 (1979).
32. A.R. Fersht and J.W. Knill-Jones, Proc. Natl. Acad. Sci. USA 78, 4251-4255 (1981).
33. L.A. Weymouth and L.A. Loeb, Proc. Natl. Acad. Sci. USA 75, 1924-1928 (1978).
34. N. Muzyczka, R.L. Poland and M.J. Bessman, J. Biol. Chem. 247, 7116-7122 (1972).
35. R.S. Mulvey and A.R. Fersht, Biochemistry 16, 4731-4737 (1977).