

# Modeling of Nucleotide Binding Domains of ABC Transporter Proteins Based on a F<sub>1</sub>-ATPase/recA Topology: Structural Model of the Nucleotide Binding Domains of the Cystic Fibrosis Transmembrane Conductance Regulator (CFTR)\*

Mario A. Bianchet,<sup>1</sup> Young Hee Ko,<sup>2</sup> L. Mario Amzel<sup>1</sup> and Peter L. Pedersen<sup>2</sup>

Received July 1, 1997; accepted November 1, 1997

Members of the ABC transporter superfamily contain two nucleotide binding domains. To date, the three dimensional structure of no member of this super-family has been elucidated. To gain structural insight, the known structures of several other nucleotides binding proteins can be used as a framework for modeling these domains. We have modeled both nucleotide binding domains of the protein CFTR (Cystic Fibrosis Transmembrane Conductance Regulator) using the two similar domains of mitochondrial F<sub>1</sub>-ATPase. The models obtained, provide useful insights into the putative functions of these domains and their possible interaction as well as a rationale for the basis of Cystic Fibrosis causing mutations. First, the two nucleotide binding domains (folds) of CFTR are each predicted to span a 240–250 amino acid sequence rather than the 150–160 amino acid sequence originally proposed. Second, the first nucleotide binding fold, is predicted to catalyze significant rates of ATP hydrolysis as a catalytic base (E504) resides near the  $\gamma$  phosphate of ATP. This prediction has been verified experimentally [Ko, Y.H., and Pedersen, P.L. (1995) *J. Biol. Chem.* **268**, 24330-24338], providing support for the model. In contrast, the second nucleotide binding fold is predicted at best to be a weak ATPase as the glutamic acid residue is replaced with a glutamine. Third, F508, which when deleted causes ~70% of all cases of cystic fibrosis, is predicted to lie in a cleft near the nucleotide binding pocket. All other disease causing mutations within the two nucleotide binding domains of CFTR either reside near the Walker A and Walker B consensus motifs in the heart of the nucleotide binding pocket, or in the C motif which lies outside but near the nucleotide binding pocket. Finally, the two nucleotide binding domains of CFTR are predicted to interact, and in one of the two predicted orientations, F508 resides near the interface. This is the first report where both nucleotide binding domains of an ABC transporter and their putative domain-domain interactions have been modeled in three dimensions. **The methods and the template used in this work can be used to analyze the structures and function of the nucleotide binding domains of all other members of the ABC transporter super-family.**

**KEY WORDS:** Nucleotide Binding Domain; Cystic Fibrosis; ABC transporters; Traffic ATPases; CFTR.

<sup>1</sup> Department of Biophysics and Biophysical Chemistry.

<sup>2</sup> Department of Biological Chemistry, The Johns Hopkins University, School of Medicine, 725 N Wolfe Street, Baltimore, MD, 21205-2185.

\* Supported by Grants from the NIH (NIDDK) and the Cystic Fibrosis Foundation to PLP. Supported by Grant GM25432 from the NIH to LMA.

## INTRODUCTION

Members of the ABC transporter super-family [Hyde et al. (1990); Higgins (1992)] share extensive sequence similarities. The family, which includes, the Cystic Fibrosis Transmembrane Conductance Regulator (CFTR), P-glycoprotein (MDR1, a permease involved in the transport of a wide variety of drugs), the Ste6 gene product, pfMDR, Histidine permease (HisP), and many others. All these proteins contain membrane domains as well as large hydrophilic domains associated with binding of ATP. Because several use hydrolysis of ATP to support substrate accumulation, they are also called "Traffic ATPases" [Doige & Ames (1993)]. Members of the ABC transporter protein family have three characteristic sequence signatures: (A)  $GX_4GK[T/S]$  (also called "P-loop"), (B)  $RX_{6-8}h_4D$  and (C)  $LSXGX[R/K]$  (X: any residue; and h: hydrophobic residue). The first two of these signatures, also called "Walker motifs" A and B, are involved in ATP binding [Walker et al. 1982].

One of the better known members of the ABC transporter super-family, the CFTR protein, is a polypeptide related to the disease Cystic Fibrosis (CF). CF, an autosomal recessive genetic disease predominant among Caucasians, affects numerous organs, including lung and airways, pancreas and sweat glands. The disease symptoms are presumably caused by a decrease in chloride ion conductance across the apical membrane in epithelial cells that results from mutations in the gene that encodes CFTR, a 1480 amino acid chloride channel [Riordan et al. (1989)]. Of the mutations found in the CF gene, 42% are missense, 21% are nonsense, and 10% alter nucleotides essential for codon splicing [Tsui (1992)]. Approximately 70% of all cases of CF are caused by a single deletion mutation ( $\Delta F508$ ) [Kerem et al. (1990)]. This mutation causes severe disease with symptomatology that includes severe lung disease, by preventing trafficking of the protein to its final position in the plasma membrane [Cheng et al. (1990), Denning et al. (1992), Dalemans et al. (1991), Qu et al. (1996)]. Another deletion close to codon 508,  $\Delta I507$  produces a very similar phenotype. Importantly, missense mutations such as F508C, I506V and I507V are benign and do not cause the disease [Kobayashi et al. (1990)]. Other missense mutations are found in or near the motif  $GX_4GK[T/S]$ , i.e., A455E and G458V. Significantly, mutation of G551, (G551D) in motif C ( $LSXGX[R/K]$ ) is associated with high chloride levels in sweat,

pancreatic insufficiency and variable lung disease, producing a phenotype very similar to  $\Delta F508$  [Hamosh et al. (1992)]. Interestingly, mutation of R553 can partially suppress the effect of  $\Delta F508$ , [Teem et al. (1993)].

Optimal activation of the CFTR channel requires phosphorylation by cAMP dependent kinases [Cheng et al. (1991)], and recently, Li et al. (1996) have presented direct evidence that ATP hydrolysis is used to gate, or to modulate the CFTR chloride channel. The A and B motifs associated with the nucleotide binding site of CFTR and ABC transporters are found also in proteins of the nucleotide binding family [Walker et al. (1982)] with known three-dimensional structures, such as adenylate kinase (ADK) [Schulz et al. (1974)], elongation factor-tu [Morikawa et al. (1978)], ras p21 [Pai et al. (1990)], transducin (a G-protein) [Noel et al. (1993)], recA [Story et al. (1992)] and the  $F_1$  sector of  $F_0F_1$ -ATPase (ATP synthase) [Abrahams et al. (1994), Blanchet et al. (1998)]. Most of these proteins use the energy of hydrolysis of (G/A)TP to bring about conformational changes that transmit signals or perform work. Of these proteins, only  $F_0F_1$ -ATPase is involved in ion transport across a biological membrane. In animal cells, this enzyme synthesizes ATP using a proton gradient across the inner mitochondrial membrane as the driving force (for reviews see [Weber and Senior (1997) and Pedersen and Amzel (1993)]. Operating in the opposite direction,  $F_0F_1$ -ATPase can pump protons across the membrane against their electrochemical gradient utilizing the free energy of ATP hydrolysis.  $F_1$ -ATPase, the soluble part of the ATP synthase complex is a multiple subunit protein made up of five different polypeptide chains. The two larger subunits,  $\alpha$  and  $\beta$ , contain the sites for synthesis or hydrolysis of ATP. Recently, atomic resolution structures of substantial parts of  $F_1$ -ATPase were elucidated, for the bovine heart  $F_1$  [Abrahams et al. (1994)], for the  $\alpha_3\beta_3$  complex of the  $F_1$  of *Thermophilium Bacterium* [Shirakihara et al. (1997)], and for the  $F_1$  of rat liver [Blanchet et al. (1998)]. The features described above make  $F_1$ -ATPase a suitable template for modeling the regions of ABC transporters associated with binding of nucleotides. This paper presents the three-dimensional models of two such regions of CFTR based on  $F_1$ -ATPase. In developing this model we utilized features common to all ABC transporters. Therefore, modeling other members of this super-family (i.e., MDR, Permeases, etc.) using the alignment and assignment of struc-

tural elements presented in this paper should be straightforward.<sup>3</sup>

## METHODS

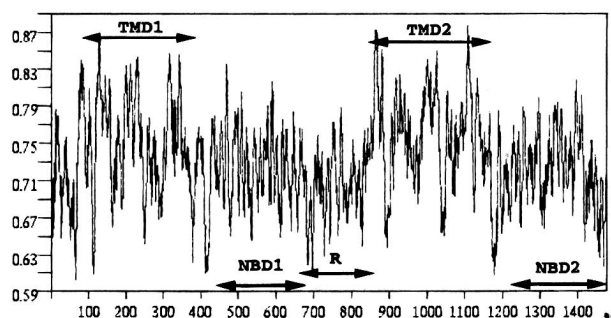
Although NBD1 and NBD2 have been defined in the literature as residues F433 through S589, and Y1218 through R1386 respectively [Riordan et al. (1989)], we define them here as residues L441 through K684 and L1227 through L1480, respectively. The justification for this is that this is the minimum number of residues that can make a complete nucleotide binding domain using the modeling approach described here. The sequences of NBD1 and NBD2 were aligned manually with the sequences of the nucleotide domains of the  $F_1\alpha$ - and  $F_1\beta$ -subunits taking into consideration the criteria discussed in the following sections. Coordinates for the equivalent atoms in the aligned residues were obtained from the template molecule, rat liver  $F_1$ -ATPase. The non common atoms were generated automatically and side chain orientations were optimized by visual inspection. Insertions were built using the program O [Jones et al. (1990)] and the molecular modeling program Quanta [MSI Inc.] based on the predicted secondary structure when it was compatible with the topology assigned. The coordinates of the initial model are optimized using several picoseconds of molecular dynamics at 300 K, with the version 2.3 of CHARMM in the Quanta 41 software. For drawing and visualization the programs Setor [Evans (1993)] and O were used. The geometry and quality of the models was assessed using the protein package of Quanta 41. Polar and apolar accessible area were calculated using the programs AREAIMOL, SURFACE and DIFAREA of the CCP4 suite of program [CCP4 (1994)]. Accessibility surfaces and electrostatics superficial potential were obtained with the program GRASP [Nicholls et al. (1991)]. The sequences of ABC transporters were aligned using the program MAXHOM, a neural network multiple sequence alignment included in the PHD package. (<http://www.embl.heidelberg.de/predictprotein>) [Rost (1996)].

<sup>3</sup> This work was first presented in abstract form at both the Ninth Annual North American Cystic Fibrosis Conference (Dallas, Texas) [Bianchet et al. (1995a)], and the Biophysical Society Meeting (Baltimore Md.) [Bianchet et al. (1995b)].

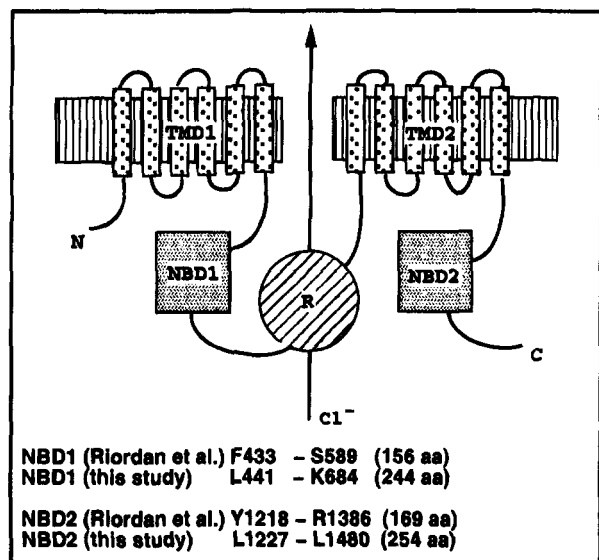
## RESULTS AND DISCUSSION

### Domain Organization of ABC Transporters

Several groups have proposed models of complete ABC transporters based on analysis of hydropathy profiles [Riordan et al. (1989), Petsko (1990)]. Significantly, CFTR is a single-chain multiple domain protein, with a hydropathy profile (Figure 1) characteristic of the ABC super-family of proteins with the expected regions of defined or dominant hydropathy. The two halves of this profile seem to be similar to each other, suggesting gene duplication. However in the case of CFTR, the lack of conservation in the relative positions of exon-intron segments argues against this [Riordan et al. (1989)]. The CFTR regions between residues 50 to 404 and 850 to 1150 are predominantly hydrophobic and associated with transmembrane spanning domains (TMD1 & TMD2 respectively). A straightforward prediction of topological transmembrane helices suggests 6 helices per transmembrane domain [for a review of the methods see White (1994)]. Two amphipathic regions between residues 441 and 687 and between residues 1227 and 1480 (the C-terminus) contain the characteristic ABC sequences. These two regions have similar size, hydropathy characteristics, and sequences homologies beyond the ABC sequences. Probably, the regions fold in a similar way forming two domains, that we refer to here as nucleotide binding domains NBD1 and NBD2. Traditionally, on the basis of sequence homology with nucleotide binding domains of other proteins, the NBDs of CFTR are considered to span shorter regions, from F433 through S589 in the case of the first NBD, and from Y1218 through R1386 in the case of the



**Fig. 1.** Hydropathy plot of the entire CFTR. Hydropathy plot calculated using the Kite-Doolittle algorithm. The double arrows point to the approximate positions of the transmembrane regions (TMD1 & TMD2), Nucleotide Binding Domains (NBD1 & NBD2) and Regulatory Domain (R domain) in the sequence.



**Fig. 2.** Schematic figure of a model for the entire CFTR. The five distinct domains are depicted. The critical  $\Delta F508$  mutation causing 70% of all cases of cystic fibrosis lies in the NBD1 domain.

second NBD [Riordan et al. (1989)]. These two nucleotide binding domains, which contain many of the sites that upon mutation give rise to CF, bind ATP and, in the case of NBD1, also hydrolyze ATP [Ko et al. (1995)]. The long exon 13 of the CF gene encodes a region of accentuated hydrophilicity bridging the N-terminal and C-terminal portions of CFTR. This region has a large number of alternating charged groups and most of the putative phosphorylation sites. Because phosphorylation by protein kinase A and protein kinase C activates the CFTR chloride channel, this region is called the "regulatory (R) domain" [Riordan et al. (1989)]. This domain makes CFTR unique within the ABC transporter family. All these observations are summarized in the simplified model of CFTR shown in the Figure 2.

### Alignment of the Nucleotide Binding Domains of CFTR with Those of other ABC Transporters

As described in Methods we modeled as NBDs 244/254 amino acid regions of CFTR. The boundaries of these regions were determined using the hydropathy profile and sequence alignment of NBD1 and NBD2. The beginning of NBD1 was set at residue L441 (18 amino acids before motif A); the end of NBD2 was set at the CFTR carboxy terminal, L1480. Alignment of NBD1 and NBD2 was used to define the end of

NBD1 at residue K684 and the beginning of NBD2 at L1227. Tables 1 & 2 show an automatic multiple sequence alignment of the two NBDs in CFTR, with non CFTR proteins found in the SwissProt database. Only proteins with more than 30% sequence identity with the target sequence were selected: this percentage is the threshold for structural similarity proposed by Sander et al. (1994). The two NBDs of CFTR have a high percentage of sequence identity with each other (33%), although, each NBD aligns with different regions of other ABC proteins (Tables 1 & 2). The alignment divides the NBDs of the ABC transporters into two classes (NBD1-like and NBD2-like), suggesting for each class a conserved distinctive function. Also, each alignment shows three common conserved regions: 1) a high identity initial region of 36 amino acids, which includes the motif A, positions 18–25 (Tables 1 and 2); 2) the motif C (positions 108–113/120–125, Table 1 and 2 respectively); and 3) the motif B (sequence positions range 120–132/132–144).

Alignments in Tables 1 and 2 also show a slightly different conserved stretch of at least eight residues (region D) after position 60, ending in an aromatic residue (F[508/1296] in CFTR). The sequences in each NBD class are ([T/S]h[K/R/Q][E/D]Nh<sub>2</sub>[F/Y]) in the NBD1-like class and Xh<sub>2</sub>XQ[D/E/Q/N/R/K]h<sub>3</sub>F in the NBD2-like class. Another conserved feature consist of a conserved glycine at position 40 in NBD1 (39 in NBD2) that is preceded or followed (NBD2 or NBD1) by a zone of high frequency of insertion. Also, a highly conserved motif is observed at positions 103–105 (underline in Table 1) or 115–117 (underline in Table 2): an acidic residue followed by a glycine (E/DXG) just before the motif C (underline in CFTR). No sequence identities are found after position 200, deep inside the exon 13 encoded segment. The alignment of a representative set of ABC transporters based on these considerations, shown in Tables 1 and 2, suggests that a model of both NBDs of CFTR can be generalized to other members of the ABC transporter family.

### Topology of Nucleotide Binding Domains

Analysis of known three-dimensional structures of proteins involved in nucleotide binding can be used to extract a set of conserved topological features. Certainly, a model for the NBDs must contain these features or have an explanation for their absence. Alignment of the secondary structural elements that form the nucleotide binding cassette (NBC) of a repre-

**Table I.** Automatic alignment [Rost (1996)] of 244 residues of NBD1 sequence of CFTR with 32 other sequences of non-CFTR proteins (listed below) found in the SwissProt sequence database. The amino acids at the beginning and the end of an insertion are indicated with lower case letters, the intervening amino acids are omitted

Code in Table 1	Identity %	No aa	No. of Ins. plus del.	Size of Ins.	Prot. Size	Swissprot Database Code	Description
cftr_human	100	244	0	0	1480	P13569	Dependent Chloride Channel (CFTR)
yhd5_yeast	40	223	4	19	1592	P38735	Probable ATP-Dependent Permease
mrp1_human	39	221	1	2	1531	P33527	Multidrug Resistance-Associate Protein
ycfi_yeast	38	226	2	8	1515	P39109	Metal Resistance Protein
mdr1_enthi	36	97	2	17	114	P16875	Multidrug Resistance Protein
yawb_schpo	34	244	2	10	1478	Q10185	Probable ATP-Dependent Permease
yfib_bacsu	34	212	4	22	573	P54718	Hypothetical ABC Transporter
mdr_leita	34	242	1	2	1548	P21441	Multidrug Resistance Protein
cydd_haein	34	204	6	31	586	P45082	Transport ATP-Binding Protein
sur_cricr	33	244	3	35	1581	Q09427	Sulfonylurea Receptor.
Yfic_bacsu	33	202	4	16	604	P54719	Hypothetical ABC Transporter
yor1_yeast	33	242	2	2	1477	P53049	Oligomycin Resistance ATP-Binding Protein
heta_anasp	33	212	4	22	607	P22638	Heterocyst Differentiation
sur_rat	32	244	3	39	1580	Q09429	Sulfonylurea Receptor
Sfuc_serma	32	204	6	40	345	P21410	Iron(III)-Transport ATP-Binding Protein
y015_mycge	32	199	4	19	589	P47261	Hypothetical ABC Transporter
yk83_yeast	32	240	4	19	1218	P36028	Probable ATP-Dependent Permease
cydc_ecoli	32	204	3	15	573	P23886	Transport ATP-Binding Protein
coma_strpn	32	201	4	19	717	Q03727	Transport ATP-Binding Protein
prtd_erwch	32	212	2	17	575	P23596	Proteases Secretion ATP-Binding Protein
y099_haein	32	219	6	33	356	P44513	Hypothetical ABC Transporter
nata_bacsu	31	200	7	27	246	P46903	ATP-Binding Transport Protein
mesd_leume	31	200	5	20	722	Q10418	Mesentericin Y105 Transporter
ybba_haein	31	178	5	23	227	P45247	Hypothetical ABC Transporter
mbpx_marpo	31	201	5	18	370	P10091	Probable Transport Protein
ndva_rhime	31	219	4	27	616	P18767	Beta-(1—>2)Glucan Export
uvra_serma	36	55	2	5	74	P25735	Excinuclease ABC Subunit
ywja_bacsu	30	201	3	15	575	P45861	Hypothetical ABC transporter
potg_ecoli	30	234	9	28	404	P31134	Putrescine Transport ATP-Binding Protein
nist_lacla	30	211	3	20	600	Q03203	Nisin Transport Atp-Binding Protein
mdl2_yeast	30	218	4	19	820	P33311	ATP-Dependent Permease MDR
hlyb_actac	30	201	4	19	707	P23702	Leukotoxin Secretion ATP-Binding Protein
spab_bacsu	30	218	3	20	614	P33116	Subtilin Transport ATP-Binding Protein

	1		-----A-----		50
cftr_human	LKDINFKIER	GQLLAVAGST	GAGKTSLLMM	IMGELEPSEG	KIKHSGRISF
yhd5_yeast	LCGLNIKFI	GKLNILIGST	GSGKSALLG	LLGELNLISG	SIIvtNSFAY
mrp1_human	LNGITFSIPE	GALVAVVGVQ	GCGKLSLLSA	LLAEMDKVEG	HVAIKGSVAY
ycfi_yeast	LKNINFQAKK	GNLTCIVGKV	GSGKTALLSC	MLGDRFRVKG	FATVHGVSAY
mdr1_enthi	.....	.....S	GCGKSTTIQL	IQRNYEPNGG	RVTLDgqIGL
yawb_schpo	LRDIDFVARR	GELCCIVGKV	GMGKSSLEA	CLGNMQKHSG	SVFRCGSIAY
yfib_bacsu	LRNVSFSAKP	RETIAILGAT	GSGKSTLFQL	IPRLYQPDGS	RIYi rRQIGY
mdr_leita	LRNVSLTIPK	GKLTMVI GST	GSGKSTLLGA	LMGEYSVESG	ELWAERSIAY
cydd_haein	.KPLNFQIPA	NHNVALVGQS	GAGKTSLMNV	ILGFL .PYEG	SLKINGhIAW
sur_cricr	LSNITIRIPR	GQLTMIVGQV	GCGKSLLLLA	TLGEMQKVSG	AVFwrGPVAY
yfic_bacsu	LKHLQFTVPA	GQSIAFVGPT	GAGKTTVTNL	LARFYEPNDG	KILIDgnMGF
yor1_yeast	FKDLNEDIKK	GEFIMITGPI	GTGKSSLLNA	MAGSMRKT DG	KVEVNGDLLM
heta_anasp	LNNITLTIER	GKTTALVGAS	GAGKTTLADL	IPRFYDPT EG	QILVDgkMAV
sur_rat	LSNITIRIPR	GQLTMIVGQV	GCGKSLLLLA	TLGEMQKVSG	AVFwrGPVAY
Sfuc_serma	LEHIDLQVAA	GSRTAIVGPS	GSGKTTLLRI	IAGFEIPDGG	QILLQggIGF
y015_mycge	LTGINFSVKH	GDVVAVIVGPT	GAGKSTIINL	LMKFYKPFEG	KIYmrEKISI
yk83_yeast	LKNISIDFKL	NSLNAIIGPT	GSGKSSLLLG	LLGELNLLSG	KIYvtNSMAY
cydc_ecoli	LKGISLQVNA	GEHIAILGRT	GCGKSTLLQQ	LTRAWDPQQG	EILl rQTISV

Table I. Continued.

coma_strpn	LSDINLTVPQ	GSKVAVFGIS	GSGKTTLAKM	MVNFYDPSQG	EISLGGyINY
prtd_erwch	LQNIHFSLQA	GETLVILGAS	GSGKSSLARL	LVGAQSPYQG	KVRLDgtIGY
y099_haein	LHDISFSLQR	GEILFLLGSS	GCGKTTLLRA	IAGFEQPSNG	EIWLKERLIF
nata_bacsu	VRDVSLTIEK	GEVVGILGEN	GAGKTTMLRM	IASLLEPSQG	VITVDggVLF
mesd_leume	IDDVSLTITA	GEKIALVGIS	GSGKSTLVKL	LVNFFQPESG	TISlrGHINY
ybba_haein	LKGVFSFMEP	AELVAIVGSS	GSGKSTLLHT	LGGLDQPSSG	EVFINGyLGF
mbpx_marpo	LDRVSLYVPK	FSLIALLGPS	GSGKSSLLRI	IAGLDNCDYG	NIWLHgrMSF
ndva_rhime	VRNVSFKAKA	GQTIAIVGPT	GAGKTTLVNL	LQRVHEPKHG	QILIDgsIAT
uvra_serma	LKNINLIIPR	DKLIVVTGLS	GSGKSSLAFD	TLY...AEG	QRRYVESLsy
ywja_bacsu	LNDINLSIQA	GETVAVFGPS	GAGKSTLCSL	LPRFYEASEG	DITirGQIGV
potg_ecoli	VDDVSLTIYK	GEIFALLGAS	GCGKSTLLRM	LAGFEQPSAG	QIMLDgpINM
nist_lacla	LKNINLSFEK	GELTAIVGKN	GSGKSTLVKI	ISGLYQPTMG	I IQyqKNISV
mdl2_yeast	FKNLNFKIAP	GSSVCIVGPS	GRGKSTIALL	LLRYYNPTTG	TITirRHIGI
hlyb_actac	LNNINLDISQ	GEVIGIVGRS	GSGKSTLTKL	IQRFYIPEQG	QVLIDggVGV
spab_bacsu	LKHINVS LHK	GERVAIVGPN	GSGKKTFIKL	LTGLYEVHEG	DILINgqIAA
	51	---D---			100
cftr_human	CSQFSWIMPG	TIKENIIFGV	SYDEYRYSV	IKACQLEEDI	SKFAEKDNIV
yhd5_yeast	CSQSAWLLND	TVKNNIIFDN	FYNEDRYNKV	IDACGLKRDL	EILPAGDLTE
mrpl_human	VPQQAWIQND	SLRENILFGC	QLEEPYRSV	IQACALLPDL	EILPSGDRTE
ycfi_yeast	VSQVPWIMNG	TVKENILFGH	RYDAEFYEKT	IKACALTIDL	AILMDGDKTL
mdr1_enth	VGGEPVLFAG	TIRENIMLGA	KEGETLSKDE	MIEcNAHEFV	SKLAEGYDTL
yawb_schpo	AAQQPWILNA	TIQENILFGL	ELDPEFYEKT	IRACCLLRDF	EILADGDQTE
yfib_bacsu	VPQEVLLFSG	TIKENIAWGK	ENASleIMDA	AKLAQIHETI	LKLPNGYDTV
mdr_leita	VPQQA WIMNA	TLRGNILFFD	EERAE DLQDV	IRCCQLEADL	AQFCGGLDTE
cydd_haein	VGQNPLLLQG	TIKENLLLG.	. . DVQANDEE	INQALMRSQA	KEFTDKLGLh
sur_cricr	ASQKPWLLNA	TVEENITFES	PFNKQRYKMV	IEACSLQPDI	DILPHGDQQT
yfic_bacsu	VLQDSFLFQG	TIRENIRYGR	LDASDqvEAA	AKTANAHSFI	ERLPGKYDTV
yori1_yeast	CG.YPWIQNA	SVRDNIIFGS	PFNKKEYDEV	VRVCSLKADL	DILPAGDMTE
heta_anasp	VSQDTFIFNT	SIRDNIAYGT	saSEAEIREV	ARLANALQFI	EEMPEGFDTK
sur_rat	ASQKPWLLNA	TVEENITFES	PFNKQRYKMV	IEACSLQPDI	DILPHGDQQT
sfuc_serma	VPQDGALEpf	TVAGNIGFGL	KGGKREKQRR	IEA...L	MEMVALDRRL
y015_mycge	VLQDSFLFSG	TIKENIRLGR	. . QDATDDEI	IAACKTadFI	MRLPKGYDTY
yk83_yeast	CSQTPWLISG	TIKDNVVGE	IFNKQKFDDV	MKSCCLDKDI	KAMTAGIRTD
cydc_ecoli	VPQRVHLSA	TLRDNLLAS	PGSSDEALSE	ILRRVGLEKL	LEDAGL.NSW
coma_strpn	LPQQPYVFNG	TILENLLGA	KEGTTqILRA	VELAIREDI	ERMPLNYQTE
prtd_erwch	LPQDVQLFKG	SLAENIARFG	DADPEKVVA	AKLAGVHELI	LSLPNGYDTE
y099_haein	GENFn1FPHL	NVYRNIA YGL	GNGKGkeKTR	IEQIMQLTGI	FELADR...
nata_bacsu	GGETGLYDRM	TAKENLQYFG	RLYGLN.RHE	IKA.RIEDLS	KRFGMRDYMN
mesd_leume	LPQEPFIFSG	SIMENLLGA	KPGTTQ.EDI	IRAVEIAedI	EKMSQGFGE
ybba_haein	VYQFHLLMAd	tALENVMMPM	L...IGHQNK	TEAKDRAEKM	.LSAVGLSHR
mbpx_marpo	VFQHYALFkm	TVYENISFGL	RLRGFSAQKI	TNKVNDLLNC	LRIAD...I
ndva_rhime	VFQDAGLMNR	SIGENIRLGR	EDASleVMAA	AEAAAASDFI	EDRLNGYDTV
uvra_serma	ARQFLSLME.	.....	.....	.....	.....
ywja_bacsu	VQQDVFLFSG	TLRENIAYGR	laSEEDIWQA	VKQAHLLELV	HNMPDGLDTM
potg_ecoli	MFQSYALFpm	TVEQNIAFGL	KQDKLpiASR	VNEMLGLVHM	QEFAKR...
nist_lacla	LFQDFVKYEL	TIRENIGLSD	LSSQWEDEKI	IKVLDlKTNN	QYVLDTQLGN
mdl2_yeast	VQQEPVLMMSG	TIRDNITYGL	TYTPTkiRSV	AKQCFCHNFI	TKFPNTYDTV
hlyb_actac	VLQDNVLLNR	SIRENIALTN	. . PCMPMEKV	IAAAKLadFI	SELREGYNTV
spab_bacsu	LFQDFMKYEM	TLKENIGFGQ	IDKLHqvLDI	VRADFLKSHS	SYQFDTQLGL
	101	-----C-----	-----B-----		150
cftr_human	LGEGGITLSG	GQRARISLAR	AVYKDADLYL	LDSFGYLDV	LTEKEIFESC
yhd5_yeast	IGEKGITLSG	GQKQRISLAR	AVYSSAKHVL	LDDCLSAVDS	HTAVWIYENC
mrpl_human	IGEKGVNLSG	GQKQRVSLAR	AVYSNADIYL	FDDPLSAVDA	HVGKHFENV
ycfi_yeast	VGEKGISLSG	GQKARLSLAR	AVYARADTYL	LDDPLAAVDE	HVARHLIEHV
mdr1_enth	IGEKGALLSG	GQRQRI....	.....	.....	.....
yawb_schpo	VGEKGISLSG	GQKARISLAR	AVYSRSDIYL	LDDILSAVDD	HVNRDLVRNL
yfib_bacsu	LGQRGVNLSG	GQKQRISLAR	ALIRKPAILL	LDDSTSALDL	QTEAKLLEAI
mdr_leita	IGEMGVNLSG	GQKARVSLAR	AVYANRDVYL	LDDPLSALDA	HVGQRIVQDV

Table I. Continued.

cydd_haein	iKDGGGLGISV	GQAQLAIAR	ALLRKGDLLL	LDEPTASLDA	QSENLVLQA.
sur_cricr	IGERGINLSG	GQRQRISVAR	ALYQQTNVVF	LDDPFSALDV	HLSDHLMQAG
yfic_bacsu	LTQNGSGISQ	GQKQLISAR	AVLADPVLLI	LDEATSNIDT	VTEVNIQEA.
yorl_yeast	IGERGITLSG	GQKARINLAR	SVYKKKDIYL	FDDVLSAVDS	RVGKHIMDEC
heta_anasp	LGDRGVRLSG	GQRQRIAIAR	ALLRDPEILI	LDEATSALDS	VSERLIQES.
sur_rat	IGERGINLSG	GQRPGISVAR	ALYQHTNVVF	LDDPFSALDV	HLSDHLMQAG
sfuc_serma	AALWPHELSG	GQQQRVALAR	ALSQQPRLML	LDEPFSALDT	GLRAATRKA
y015_mycge	ISNKADYLSV	GERQLLTAR	AVIRNAPVLL	LDEATSSVDV	HSEKLIQES.
yk83_yeast	VGDGGFSLSG	GQQQRIALAR	AIYSSSRYL	LDDCLSVDV	ETALYIYEEC
cydc_ecoli	LGEGGRQLSG	GELRRLAIAR	ALLHDAPLVL	LDEPTEGLDA	TTESQILEL.
coma_strpn	LTSDGAGISG	GQRQRIALAR	ALLTDAPVLI	LDEATSSLDI	LTEKRIVDN.
prtd_erwch	LGDDGGGLSG	GQRQRIGLAR	AMYGDPCLLI	LDEPNASLDS	EGDQALMQAI
y099_haein	..FPHQLSG	GQQQRVALAR	ALAPNPELIL	LDEPFSALDE	HLRQQIRQEM
nata_bacsu	RRVGG..FSK	GMRQKVAIAR	ALIHDPDIIL	FDEPTTGLDI	.TSSNIFREF
mesd_leume	LAESG.NISG	GQKQRIALAR	AILVDSFVLI	LDESTSNLDV	LTEKKIIDNL
ybba_haein	ITHRPSALSG	GERQVAIAR	ALVNNPSLVL	ADEPTGNLDH	KTESIFELI
mbpx_marpo	SFEYPAQLSG	GQKQRVALAR	SLAIQPDFLL	LDEPFGALDG	ELRRHLSKWL
ndva_rhime	VGERGNRLSG	GERQVAIAR	AILKNAPILV	LDEATSALDV	ETEARKVDA.
uvra_serma	.....	.....	.....	.....	.....
ywja_bacsu	IGERGVKLSG	GQKQLSIAR	MFLKNPSILI	LDEATSALDT	ETEAAIQKA.
potg_ecoli	..KPHQLSG	GQRQRVALAR	SLAKRPKLL	LDEPMGALDK	KLRDRMQLEV
nist_lacla	WFQEGHQLSG	GQWQKIALAR	TFFKKASIYI	LDEPSAALDP	VAEKEIFDYF
mdl2_yeast	IGPHGTLLSG	GQKQRIAIAR	ALIKKPTILI	LDEATSALDV	ESEGAINYT.
hlyb_actac	VGEQGAGLSG	GQRQRIAIAR	ALVNNPRILI	FDEATSALDY	ESENIIMHN.
spab_bacsu	WFDEGRQLSG	GQWQKIALAR	AYFREASLYI	LDEPSSALDP	IAEKETFDTF
	151				200
cftr_human	VCKLMANKTR	ILVTSKMEHL	KKADKILILH	EGSSYFYGTF	SELQNLQPDF
yhd5_yeast	ItpLMKNRTC	ILVTHNVstL	RNAHFAIVLE	NGKVKNQGTI	TELQS.KGLF
mrpl_human	IggMLKNKTR	ILVTHMSVYL	PQVDVIVMS	GKISEMGSY	QELLARDGAF
yfi_yeast	LggLLHTTKTK	VLATNKVSAL	SIADSIALLD	NGEITQQGTY	DEITKDADSP
mdrl_enth	.....	.....	.....	.....	.....
yawb_schpo	LggLLRSRCV	ILSTNSLTVL	KEASMIYMLR	NGKIIESGSF	TQLSsllSEF
yfib_bacsu	S...TYHCTT	LIITQKITTA	MKADQILLE	DGELIEKGTH	SELLS....
mdr_leita	ILGRLRGKTR	VLATHQIHL	PLADYIVVLQ	HGSIVFAGDF	AAFSaLEETL
cydd_haein	LNESASQHQT	LMITHRIEDL	KQCDQIFVMQ	RGEIVQQGKF	TELQHE...W
sur_cricr	ILELldKRTV	VLVTHKLQYL	PHADWIIAMK	DGTIQREGTL	KDFQRSECQL
yfic_bacsu	LARLMEGRTS	VIIAHLRNTI	QRADQIVVLK	NGEMIEKGS	DELIR.QKGF
yorl_yeast	LTGMLANKTR	ILATHQLSLI	ERASRVIVLG	TDGQVDIGTV	DELKARNQTL
heta_anasp	IEKLSVGRTV	IAIAHLRSTI	AKADKVVVME	QGRIVEEQNY	QELLEQRGKL
sur_rat	ILELldKRTV	VLVTHKLQYL	PHADWIIAMK	DGTIQREGTL	KDFQRSECQL
sfuc_serma	AELLtAKVAS	ILVThqSEAL	SFADQVAVMR	SGRLAQVGAP	QDL.....
y015_mycge	IGRLMKNKTS	FIISHRLSII	RDATLIMVIN	DGKVLMEGNH	DQLMKQNGFY
yk83_yeast	LcpMMKGRTC	IITSHNISLv	KRADWLIVLD	RGEVKSQGKP	SDLIKSN.EF
cydc_ecoli	LAEMMREKTV	LMVTHRLRGL	SFRFQIIVMD	NGQIEEGQTH	AELLARQGRY
coma_strpn	..LIALDKTL	IFIAHRLTIA	ERTEKVVVLD	QGKIVEEGKH	ADLLA.QGGF
prtd_erwch	VALQKRGATV	VLITHRPALT	TLAQKILILH	EGQQQRMG11	TELQQRSAAN
y099_haein	LQALrsGASA	IFVThrDESL	RYADKIAIIQ	QGKILQIDTP	RTLY...W
nata_bacsu	IQQLKREQKT	ILFSSHieVQ	ALCDSVIMIH	SGEVIYRGAL	ESLYESErIF
mesd_leume	M..KLTEKTI	IFVAHLRTIS	QRVDRILTMQ	SGKIIEDGTH	DTLLKAGGFY
ybba_haein	QqnQEQNIAF	LLVTHDMGLA	EKLSRRLVMQ	DG.....	.....
mbpx_marpo	KRYLQDNktT	IMVTHDqeAI	SMADEVILK	EGRLLQQGKP	KNLYDQPINF
ndva_rhime	IDALRKDRTT	FIIAHLRSTV	READLVIFMD	QGRVVMEMGF	HELSSQNGRF
uvra_serma	.....	.....	.....	.....	.....
ywja_bacsu	LQELSEGRTT	LVI AHLRATI	KDADRIVVVT	NNGIEEQGRH	QDLIEAGGLY
potg_ecoli	VDILevGVTC	VMVTHDQeam	TMAGRIAIMN	RGKFVQIGEP	EIIEHPtrY
nist_lacla	V.ALSENNIS	IFISHSLNAA	RKANKIVVMK	DGQVEDVGSH	DVLLRRCQYY
mdl2_yeast	FGQLMKSKSM	TIVSIAHRLr	RSENVIVLGH	DGSVVMEMGF	KELANPTSA
hlyb_actac	MHKICQNRTV	LI AHLRSTV	KNADRIIVMD	KGEIEEQGKH	QELLKDEKGL
spab_bacsu	F.SLSKDKIG	IFISHRLVAA	KLADRIVVMD	KGEIVGIGTH	EELLKTCPY

Table I. Continued.

	201				244
cftr_human	SSKLMGCDSF	DQFSAERRNS	ILTETLHRFS	LEGDAPVSWT	ETKK
yhd5_yeast	KEKYVQLSSR	DSINEKNANR	LKAP.....	.....	....
mrpl_human	AEFLRTYAST	EQEQDAEENG	V.....	.....	....
ycfi_yeast	LWLLNngKS	NEFGDSESS	VRESSI....	.....	....
mdrl_enth1	.....	.....	.....	.....	....
yawb_schpo	SKKDTASSTG	ADTPLRSRQS	VITSSTDVTS	SASRSSDTVS	NYPK
yfib_bacsu	ESQLYKRIYE	SQFGREGSES	.....	.....	....
mdr_leita	RGELKGSKDV	ESCSSDVDTE	SATAETAPYV	AKAKGLNAEQ	ET..
cydd_haein	.....GFF	AELLAQRQQD	I.....	.....	....
sur_cricr	FEHWKTLMNR	QDQemERKAS	EPSQGLPRAM	SSRDGLLDE	EEEE
yfic_bacsu	YSDL.....	.....	.....	.....	....
yor1_yeast	INLLQFSSQN	SEKEDEEQEA	VVAGELGQLK	YESEVK.ELT	ELKK
heta_anasp	W.....KY	HQMQUESGQT	.....	.....	....
sur_rat	FeeLEKETVM	ERKAPEPSQG	LPRAMSSRDG	LLLDEDEEEE	EAAE
sfuc_serma	.....Y	LRPVDEPTAS	FLGETL....	.....	....
y015_mycge	AR.....	.....	.....	.....	....
yk83_yeast	LRESINNSDK	NTTHNQIDLK	RSTTSKKTKN	GDPEGGNSQD	E...
cydc_ecoli	YQFKQG....	.....	.....	.....	....
coma_strpn	YHLV.....	.....	.....	.....	....
prtd_erwch	QARMNPTAAM	PQ.....	.....	.....	....
y099_haein	SPNHLETAKF	MGESIVLPAN	LLDENTAQCQ	L.....	....
nata_bacsu	MSKLV.....	.....	.....	.....	....
mesd_leume	ASLF.....	.....	.....	.....	....
ybba_haein	.....	.....	.....	.....	....
mbpx_marpo	FVGIF.....	.....	.....	.....	....
ndva_rhime	AALL.....	.....RASG	ILTDEDVRKS	LT.....	....
uvra_serma	.....	.....	.....	.....	....
ywja_bacsu	SR.....	.....	.....	.....	....
potg_ecoli	SAEFIGsnVF	EGVLKERqgL	VLDSPLVHP	LKVDADASVV	D...
nist_lacla	QELYSEQYE	DN.....	.....	.....	....
mdl2_yeast	LSQLLNEKAA	PGPSDQQLQ.	.....	.....	....
hlyb_actac	YSYL.....	.....	.....	.....	....
spab_bacsu	KKMDESENYM	NPLEEESGSK.	.....	.....	....

sentative set of nucleotide binding proteins (Table 3) shows a strong degree of conservation of topological motifs. With the increased number of reported structures of NB proteins it is now possible to cluster these structures in at least three topologically related types: (1) mitochondrial  $F_1$ -ATPase/recA, (2) ADK and (3) p21/G-proteins (Figures 3a,b,c). They all share an  $\alpha/\beta$  motif (known as the Rossman fold; [Rossman et al. (1974)]), consisting of 5 to 7, predominantly parallel  $\beta$ -strands, surrounded by 3 to 7 helices. The topology of these proteins (Figure 3a,b,c) shows strong conservation of a central motif, composed of a  $\beta$ -strand followed by an  $\alpha$ -helix, a  $\beta$ -strand that can be far away in the sequence, and an additional  $\alpha$ -helix and  $\beta$ -strand ( $\beta 1\alpha A \dots \beta 4 \dots \alpha F\beta 6$ , boxed in Figure 3a). Topological insertions or changes are concentrated in the regions between helix  $\alpha A$  and strand  $\beta 4$ , and between  $\beta 4$  and  $\alpha F$ . This conserved core motif contains two

important consensus sequences  $GX_4GK[T/S]$ , motif A (also known as the "P-loop"), and  $RX_{6-8}h_4D$ , motif B. The P-loop forms the binding site for triphosphate (or diphosphate) nucleotides and for  $Mg^{++}$ . Strand  $\beta 1$ , loop ( $GX_4G$ ) and helix  $\alpha A$  (Table 3 and Figures 3) define the motif containing the P-loop.

Motif B is located in strand  $\beta 4$  and ends with a conserved negative charged residue (aspartate, D). The third consensus signature, motif C, is not equally general; for example it is not present in the  $F_1$ /recA NBD. In the case of motifs A and B, not only the secondary structure is conserved, but also their spatial arrangement (Figure 4). Clearly, any structural model of a protein containing a nucleotide binding cassette must account for the facts described above.

We have chosen a member of the first topological class (Figure 3a) as a template for modeling NBDs.  $F_1$ -ATPase is a five subunit protein with stoichiometry



**Table II.** Automatic alignment [Rost (1996)] of 250 residues of NBD2 sequence of CFTR with 25 other sequences of non-CFTR proteins (listed below) found in the SwissProt sequence database. The amino acids at the beginning and the end of an insertion are indicated with lower case letters, the intervening amino acids are omitted

Code in Table 1	Identity %	No aa	No. of Ins. plus del.	Size of Ins.	Prot. Size	Swissprot Database Code	Probable Transport ATP-Binding Protein
cftr_human	100	250	0	0	1480	P13569	Dependent Chloride Channel (CFTR)
mrpl_human	38	221	1	1	1531	P33527	Multidrug Resistance-Associate
yawb_schpo	38	221	1	1	1478	Q10185	Probable ATP-Dependent Permease
ycfi_yeast	37	222	3	3	1515	P39109	Metal Resistance Protein
sur_human	36	220	1	1	395	Q09428	Sulfonylurea Receptor
yorl_yeast	36	220	4	23	1477	P53049	Oligomycin Resistance ATP
sur_rat	36	220	1	1	1580	Q09429	Sulfonylurea Receptor.
sur_cricr	35	220	1	1	1581	Q09427	Sulfonylurea Receptor.
yk84_yeast	35	230	3	22	306	P36171	Probable ATP-Dependent Transporter
yhd5_yeast	34	231	4	19	1592	P38735	Probable ATP-Dependent Permease
msba_haein	33	216	2	3	587	P44407	Probable Transport ATP-Binding Protein
lcnc_lacla	33	211	3	5	715	Q00564	Lactococcin A Transport
yfib_bacsu	32	225	4	9	573	P54718	Hypothetical ABC Transporter
ste6_yeast	32	217	3	5	1290	P12866	Protein Homolog) (P-Glycoprotein)
cydc_ecoli	32	211	3	4	573	P23886	Transport ATP-Binding Protein
ywja_bacsu	32	218	2	2	575	P45861	Hypothetical ABC Transporter
cydd_haein	32	203	3	7	586	P45082	Transport ATP-Binding Protein
mdr_leita	32	241	3	10	1548	P21441	Multidrug Resistance Protein
cyab_borpe	31	213	2	2	712	P18770	Cyclolysin Secretion ATP-Binding
msba_ecoli	31	213	2	3	582	P27299	Probable Transport ATP-Binding
heta_anasp	31	213	2	2	607	P22638	Heterocyst Differentiation
nist_lacla	31	219	3	7	600	Q03203	Nisin Transport ATP-Binding
cydc_haein	30	218	2	3	576	P45081	Transport ATP-Binding Protein
yabj_haein	30	182	7	25	229	P44986	Hypothetical ABC Transporter
hst6_canal	30	213	3	51	1323	P53706	ATP-Dependent Permease HS
yfic_bacsu	30	218	2	2	604	P54719	Hypothetical ABC Transporter

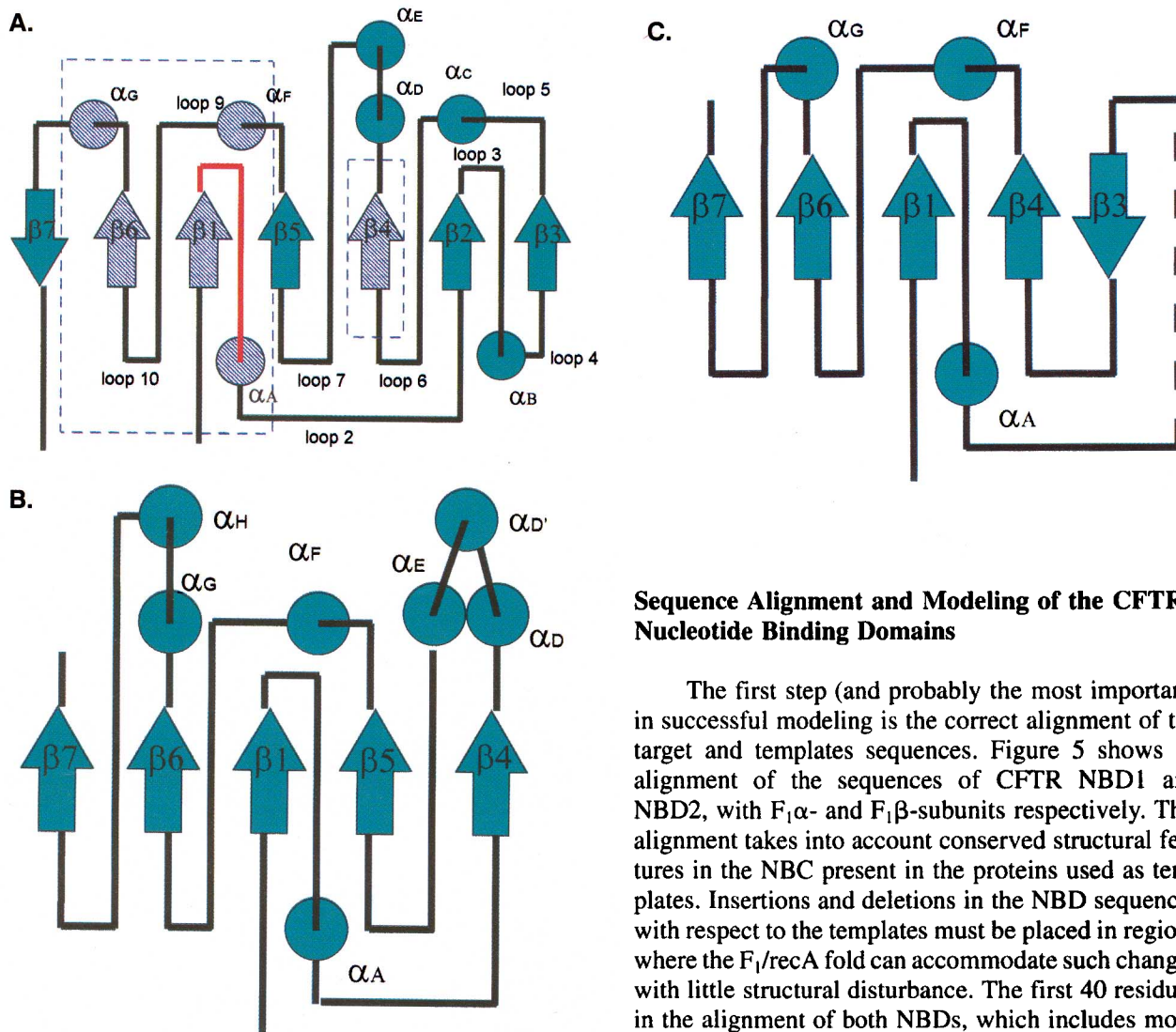
	1	-----A-----		50	
cftr_human	LENISFSISP	QQRVGLLGR	GSGKSTLLSA	FLRLLNTEGE	IQIDGVSWDS
mrpl_human	LRHINVTING	GEKVGIVGR	GAGKSSLTLG	LFRINEseGE	IIDGINIAK
yawb_schpo	LNDISVNIKP	QEKIGIVGR	GAGKSTLTLA	LFRLIetSGD	QLDDINITS
ycfi_yeast	LKHINIHIKP	NEKVGIVGR	GAGKSSLTLA	LFRMIEaeGN	IVIDNIAINE
sur_human	LKHVNALISP	GQKIGICGR	GSGKSSFSLA	FFRMVDteGH	IIDGIDIK
yorl_yeast	LKNLNLNIKS	GEKIGICGR	GAGKSTIMSA	LYRLNetAGK	ILIDNVDISQ
sur_rat	LKHVNALISP	GQKIGICGR	GSGKSSFSLA	FFRMVDmeGR	IIDGIDIK
sur_cricr	LKHVNTLISP	GQKIGICGR	GSGKSSFSLA	FFRMVDmeGR	IIDGIDIK
yk84_yeast	LDNVSPFKVKA	GTKVGIVGR	GAGKSSIAAA	IYRLSDwnGT	ITIDNKDIKH
yhd5_yeast	IRNVSPKVDP	QSKIGIVGR	GAGKSTIITA	LFRLLEptGC	IKIDGQDISK
msba_haein	LNNISFSVPA	GKTVALVGRS	GSGKSTIANL	VTRFYDiqGE	ILLDGVNIQD
lcnc_lacla	LSEIELSIKE	NEKLTIVGMS	GSGKSTLVKL	LVNFFqtSGT	ITLGGIDLQQ
yfib_bacsu	LRNVSFSAKP	RETIAILGAT	GSGKSTLFQL	IPRLYQpsGR	IYIDEKPVQD
ste6_yeast	.KNMNFDMFC	GQTLGIGES	GTGKSTLVLL	LTKLYNcGgK	IKIDGTDVND
cydc_ecoli	LKGISLQVNA	GEHIAILGR	GCGKSTLLQQ	LTRAWDppGE	ILLNDSPIAS
ywja_bacsu	LNDINLSIQA	GETVAFVGPS	GAGKSTLCSL	LPRFYEaeGD	ITIDGISIKD
cydd_haein	. . .LNFQIPA	NHNVALVGQS	GAGKSTLMNV	ILGFLPYEGS	LKINGQELRE
mdr_leita	LRGVSFQIAP	REKVGIVGR	GSGKSTLLLT	FMRMVEvgGV	IHVNGREMSA
cyab_borpe	LRNVSLRIAP	GEVGVVGRS	GSGKSTLTRL	IQRMFVarGR	VLIDGHDIGI
msba_ecoli	LRNINLKIPA	GKTVALVGRS	GSGKSTIASL	ITRFYDieGE	ILMDGHDLRE
heta_anasp	LNNITLTIER	GKTTALVGAS	GAGKTTLADL	IPRFYdtEGQ	ILVDGLDVQY
nist_lacla	LKNINLSFEK	GELTAIVGKN	GSGKSTLVKI	ISGLYqtMGI	IQYDKMRSSL
cydc_haein	LKNLTLDEQ	GKKIAILGKT	GSGKSSLLQL	LVRNYDaqGE	LLLAEKPIISA
yabj_haein	. . .NLSVNA	GERVAIIGES	GAGKSTLLNL	IAGFEFpqGE	IWLNDKNH . .

Table II. Continued.

hst6_canal	LKSISLDVVK	FTTIGIVGQS	GSGKSTILKI	LFRLYDIdQT	VKIFNQNLyL
yfic_bacsu	LKHLQFTVPA	GQSIaFVGPT	GAGKTTVTNL	LARFYEpDGK	ILIDGTDIKT
	51	---D---			100
cftr_human	ITLQQWRKAF	GVIPQKVFI	SGTFRKNLDP	YEQWSDQEIW	KVADEVGLRS
mrp1_human	IGLHDLRFKI	TIIPQDPVLF	SGSLRMNLDP	FSQYSDEEVW	TSLELAHLKD
yawb_schpo	IGLHDLRSRL	AIIPQENQAF	EGTIRENLDP	NANATDEEIW	HALEAASLKQ
ycfi_yeast	IGLYDLRHKL	SIIPQDSQVF	EGTVRENIDP	INQYTDEAIW	RALESHLKe
sur_human	LPLHTLRSRL	SIILQDPVLF	SGTIRFNLDP	ERKCSDSLW	EALEIAQLKL
yor1_yeast	LGLFDLRRKL	AIIPQDPVLF	RGTIRKNLDP	FNERTDDELw	qKPDENGTHG
sur_rat	LPLHTLRSRL	SIILQDPVLF	SGTIRFNLDP	EKKCSDSLW	EALEIAQLKL
sur_cricr	LPLHTLRSRL	SIILQDPVLF	SGTIRFNLDP	EKKCSDSLW	EALEIAQLKL
yk84_yeast	IPLERLRNSI	SCIPQDPTLF	DGTVRSNLDP	FDRYSDVQIY	GVLSKVGLIe
yhd5_yeast	IDLVTLRRSI	TIIPQDPILF	AGTIKSNVDP	YDEYDEKKIF	KALSQVNLIS
msba_haein	YRLSNLRENC	AVVSQQVHLF	NDTIANNIAY	AAqySREEII	AAAKAAyALE
lcnc_lacla	FDKHLRRLI	NYLPQQPYIF	TGSILDNLLI	nENASQEEIL	KAVELAEIRA
yfib_bacsu	IPAEGLRRQI	GYVPQEVLLF	SGTIKENIAw	kENASLDEIM	DAAKLAQIHE
ste6_yeast	WNLTSLRKEI	SVVEQKPLLF	NGTIRDNLTy	qDEILEIEMY	DALKYVGIHD
cydc_ecoli	LNEAALRQTI	SVVPQRVHLF	SATLRDNLII	sPGSSDEALS	EILRRVGLK
ywja_bacsu	MTLSSLRGQI	GVVQQDPVLF	SGTLRENIAY	GRlaSEEDIW	QAVKQAHLLE
cydd_haein	SNLADWRKHI	AWVGQNPILL	QGTIKENLLL	GdqANDEEI.	...NQALMRS
mdr_leita	YGLRDVRRHF	SMIPQDPVLF	DGTVRQNVDP	FLEASSAEVW	AALELVGLRE
cyab_borpe	VDSASLRRQL	GVVLQESTLF	NRSVRDNIAI	rPGASMHEVV	AAARLAGAHE
msba_ecoli	YTLASLRNQV	ALVSQNVHLF	NDTVANNIAY	aeQYSREQIE	EAARMAYAMD
heta_anasp	FEINSLRRKM	AVVSQDTFIF	NTSIRDNIAY	GTsaSEAEIR	EVARLANALQ
nist_lacla	MPEEFYQKNI	SVLFQDFVKY	ELTIRENIGL	ssQWEDEKII	KVLDNLGLDF
cydc_haein	YSEETLRHQI	CFLTQRVHVF	SDTLRQNLQF	adKISDEQMI	EMLHQVGLSK
yabj_haein	TRSAPYERPV	SMLFQENNLf	pITVQQNLAI	ITALEQEKEIE	QVACSVGLGD
hst6_canal	INSGLLCQTI	AIVPQFPKFF	SGTIYDNLTy	sSSVSDSEII	KILKLVNLHQ
yfic_bacsu	LTRASLRKNM	GFVLQDSFLF	QGTIRENIRy	rLDASDQEVE	AAAKTANAHS
	101	--C-----		-----B-----	150
cftr_human	VIEQFP GKLD	FVLVDGGCVL	SHGHKQLMCL	ARSVLSKAKI	LLLDEPSAHL
mrp1_human	FVSALPDKLD	HECAEGGENL	SVGQRQLVCL	ARALLRRTKI	LVLDEATAAV
yawb_schpo	FIQTLDGGLY	SRVTEGGANL	SSGQRQLMCL	TRALLTPTRV	LLLDEATAAV
ycfi_yeast	vLSMSNDGLD	AQLTEGGGNL	SVGQRQLLCL	ARAMLVPSKI	LVLDEATAAV
sur_human	VVKALPGGLD	AIITEGGENF	SQGQRQLFCL	ARAFVRKTSI	FIMDEATASI
yor1_yeast	KMHKF..HLD	QAVEEEGSNF	SLGERQLLAL	TRALVRQSKI	LILDEATSSV
sur_rat	VVKALPGGLD	AIITEGGENF	SQGQRQLFCL	ARAFVRKTSI	FIMDEATASI
sur_cricr	VVKALPGGLD	AIITEGGENF	SQGQRQLFCL	ARAFVRKTSI	FIMDEATASI
yk84_yeast	kLRNRFIDLN	TVVKSGGSNL	SQGQRQLLCL	ARSM LGARNI	MLIDEATASI
yhd5_yeast	SHEFEEvLH	TEIAEGLNL	SQGERQLLFI	ARSLLRPKI	ILLDEATSSI
msba_haein	FIEKLPQVFD	TVIGENGTSL	SGGQRQLLAI	ARALLRNSPV	LILDEATSAL
lcnc_lacla	DIEQMQLGYQ	TESSDASSL	SGGQRQLIAL	ARALLSPAKI	LILDEATSAL
yfib_bacsu	TILKLPNGYD	TVLGQRGVNL	SGGQKQRISI	ARALIRKPAI	LLLDDTSAL
ste6_yeast	FVISSPQGLD	TRI..DTLL	SGGQAQRICI	ARALLRKSki	LILDECTSAL
cydc_ecoli	LLE..DAGLN	SWLGEGRQL	SGGELRRLAI	ARALLHDAPL	VLLDEPTEGL
ywja_bacsu	LVHNMPPGLD	TMIGERGVLK	SGGQKQRISI	ARMFLKNPSI	LILDEATSAL
cydd_haein	QAKEFTDKLG	leIKDGGLGI	SVGQAQR LAI	ARALLRKGDL	LLLDEPTASL
mdr_leita	RVASESEGID	SRVLEGSNY	SVGQRQLMCM	ARALLKrsGF	ILMDEATANI
cyab_borpe	FICQLPEGYD	TMLGENGVGL	SGGQRQRIGI	ARALIHRPRV	LILDEATSAL
msba_ecoli	FINKMDNGLD	TVIGENGVLL	SGGQRQR IAI	ARALLRDSPI	LILDEATSAL
heta_anasp	FIEEMPEGFD	TKLGDRGVRL	SGGQRQR IAI	ARALLRDP EI	LILDEATSAL
nist_lacla	LKTNNQYVLD	TQlfQEGHQL	SGGQWQKIAL	ARTFFK KASI	YILDEPSAAL
cydc_haein	LLEQEGKGLN	LWLGDGGRPL	SGGEQRRLGL	ARILLNNASI	LLLDEPTEGL
yabj_haein	YLERLP...	.....NSL	SGGQKQRVAL	ARCLLRDKPI	LLLDEPFSAL
hst6_canal	FIVSLPQGLL	TIMNDSntF	SGGQLLLAI	ARALLRNPKI	LLLDECTSAL
yfic_bacsu	FIERLPKGYD	TVLTQNGSGI	SGGQKQLISI	ARAVLADPVL	LILDEATSNI

Table II. Continued.

	151				200
cftr_human	DPVTYQIRR	TLKQAFADCT	VILCEHRIEA	MLECQQFLVI	EENKVRQYDS
mrpl_human	DLETDDLQSQ	TIRTQFEDCT	VLTIAHRLNT	IMDYTRVIVL	DKGEIQEYGA
yawb_schpo	DVETDAIVQR	TIRERFNDRT	ILTIAHRINT	VMDSNRILVL	DHGKVVVEFDS
ycfi_yeast	DVETDKVVQE	TIRTAFKDRT	ILTIAHRLNT	IMDSDRIVL	DNGKVAEFDS
sur_human	DMATENILQK	VVMTAFADRT	VVTIAHRVHT	ILSADLVIVL	KRGAILEFDK
yorl_yeast	DYETDGKIQT	RIVEEFGDCT	ILCIAHRLKT	IVNYDRILVL	EKGEVAEFDT
sur_rat	DMATENILQK	VVMTAFADRT	VVTIAHRVHT	ILSADLVMVL	KRGAILEFDK
sur_cricr	DMATENILQK	VVMTAFADRT	VVTIAHRVHT	ILSADLVMVL	KRGAILEFDK
yk84_yeast	DYISDAKIQQ	TIRETMKNTT	ILTIAHRLRS	VIDYDKILVM	EMGRVKEYDH
yhd5_yeast	DYDSDHLIQG	IIRSEFNKST	ILTIAHRLRS	VIDYDRIVM	DAGEVKEYDR
msba_haein	DTESERAIQS	ALEELKKDRT	VVVIARHLS	IENADEILVI	DHGEIRERGN
lcnc_lacla	DMITEKKILK	NLLP.LDKT	IIFIAHRLSV	AEMSHRIIV	DQGVVIESGS
yfib_bacsu	DLQTEAKLLE	AIST.YHCT	TLIITQKITT	AMKADQILL	EDGELIEKGT
ste6_yeast	DSVSSSIINE	IVKKGPPALL	TMVITHSEQM	MRSNCNSIAVL	KDGKVVVERGN
cydc_ecoli	DATTESQILE	LLAEMMREKT	VLMVTHRLRG	LSRFQIIVM	DNGQIIEQGT
ywja_bacsu	DTETEAAIQK	ALQELSEGRT	TLVIAHRLAT	IKDADRIVVV	TNNGIEEQGR
cydd_haein	DAQSENVLVQ	ALNEASQHQ	TLMITHRIED	LKQCDQIFVM	QRGEIVQQGK
mdr_leita	DPALDRQIQA	TVMSAFSAYT	VITIAHRLHT	VAQYDKIIVM	DHGVAEMGS
cyab_borpe	DYESEHIQIR	NMRDICDGR	VIIIAHRLSA	VRCADRIVVM	EGGEVAECGS
msba_ecoli	DYESERAIQ	ALDELQKNRT	SLVIAHRLST	IEKADIVVV	EDGVIVERGT
heta_anasp	DSVSERLIQE	SIEKLSVGRT	VIAIAHRLST	IAKADKVVVM	EQGRIVEQGN
nist_lacla	DPVAEKEIFD	YFVALSENNI	SIFISHSLNA	ARKANKIVVM	KDGQVEDVGS
cydc_haein	DRETERQILR	LILQHAENKT	LIIVTHRLSS	IEQFDKICVI	DNGRLIEEGD
yabj_haein	DQKLRVEMLA	LIAKLeKDLT	LLLVTQHPSE	LisIDQVLVV	ENGQISQLQ.
hst6_canal	DPITTKI IIN	VIKSLHGKLT	ILFVTHDKEL	MRIADNLIVM	KDGQIVEQGD
yfic_bacsu	DTVTEVNIQE	ALARLMEGRT	SVIIAHLNT	IQRADQIVVL	KNGEMIEKGS
	201				250
cftr_human	IQKLLNERSL	FRQAISPSDR	VKLFPHRNSS	KCKSKPQIAA	LKEETEEEVQ
mrpl_human	PSDLLQQRGL	FYSMAKDAGL	V.....	.....	.....
yawb_schpo	TKKLLNKAS	LFYSLAKESG	L.....	.....	.....
ycfi_yeast	PGQLLSdKSL	FYSLCMEAGL	VN.....	.....	.....
sur_human	PEKLLSRKDS	VFASFVRADK	.....	.....	.....
yorl_yeast	PWTLFSQesI	FRSMCSRSGI	VE.....	.....	.....
sur_rat	PEKLLSQKDS	VFASFVRADK	.....	.....	.....
sur_cricr	PETLLSQKDS	VFASFVRADK	.....	.....	.....
yk84_yeast	PYTLISRnf	YRLCRQSGEF	ENLFELAKVS	.....	.....
yhd5_yeast	PSELLkeRGI	FYSMCRDSGG	LELLKqkQSS	K.....	.....
msba_haein	HKTLLQNGA	YKQLHS....	.....	.....	.....
lcnc_lacla	HVDLLAQNGF	YEQ.....	.....	.....	.....
yfib_bacsu	HSELLSESQL	YKRIYE....	.SQFGREGSE	SC.....	.....
ste6_yeast	FDTLYNRRGE	LFQIVSNQSS	.....	.....	.....
cydc_ecoli	HAELLARQGR	YYQ.....	.....	.....	.....
ywja_bacsu	HQDLIEAGGL	YSRLHQAQ..	.....	.....	.....
cydd_haein	FTELQHEGFF	.....	.....	.....	.....
mdr_leita	PRELVMNHQS	MFHS.....	.MVESLGSR	GSKDFYELLM	GRRIVQPAV.
cyab_borpe	HETLLAAGGL	YAR.....	.....	.....	.....
msba_ecoli	HNDLLEHRGV	YAQ.....	.....	.....	.....
heta_anasp	YQELLEQRGK	LWK.....	.....	.....	.....
nist_lacla	HDVLLRRCQY	YQELYSEQY	.....	.....	.....
cydc_haein	YNSLITKENG	FFKRLIER..	.....	.....	.....
yabj_haein	.....	.....	.....	.....	.....
hst6_canal	FQQLISNDGE	FTK.....	.....	.....	.....
yfic_bacsu	HDELIRQKGF	YSDLYESQ..	.....	.....	.....



**Fig. 3.** Diagrams of the three different types of topologies found in Nucleotide Binding Cassette proteins. a) Mitochondrial F<sub>1</sub>-ATPases/recA fold; b) Adenylate Kinase fold; c) Ras P21/EF-Tu/G-Protein fold. In part a of the figure, conserved structural elements common to all three proteins are boxed and striped filled and enclosed in the box; also, the P-loop is highlighted in red. In part c of the figure the dashed line shows a zone where insertions are present; examples of such insertions are: the large  $\alpha$ -helical domain (6 helices) in heterotrimeric G-Proteins, or a “downward”  $\alpha$ -helix in myosin.

$\alpha_3\beta_3\gamma\delta\epsilon$ . The larger subunits  $\alpha$  and  $\beta$  have sequence similarities, and share the same overall fold. The fold can be divided into three domains: a N-terminal beta barrel domain, a nucleotide binding domain, and a carboxy-terminal helical domain. The nucleotide binding domain of this protein was used as a template.

### Sequence Alignment and Modeling of the CFTR Nucleotide Binding Domains

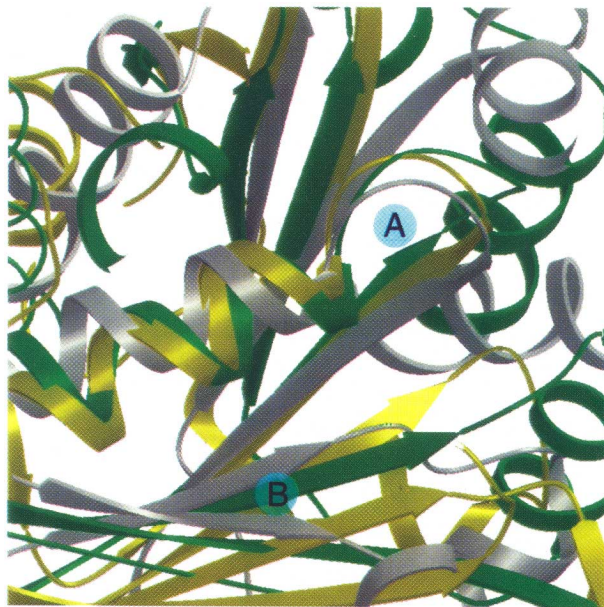
The first step (and probably the most important) in successful modeling is the correct alignment of the target and templates sequences. Figure 5 shows an alignment of the sequences of CFTR NBD1 and NBD2, with F<sub>1</sub> $\alpha$ - and F<sub>1</sub> $\beta$ -subunits respectively. This alignment takes into account conserved structural features in the NBC present in the proteins used as templates. Insertions and deletions in the NBD sequences with respect to the templates must be placed in regions where the F<sub>1</sub>/recA fold can accommodate such changes with little structural disturbance. The first 40 residues in the alignment of both NBDs, which includes motif A, shows strong similarity with the same region of the NBDs of F<sub>1</sub>. The first large insertion appears after this region, at the end of  $\alpha$ A (Figure 3a; residues 482–494 in NBD1 and 1267–1280 in NBD2). This insertion correlates with the high frequency of insertions observed in this zone (positions from 44–48/38–40 in Table 1 and 2 respectively) in ABC transporters. This insertion is compatible with the F<sub>1</sub>/recA fold which has a variable size loop 2 (Figure 3a) in that region. Strand  $\beta$ 2 in the core of the fold, following the insertion, aligns well with the hydrophobic regions 495–503 of NBD1 (1284–1290 of NBD2). The last large insertion is observed in NBD2 before the region that aligns with the strand  $\beta$ 3. This hydrophilic stretch of 14 residues forms a topological insertion which is probably exposed. The equivalent region in F<sub>1</sub> subunits is mainly hydrophobic, residing at the interface between

**Table III.** Alignment of the main secondary elements ( $\alpha$  helices or  $\beta$  strands) found in a representative set of Nucleotide Binding proteins. The minus sign indicates anti-parallel  $\beta$  strands

Protein	Nucleotide Binding Domain Topology						Residue range
F <sub>1</sub> -ATPase	$\beta 1\alpha A\beta 2\alpha B$	$\beta 3\alpha C$	$\beta 4\alpha D$	$\alpha E \beta 5$	$\alpha F\beta 6\alpha G$	( $-\beta 7$ )	150–336
recA	$\beta 1\alpha A\beta 2\alpha B$	$\beta 3\alpha C$	$\beta 4$	$\alpha E \beta 5$	$\alpha F\beta 6$	( $-\beta 7$ )	149–326
ADK	$\beta 1\alpha A$		$\beta 4 \alpha D\alpha D'$	$\alpha E \beta 5$	$\alpha F\beta 6\alpha G\alpha H$	$\beta 7$	9–194
Ras p21	$\beta 1\alpha A$	( $-\beta 3$ )	$\beta 4$		$\alpha F\beta 6\alpha G$	$\beta 7$	1–146
EF-tu	$\beta 1\alpha A$	( $-\beta 3$ )	$\beta 4$		$\alpha F\beta 6\alpha G$	$\beta 7$	99–173
Transducin	$\beta 1\alpha A (\alpha)^6$	( $-\beta 3$ )	$\beta 4$		$\alpha F\beta 6\alpha G$	$\beta 7$	59–323
Myosin	$\beta 1\alpha A (-\alpha B)$	( $-\beta 3$ )	$\beta 4 (\alpha)^8$	$\alpha E \beta 5$	$\alpha F\beta 6\alpha G$		172–680

the NBD and the  $\beta$ -barrel domain that forms the N-terminal domain. Inspection of the sequence alignments of a representative group of ABC transporters proteins (Tables 1 & 2) suggests the presence at position 64 of the fold (504 of CFTR) of a highly conserved negatively charged residue preferentially a glutamic acid followed by a conserved asparagine (in region D), placed between motifs A and B. Mutations such as E504Q do cause CF. Sequence analysis (Fig. 5) aligns this region with an equivalent region of F<sub>1</sub>: loop 3 and helix  $\alpha B$  (Fig. 3a). This suggests that glutamic 504 of CFTR is equivalent to glutamic 188 of the  $\beta$  subunits of F<sub>1</sub>-ATPase. This charged residue is thought to be the general base for the activation of the water

molecule that attacks the  $\gamma$ -phosphate during ATP hydrolysis by F<sub>1</sub> [Abrahams et al. (1994)]. In the  $\alpha$ -subunit of F<sub>1</sub>, which is not active in ATP hydrolysis, a glutamine residue (Q208) is present in this position. An equivalent residue, Q1291, is present in NBD2 and highly conserved in the NBD2-like class (Table 2), suggesting a role for NBD2 similar to that of the F<sub>1</sub>  $\alpha$ -subunit. Based on these and other considerations, NBD1 was modeled based on F<sub>1</sub> $\beta$  and NBD2 based on F<sub>1</sub> $\alpha$ . The sequence of loop 5 of F<sub>1</sub> (EPP or DAA in  $\alpha$  and  $\beta$ -subunits) can be aligned with conserved residues [E/D]XG, just before motif C, which is highly conserved in ABC transporters (Tables 1 & 2). A sequence analogous to motif C seems to be absent in F<sub>1</sub>.



**Fig. 4.** View of a spatial alignment of Nucleotide Binding (NB) cassettes of a representative set of NB proteins: ADK is shown in silver, P21 in gold and recA in green.

### Three-dimensional Modeling of the CFTR Nucleotide Binding Domains

The detailed three-dimensional models of the NBDs of CFTR are presented in Figures 6 (NBD1) and 7 (NBD2). Figures 8 and 9 show the secondary structure predictions [Chou & Fasman (1978)] in comparison with the suggested models based on the F<sub>1</sub>/recA topology. Only two insertions, positions 44–53 in NBD1 and position 90 in NBD2, required careful adjustment. The structural alignment of F<sub>1</sub> subunits  $\alpha$  and  $\beta$  to each other was based on the structure of rat liver F<sub>1</sub>-ATPase [Bianchet et al. (1998)]. The two NBD structural models are similar to each other as expected from the similarity of their templates, the  $\beta$  and  $\alpha$  subunits of F<sub>1</sub>. The starting motif, N-terminus— $\beta 1\alpha A$ , was predicted to be similar in size to that of F<sub>1</sub>. The insertion found in ABC transporters (after position 44 and 38 in Tables 1 and 2 respectively) was modeled as a longer  $\alpha A$ , a loop and a short anti-parallel  $\beta$ -strand in the NBD1. A loop containing the D region

				-- A --				
Beta	130	QEILVTGIVK	VDLLAPYAKG	GKIGLFGGAG	VGKTVLIMEL	INNV_AKAH	_____G_____	
Alpha	143	REPMQTGIKA	VDSLVPPIGRG	QRELIIGDRQ	TGKTSIAIDT	IIN_QKRF	N_DGTDE	
NBD1	441	_____L	KDINFKIERG	QLLAVAGSTG	AGKTSLLMMI	MGELEPSEK	IKHSGRIS	
NBD2	1227	_____L	ENISFSISPG	QRVGLLGRGTG	SGKSTLLSAF	LRLLENT_EGE	IQIDGVSWDS	
			---D---					
Beta	178	_____GGYSVF	AGVGERTRREG	NDLYHEMIES	GVINLKDAT	_____	_____SKVAL_V	
Alpha	196	_____KKKLYCIY	VAIGQKRSTV	AQLVKRLTDA	DAM_____	_____	_____KYTIVV	
NBD1	490	FCSQFSWIMP	GTIKENIIFG	VSYDEYRYS	VIKACQLEED	_____	_____ISKFAEKD	
NBD2	1277	ITLQQWRKAF	GVIPQKVFIF	SGTFRKNDLP	YEQWS_DQEI	WKVADEVGLR	SVIEQFPQKL	
			^	^				
			--C-		-----B-----			
Beta	219	YGQMNPEP	_____GARARVA	LTGLTVAEYF	RDQEGQDVLL	FIDNIFRFTQ	AGSEVSALLG	
Alpha	233	SATASDAA	_____PLQYLAP	YSGCSMGYF	RDNGKHALII	YDDLKQAVA	YRQMSLLLR	
NBD1	538	NIVLGEGGIT	LSGGQRARIS	L_____A	RAVYKDADLY	LLDSPFGYLD	VLTEKEIFES	
NBD2	1336	DFVLVDGGCV	LSHGKQLMC	L_____A	RSVLSKAKIL	LLDEPSAHL	PVT_YQIIRR	
Beta	274	RIPSAVGYQP	TLATDMGTMQ	ERITTTKK	_____GSITSV	QAIYVPADDL	TDPAPATFFA	
Alpha	288	PPGREAYPGD	VFYLHSRLL	ERAAMNDSF	GG_GSLTAL	PVIETQAGDV	SAYIPTNVIS	
NBD1	590	CVCKLMANKT	RILVTSKM	EHLKKADKIL	ILNEGSSYFY	GTFSELQNLQ	PDFSSKLMGC	
NBD2	1387	TLRQAFADCT	VILCEHRI	EAMLECCQFL	VIENKVRQY	DSIQKLLNER	SLFRQAISPS	
Beta	328	HLDATTVLSR	AIAELGIYPA	VDPLDSTSR	MDPNIVGS			
Alpha	345	ITDQIFLET	ELFYKGIKIRPA	INVGLSVSRV	GSAAQTRA			
NBD1	648	DSFDQFSAER	RNSILTETLH	RF_SLEGDAP	VSWTETKK			
NBD2	1444	DRVKLF__PH	RNSSKCKSKP	QIAALKETE	EEVDTRL			

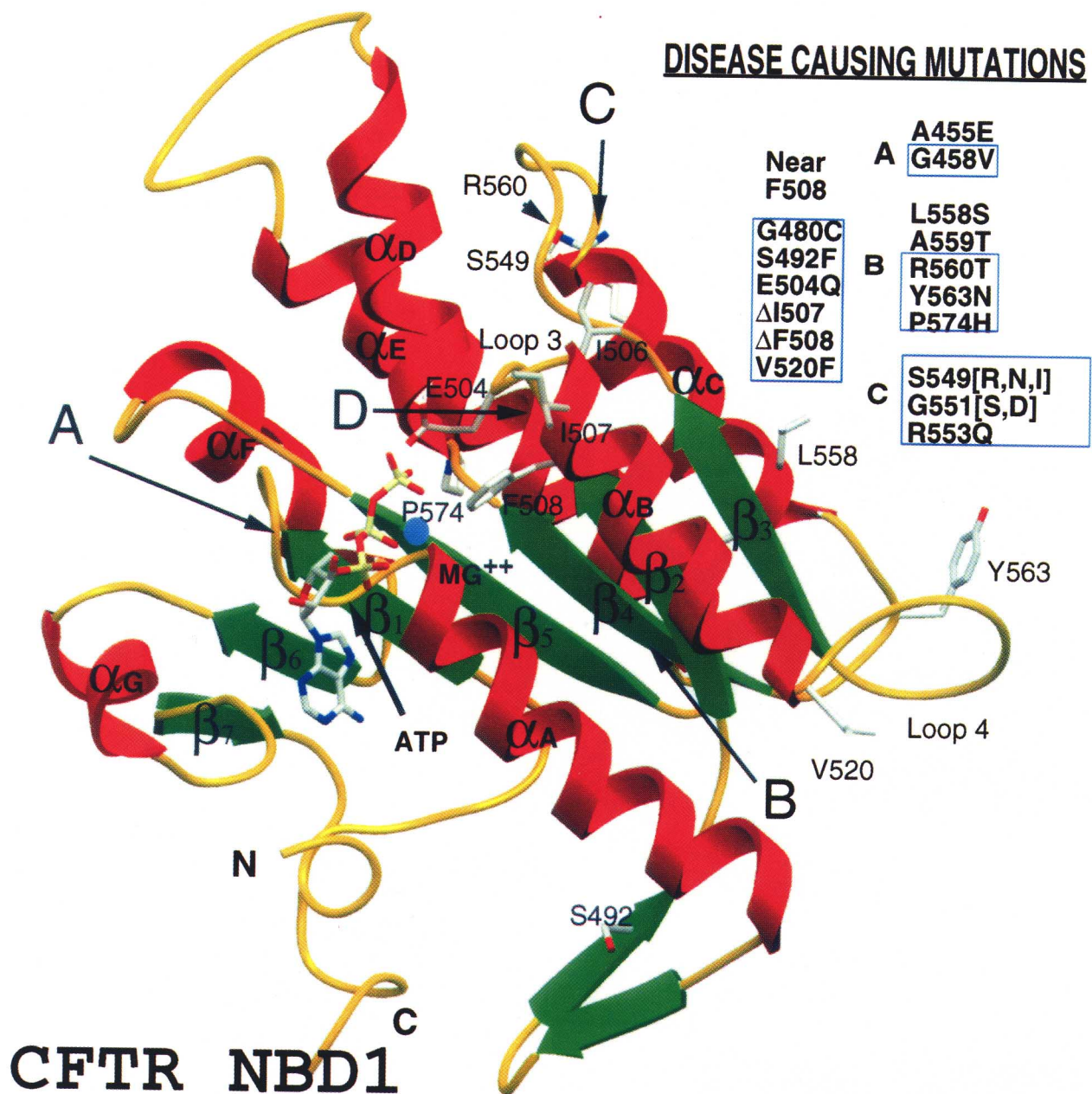
**Fig. 5.** Alignment of the NBD of CFTR with  $\alpha$ - and  $\beta$ -subunits of rat Liver  $F_1$ -ATPase. Residues are represented by a single letter code, colored by hydrophathy; black, hydrophobic; red, acidic hydrophilic; green, neutral hydrophilic; and blue, basic hydrophilic. Walker motifs A and B and regions C and D are indicated by brackets above the sequences. The positions of the proposed catalytic base and of F508 of CFTR are indicated by arrows below the sequences. The sequence alignment of both subunits of  $F_1$  are based on the alignment of the three-dimensional structures of  $\alpha$ - and  $\beta$ -subunits of rat Liver  $F_1$ -ATPase.

was placed close to the P-loop. The position of E504 in NBD1, and Q1291 in the NBD2 are, as expected, near the  $\gamma$ -phosphate binding pocket. Helix  $\alpha$ B was modeled longer than the one observed in  $F_1$  and recA. F508 is present at the beginning of  $\alpha$ B, close to the adenosine binding pocket. The overall folds of the motifs A and B in the model presented here are similar to those proposed in previous studies [Hyde et al. (1990), Mimura et al. (1991), Carson et al. (1995)]. However, the sequence of region D (501-TIKENIIF-508 in CFTR NBD1 or 1287-GVIPQKVFIF-1296 in NBD2) is placed in a loop and in the first turn of helix  $\alpha$ B, forming part of the binding site and probably interacting with the adenosine moiety. The end of region D in our model is associated with a hydrophobic patch close to the binding site and probably at an interface region. The model suggests that E504 has a function similar to that of E188 in  $F_1$ -ATPase. Certainly, the disease causing mutation E504Q, and probably the  $\Delta$ I507 and  $\Delta$ F508 mutations have an effect on the function, position, and in the geometry of the loop that contains E504. Significantly, a synthetic peptide that includes the carboxyl end of region D can bind ATP (unpublished observation), suggesting the partici-

pation of this region in ATP binding, as predicted by our model.

In the present model of NBD1 the adenosine and the sugar moiety are close to F508. Deletion of one residue at the beginning of helix  $\alpha$ B could produce a  $100^\circ$  rotation of the helix if the loop containing E504 remains fixed; this would produce a major change in a region on the surface of the model. This change suggests an explanation for why  $\Delta$ I507 and  $\Delta$ F508 manifest themselves as folding mutants (Thomas et al. (1993) and references within), with temperature dependent membrane traffic and targeting problems [Cheng et al. (1990), Rich et al. (1991)]. Our model places motif C near the loop containing E504 and region D. The G551D mutation, which occurs in this region, may have charge and steric conflicts with E504 which account for the effects of the mutation. It is also possible that the G551D mutation impairs critical NBD1 interactions with other CFTR domains.

In the second nucleotide domain the putative catalytic base is a glutamine (Q1291). Unlike E504 in NBD1, Q1291 is predicted to be a poor catalytic base. Therefore, NBD2 is predicted to have a much lower catalytic capacity than NBD1. It is possible also that

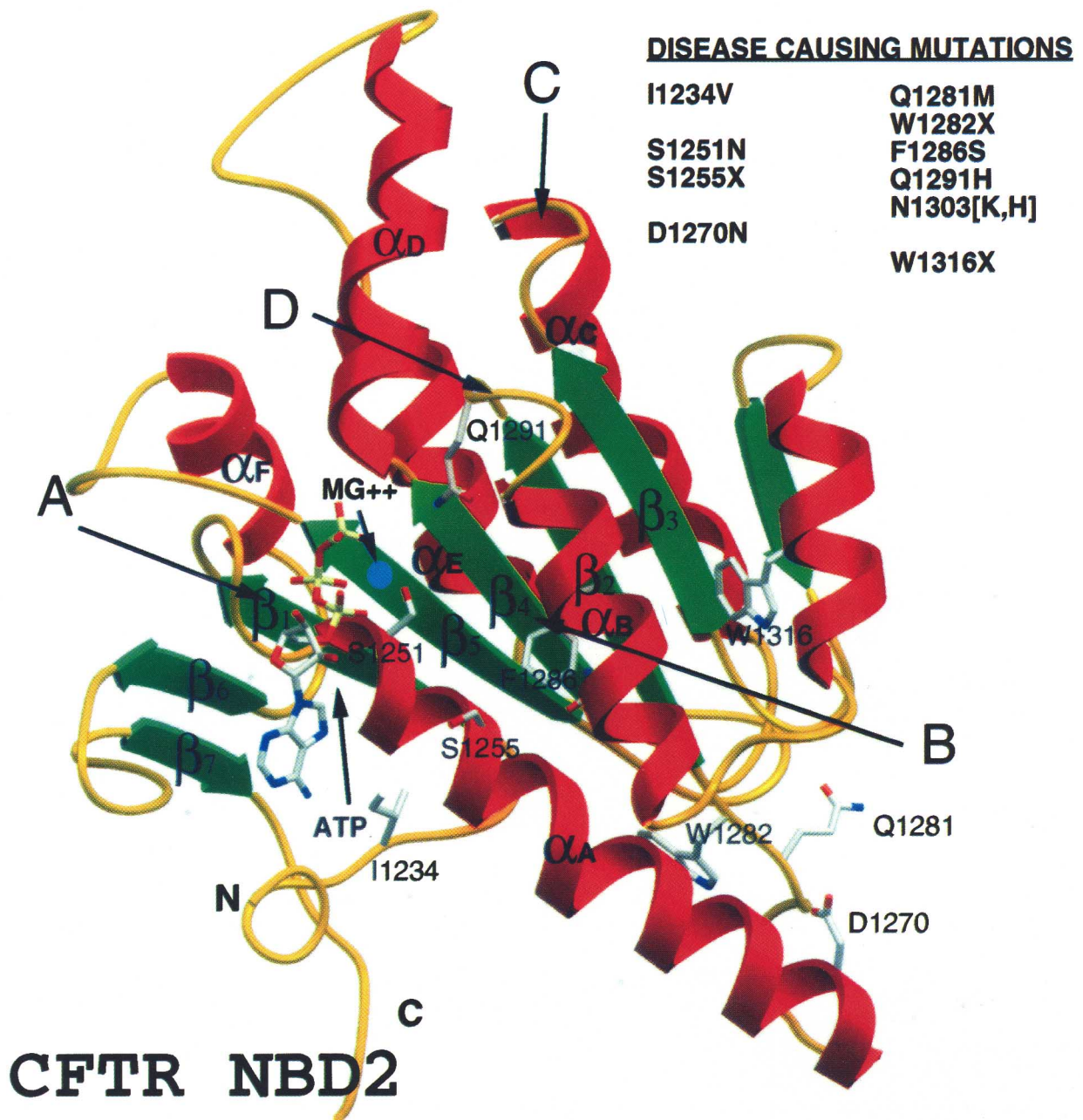


**Fig. 6.** *F<sub>1</sub>/recA* based model of the first Nucleotide Binding Domain (NBD1) of CFTR. An ATP molecule is shown bound in the proposed binding site. Residues which when mutated cause CF are represented and tabulated in the Figure.

NBD2 has no catalytic capacity and only serves an ATP-dependent regulatory role. The same pattern is found in the alpha subunit of F<sub>1</sub>-ATPase in which the equivalent residue is also a glutamine (Q208, F<sub>1</sub> numbering). To date, the intact  $\alpha$ -subunit is believed to have no catalytic capacity. It should be noted, however, that in NBD2-like domains frequently the charged residue that follows the conserved glutamine of region

D is an acidic one, usually aspartate (with lower frequency E), which may serve as catalytic base, e.g. MDR1 (*mrd1\_human* in Table 2). This is not the case in all ABC transporters, e.g., CFTR and Ste6, where a basic residue (K and R respectively) follow the conserved glutamine.

Overall the three-dimensional models contain 35/37%  $\alpha$ -helix and a 24/25%  $\beta$ -strand, for NBD1/NBD2.



**Fig. 7.** *F<sub>1</sub>/recA* based model of the second Nucleotide Binding Domain (NBD2) of CFTR. An ATP molecule is shown bound in the proposed binding site. Residues which when mutated cause CF are represented and tabulated in the Figure.

These values are in good agreement with values obtained using circular dichroism spectroscopy that show 40%  $\alpha$ -helix and 24%  $\beta$ -strand using recombinant NBD1 of the traditional size (F443-S589) [Ko et al. (1994)]. The model of NBD1 presented here shows 39%  $\alpha$ -helix and 24%  $\beta$ -strand in the same region.

#### Comparison with Previous Models

Several models of ABC transporter NBDs have been discussed in the literature, three of CFTR [Hyde et al. (1990), Carson et al. (1995), Annereau et al. (1997)] and one of the periplasmic histidine permeases



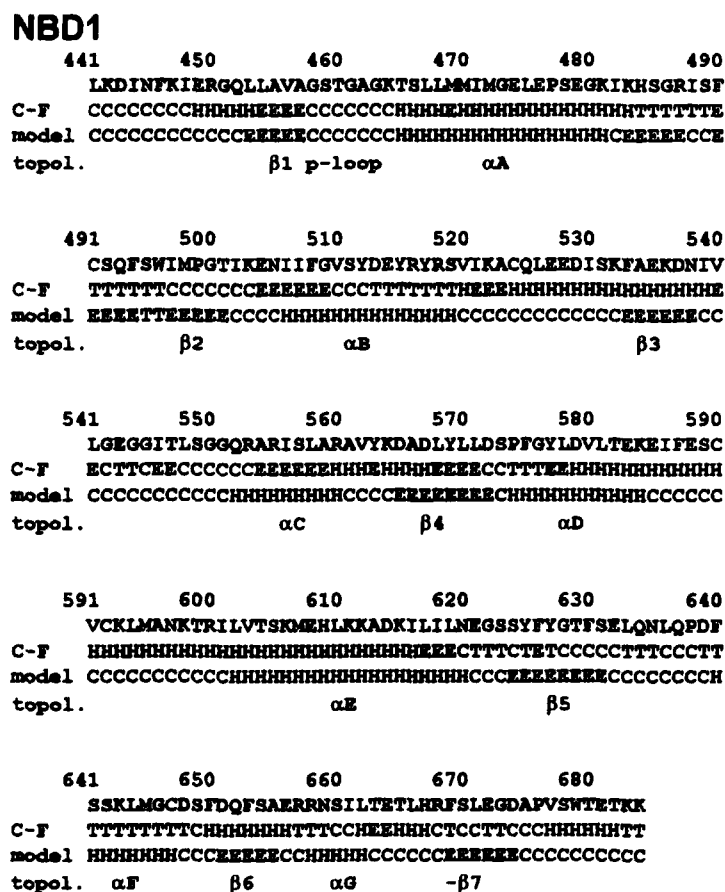


Fig. 8. Secondary structure assignment of the first nucleotide binding domain of CFTR. The sequence of NBD1 is aligned with the predicted secondary structure using a Chou and Fassman method (second line), and with the proposed secondary structure (third line). H, T, C and E are helix, turn, coil and beta strand respectively. The name of the secondary structural elements are shown in the last line.

[Mimura et al. (1991)]. All these models are based on the three folding motifs described previously (Fig. 3 a,b,c). The NBD of Hyde et al. (1990) and Mimura et al. (1991) were modeled using the ADK-fold and the Carson et al. (1995b) model was based on G-proteins. Recently, Annereau et al. (1997) reported a model based on  $F_1$ -ATPase. The NBDs in all the models have motifs A and B present in the secondary structural motif  $\beta$ 1 $\alpha$ A ...  $\beta$ 4 (boxed in Figure 3a). Table 4 summarizes the similarities and differences among the different NBD models.

Hyde et al. (1990) and Mimura et al. (1991) made the first attempt to obtain a model of the nucleotide binding domain of ABC transporters. The model is based on the nucleotide binding domain of adenylate kinase, one of the few NBC proteins with a known three-dimensional structure at the time. This enzyme

catalyzes the conversion of ATP and AMP to two ADP molecules and contains two nucleotide binding sites. An ADK-like fold, such as the one proposed by Hyde et al. (1990), assigns the region of F508 to an alpha-helical subdomain (AMP site in ADK), far away from the nucleotide site. Motif C (LSXGX[R/K]) has no role in this model. In contrast, motif C participates in ATP (or GTP) binding in models based on the G-protein-fold such as the one proposed by Carson et al. (1995b). It has been suggested that motif C is associated with nucleotide binding based on a proposed homology between CFTR and a heterotrimeric G-protein [Manavalan et al. (1995)]. The sequence DX[G/A]GQR, in heterotrimeric G-proteins, participates in binding of the GTP phosphates,<sup>4</sup> and Dearborn and Manavalan (1994) suggest that this sequence is equivalent to motif C of the ABC transporters. Supporting this

## NBD2

```

1227      1240      1250      1260      1270
LENISFSISPGQRVGLLGRVSGKSTLLSAFLRLLNTEGEIQIDGVSWDS
C-F CCCCCCCCCCTTEEEEEETCTCCCCEEEEHHHHHCCCCHEEECCCCC
model HHHHCCCCCEEEEEEECCCCCHHHHHHHHHHHHHCCCCCCCCCCCCC
topol.           $\beta$ 1 p-loop           $\alpha$ A

1280      1290      1300      1310      1320
ITLQGNKAFGVIPQKVFIFSGTFRKRLDPYEQWSDQEINKVADEVGLRS
C-F CECTTCTTTTCCCEEEEEETTTTTCCTTTTCHHHHHHHHHHHHHHE
model CCCCCCEEEEEETTCCHHHHHHHHHHHHHCCCCCHHHHHHHHHHHHE
topol.           $\beta$ 2           $\alpha$ B           $\alpha$ B'

1330      1340      1350      1360      1370
VIEQFPKLDVFLVDGGCVLSHGHRQLMCLARSVLSKAKILLLDEPSAHL
C-F HEEECTTCCEEEEETTTHEEHTTTHHEEHHHHHHHHHHHEEHHHTHCC
model EETTEEEEEEECCCCCCCCCHHHHHHHHHHHHCCCEEEEEEEHHHHH
topol.           $\beta$ 3           $\alpha$ C           $\beta$ 4           $\alpha$ D

1380      1390      1400      1410      1420
DFVTYQIIRRTLKQAFADCTVILCEHRIEAMLECCQFLVIEENKVRQYDS
C-F CTTEEEEEETCTTTTTHHHHHHHHHHHHHHHHHHHHHHHHHHHHTTE
model HHHHHHHHCCCCCCCCCHHHHHHHHHHHHHHHHHHHHHCCCCCCCCCEEEEE
topol.  $\alpha$ D           $\alpha$ E           $\beta$ 5

1430      1440      1450      1460      1470
IQKLLNERSLFRQAI SPSDRVLFPHRNSKCKSRKQIAALKEETEEVQ
C-F EEECCCHHHHEETCCCTTEEEETTTTPTTTTCHHHHHHHHHHHHHHH
model ECCCCCCHHHHHHHHHCCCCCEEEEECCCCCCCCCEEEEECCCCC
topol.           $\alpha$ F           $\beta$ 6          - $\beta$ 7

1480
DTRL
C-F HHHH
model CCCC
topol.

```

Fig. 9. Secondary structure assignment of the second nucleotide binding domain of CFTR. The sequence of the NBD2 of CFTR is aligned with the secondary structure predicted using a Chou and Fassman method (second line), and the proposed secondary structure (third line). H, T, C and E are helix, turn, coil and beta strand respectively. The name of the secondary structural elements are shown in the last line.

view, is the fact that GTP binding has been observed in a synthetic peptide of the NBD2 of CFTR [Randak et al. (1996)]. However, this is not compelling because GTP binding and even GTPase activity, are usually found in ATPases. Arguing against the hypothesis depicting the participation of motif C in nucleotide binding is the observation that a shorter synthetic peptide that does not include motif C can still bind ATP [Ko et al. 1994]. Also, the suggested equivalent of

motif C, DX[G/A]GQR, occurs just after motif B in G-proteins, and not before as in the ABC transporters. Topologically, this sequence is inserted in the loop before  $\beta$ 4 in ras P21/G-proteins (Fig. 3c). However, the G-protein based model suggests a direct explanation for why the mutation G551D, in the core of motif C, affects the ATP binding function. Supportive of these finding is the corrective mutation R555K. Nevertheless, a G-protein like fold for CFTR requires a large topological insertion between  $\alpha$ A and  $\beta$ 4 as is observed in the heterotrimeric G-proteins. The highly conserved feature of the motif B, located in this insertion, would be affected. Because of this, no explanation for other

<sup>4</sup> The proposed general base, E202 (G-Protein numbering), comes just after this sequence.

**Table IV.** Comparison of the predicted position of motif C and F508 in various models for CFTR

Model	Protein used in modeling	Position of motif C		Position of F508	
		Secondary structure	Relative to nucleotide binding site	Secondary structure	Relative to nucleotide binding site
This paper	F <sub>1</sub>	loop	near	$\alpha$ -helix	within
Annereau et al. (1997)	F <sub>1</sub>	loop	near	$\beta$ -strand	outside
Hyde et al. (1990)	ADK	?	?	$\alpha$ -helix	outside
Carson et al. (1995)	G-proteins	loop	within	?	?

conserved features (such as the mutations in the region D) can be given based on this model.

Although, Annereau's et al. (1997) model also uses F<sub>1</sub>-ATPase as template, their NBD sequence alignment differs with that presented here (Fig. 5). This model localizes F508 on a  $\beta$ -sheet buried in the core of the nucleotide binding fold.

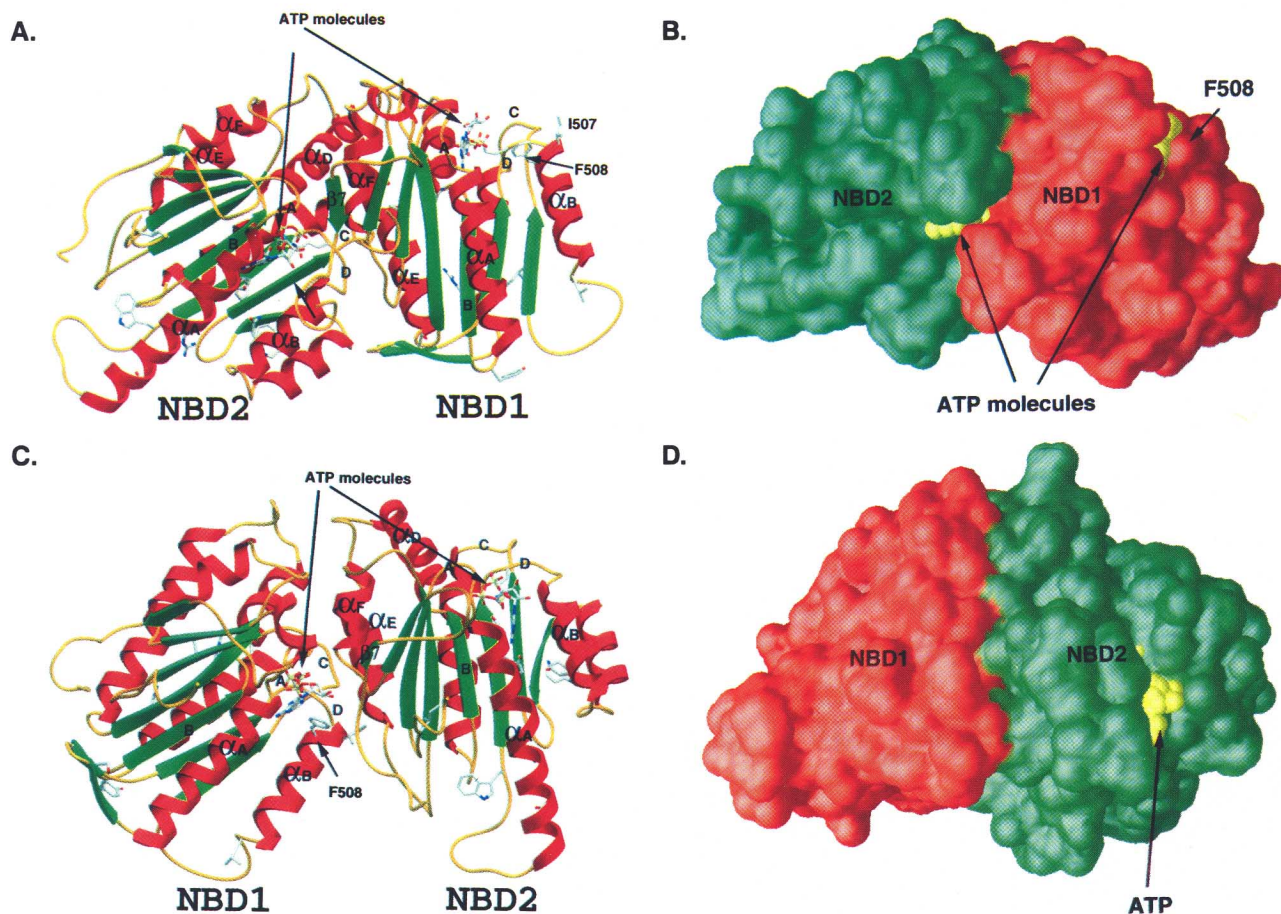
#### Quaternary Association of NBD1 and NBD2 of CFTR

Modification of only one NBD in ABC transporters affects the total ATPase function. In particular, the observed behavior of CFTR in ATP hydrolysis [Carson et al. (1995)] suggests a close interaction between the two different NBDs. ATP hydrolysis in the second domain is believed to close the channel opened by the hydrolysis of ATP in the first domain. Also, mutation of lysines K458 in NBD1 and K1250 in NBD2 produce variations in the start and in the duration of the channel activity burst, respectively. These experiments assign a distinct function to each NBD and suggest a proximity relation between both domains. The assembly of recA into a recA polymer [Story et al. (1992), Yu and Egelman (1997)] and the  $\alpha_3\beta_3$  association of subunits in F<sub>1</sub> show similarities and suggest a mode of association of the two NBDs in CFTR. The NBDs in F<sub>1</sub>-ATPase, i.e.  $\alpha$ - and  $\beta$ -subunits, have two interfaces with NBDs of adjacent subunits. Two different contact interfaces between molecules are present in the F<sub>1</sub> trimer of  $\alpha\beta$  pairs:  $\alpha/\beta$  and  $\beta/\alpha$ . Each type of interface has one of the ATP binding sites [Abrahams et al. (1994), Bianchet et al. (1997)]. This molecular quaternary association NBD1/NBD2 may also occur in CFTR, and two types of interactions are possible, one similar to the  $\alpha/\beta$  interaction in F<sub>1</sub>, in which the nucleotide binding site of NBD2 is at the interface and another similar to the

$\beta/\alpha$  interface in F<sub>1</sub>, in which the nucleotide binding site of NBD1 is at the interface. In both modes, the hydrolysis of ATP, as in F<sub>1</sub>-ATPase could be associated with a conformational change [Abrahams et al. (1994)] easily transmitted to the second domain through the interface.

Helices  $\alpha$ E and  $\alpha$ F and strand  $\beta$ 7 in NBD1 loop5 (motif C), helix  $\alpha$ D and loop 3 (P1290-F1294, part of motif D) in NBD2, are placed in the buried interface in the  $\alpha/\beta$ -like association (Fig. 10a). Symmetrically,  $\alpha$ E,  $\alpha$ F and strand  $\beta$ 7 of NBD2 and loop 5 (motif C),  $\alpha$ D and loop 3 (region D) of the NBD1 are in the buried interface in the  $\beta/\alpha$ -like association (Fig. 10b); F508 is buried in this mode of association. Relevant regions such as motif C of NBD1, the locus for several of the critical disease causing mutations in CF, are localized at the proposed interfaces in the  $\beta/\alpha$ -like association, and more "exposed" in the  $\alpha/\beta$ -like interface. Residues at the onset of the motif B, where mutations cause CF, are in all the cases at the "exposed" face of helix  $\alpha$ C in both NBDs. The effect of mutations in this region, are likely to affect the interaction of the NBDs with others portions of the protein.

The extent and the characteristics of the accessible surface area (ASA) buried in each mode of association (Table 5) can be used to rank the two possibilities. ASA calculations in both types of association shows the NBD2 buries more polar and less apolar area and the reciprocal for NBD1. Also, there are substantial differences in the amount of buried apolar and polar surfaces associated with the two possible interfaces. The putative  $\alpha/\beta$ -like interface buries 1906 Å<sup>2</sup> of accessible apolar area and 475 Å<sup>2</sup> of polar accessible area, 245 Å<sup>2</sup> more of apolar and 304 Å<sup>2</sup> less of polar area than the  $\beta/\alpha$ -like interface. From a thermodynamical point of view, the  $\alpha/\beta$ -like association is favored. Figures 11a,b show the accessible surface of the  $\alpha/\beta$ - and  $\beta/\alpha$ -like quaternary association respectively, with



**Fig. 10.** Modes of Quaternary association of the nucleotides Binding Domains. a) Ribbon cartoon and b) space filling model of the  $\alpha/\beta$ -like association, with each domain of different color, red NBD1 and green NBD2; c) ribbon cartoon and d) space filling model of the  $\beta/\alpha$ -like association, with each domain of different color, red NBD1 and green NBD2.

the hydropathy of each surface atom color coded (green hydrophobic and red hydrophilic). The figures depict the hydropathy pattern of a typical soluble protein with no large hydrophobic patches. An electrostatic analysis (Figure 11c,d) displays a more charged NBD2 with

concentrated positive charge (colored blue) near the nucleotide binding site.

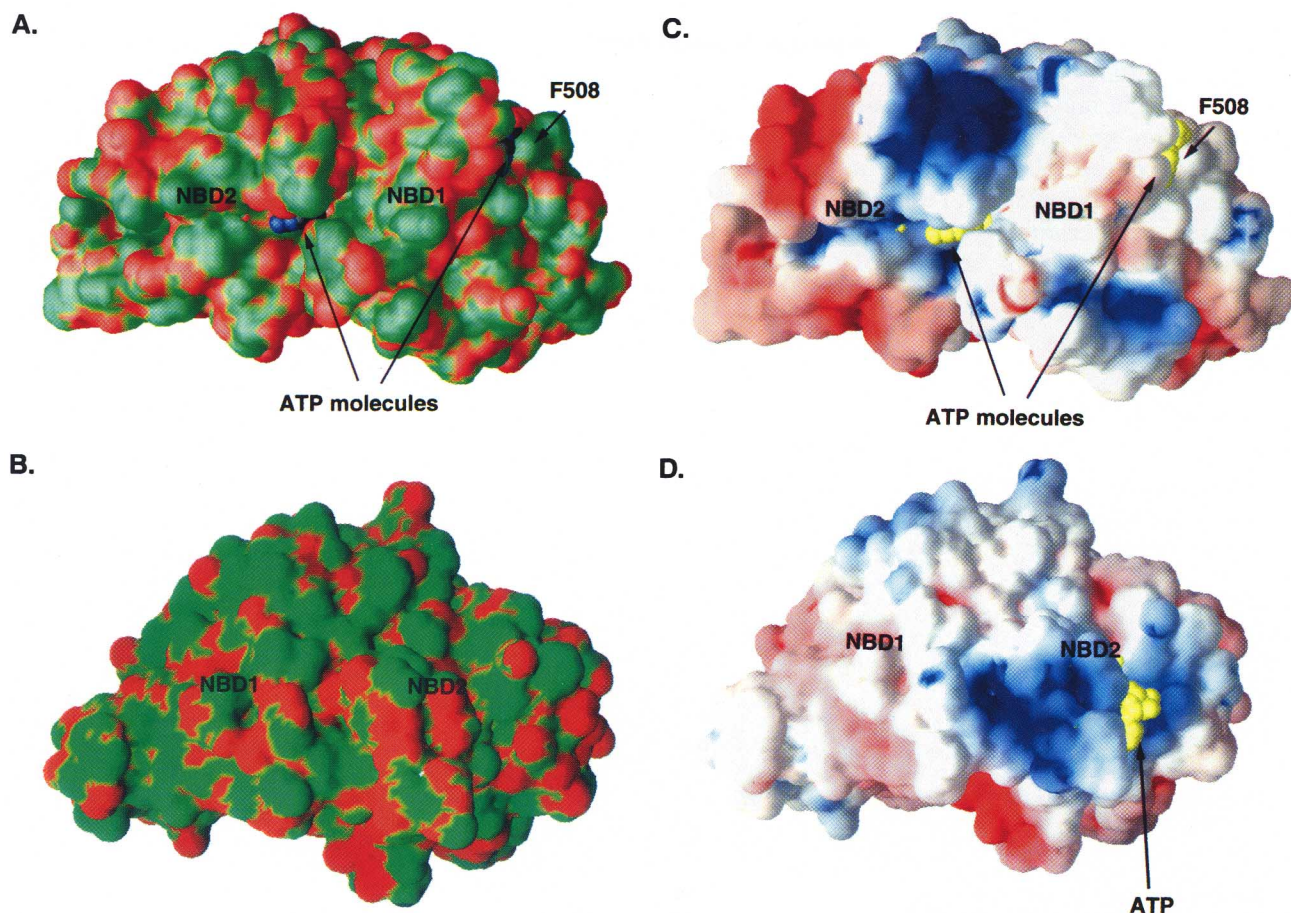
The end of region D where the critical F508 residue lies, together with 4–6 hydrophobic residues upstream of motif C, are associated with a hydrophobic patch, which is relatively exposed in the favored  $\alpha/\beta$ -like model. The concentration of mutations that cause CF in this zone and the intracellular traffic problems associated with these mutations, suggest that this hydrophobic patch may have a function in sorting and membrane trafficking of the CFTR protein to its final location in the plasma membrane.

**Table V.** Accessible surface area for the two different modes of association proposed. Contributions of each domain alone in each mode are also shown. All the areas are in  $\text{\AA}^2$

		NBD1	NBD2	Total
not associated	apolar	9047.6	8692.6	17740.2
	polar	4296.9	4582.0	8878.9
in the $\alpha/\beta$ -like association	apolar	8096.4	7737.7	15834.1
	polar	4022.1	4382.1	8404.2
in the $\beta/\alpha$ -like association	apolar	8132.4	7946.8	16079.2
	polar	3926.9	4173.7	8100.6

## CONCLUSIONS

The structure of  $F_1$  provides a useful template for modeling the nucleotide binding domains of ABC



**Fig. 11.** Views of the accessibility surface of the NBD complexes. Carbon and sulfur atoms depicting hydrophobic surface are colored green and nitrogen and oxygen atoms depicting hydrophilic surface are in red. a)  $\alpha/\beta$ -like, b)  $\beta/\alpha$ -like, c) and d) their respective molecular surfaces colored with the calculated electrostatic potentials in KT units, blue 5KT (positive), red  $-5KT$  (negative).

transporters and in particular those of CFTR. The NBDs alignment, hydropathy profile, and alignment against the  $F_1$  subunits, suggest NBDs longer than that traditionally accepted. The CFTR NBD1 model that results (Fig. 6) gathers the disease causing mutations in three different clusters: (1) mutations affecting the nucleotide binding pocket and the putative general base: A455E, G458V, E504Q  $\Delta$ I507  $\Delta$ F508 P574H; (2) mutations in motif C which are probably related to an interaction with region D: S549[R,N,I] G551[S,D], R553Q; and (3) mutations within or near motif B, L558S, A559T, R560T, Y563N and mutations S492F and G480C. Sequence alignments suggest an association of the two domains similar to that of the  $\alpha$  and  $\beta$  subunits of  $F_1$ . Deletions near the nucleotide binding pocket are certainly predicted to affect the proposed fold around the putative interface, and impair nucleotide associated functions. Changes in

other regions predicted to be at the interface between CFTR domains might affect assembly of the protein. The proposed domain association also allows for a cooperative behavior between NBD1 and NBD2 as observed in intact CFTR [Carson et al. (1995a)]. Therefore, it seems possible, that the hydrolysis of ATP in NBD1 produces a conformational change that opens the channel and makes the second NBD active, until binding or hydrolysis of ATP in NBD2 reverts the opening and returns the NBD1 domain to its original state.

The models presented here indicate that NBD2 may exhibit a much lower ATP hydrolysis rate than NBD1 or perhaps no hydrolytic rate at all. Therefore, more ATP utilization would be predicted to be involved in channel opening than closing. It remains possible that channel closing requires preferentially binding of ATP to NBD2 rather than hydrolysis *per se*, and that

this ATP binding to NBD2 induces a conformational change in NBD1.

Although, this paper focuses on CFTR, most of the modeling procedures can be generalized to other ABC transporters. Here it should be noted that simply by examining the region D to establish whether it has a glutamic acid residue (E) or a glutamine residue (Q), one can predict which of the two NBDs (or both) may be functioning primarily in ATP hydrolysis.

## ACKNOWLEDGMENTS

We thank Dr. Margarita Faig for reading the manuscript and for many helpful discussions during its preparation.

## REFERENCES

- Annereau J.P., Wulbrand U., Vankeerberghen A., Cuppens H., Bon-temis F., Tummler B., Cassiman J.J. and Stoven V. (1997) *FEBS Lett.* **407**, 303–308.
- Abrahams J.P., Leslie A.G.W., Lutter R. and Walker J.E. (1994) *Nature* **370**, 621–628.
- Bianchet M.A., Hüllihnen J., Pedersen P.L. and Amzel L.M. (1998) In Review.
- Bianchet M., Ko Y.H. and Pedersen P.L. (1995a) *Ped. Pulmon. Suppl.* **12**, abstract 32.
- Bianchet M.A., Ko Y.H., Amzel L.M., Pedersen P.L. (1995b) *Biophysical J.* **70**(2):A214.
- Carson M.R., Travis S.M. and Welsh M.J. (1995a) *J. Biol. Chem.* **270**(4):1711–1717.
- Carson M.R., Travis S.M. and Welsh M.J. (1995b) *Biophysical J.* **69**, 2443–2448.
- Cheng S.H., Gregory R.J., Marshall J., Paul S., Souza D.W., White G.A., O'Riordan C.R. and Smith A.E. (1990) *Cell* **63**, 827–829.
- Cheng S.H., Rich D.P., Marshall J., Gregory R.J., Welsh M.J., Smith A.E. (1991) *Cell* **66**, 1027–1036.
- Chou P.Y. and Fasman G.D. (1978) *Adv. Enzymol.* **47**, 45–148.
- Denning G.M., Ostergard L.S., Cheng S.H., Smith A.E. and Welsh M.J. (1992) *J. Clin. Invest.* **89**, 339–349.
- Dalemans W., Barbry P., Champigny G., Jallat S., Dott K., Dreyer D., Crystal R.G., Pavirani A., Lecocq J.P. and Lazduski M. (1991) *Nature* **354**, 526–528.
- Doige C.A. and Ames G.F.-L. (1993) *Annu. Rev. Microbiol.* **47**, 291–319.
- Evans S.V. (1993) *J. Mol. Graphics* **11**, 134–138.
- Higgins C.F. (1992) *Annu. Rev. Cell Biol.* **267**, 6455–6458.
- Hyde S.C., P. Emsley, M.J. Harsthorn, M.M. Mimmack, U. Gileadi, S.R. Pearce, R.E. Hubbard & C.F. Higgins. (1990) *Nature* **346**, 362–365.
- Kerem E., Corey M., Kerem B.S., Rommens J., Mannewitz D., Kobayashi K., Knowles M.R., Boucher R.C., O'Brien W.E., Beaudet A.L. (1990) *J. Hum. Genet.* **110**, 599–605.
- Ko Y.H. and Pedersen P.L. (1995) *J. Biol. Chem.* **268**, 24330–24338.
- Ko Y.H., Thomas P.J. and Pedersen P.L. (1994) *J. Biol. Chem.* **269**, 14584–14588.
- Ko Y.H., Delannoy M. and Pedersen P.L. (1997) *Biochem.* **36**, 5053–5064.
- Levison H., Tsui L.C., Durie P. (1990) *N Engl. J. Med.* **323**, 1517–1522.
- Li C.H., Ramjeesingh Wang W., Garami E., Hewryk M., Lee D., Rommens J.M., Galley K., Bear C.E. (1996) *J. Biol. Chem.* **271**, 28463–28468.
- Manavalan P., Dearborn D.G., McPherson J.M. and Smith A.E. (1995) *FEBS Lett.* **366**, 87–91.
- Morikawa K., la Cour T.F., Nyborg J., Rasmussen K.M., Miller D.L., Clark B.F. (1978) *J. Mol. Biol.* **5**:125, 325–338.
- Nicholls A., Sharp K. and Honig B. (1991) *PROTEINS: Structure, Function and Genetics*, **11**, 281–296.
- Noel J.P., Hamm H.E., Sigler P.B. (1993) *Nature*, **366**, 654–663.
- Pai E.F., Krengel U., Petsko G.A., Goody R.S., Kabsch W., Wittinghofer A. (1990) *EMBO J.*, **9**, 2351–2359.
- Pedersen P.L. and Amzel L.M. (1993) *J. Biol. Chem.* **268**, 9937–9940.
- Qu B.-H. and Thomas P. (1996) *J. Biol. Chem.* **271**, 7261–7264.
- Riordan J.R., Rommens J.M., Kerem B., Alon N., Rozmahel R., Grzelczak Z., Zielenski J., Lok S., Plavsky N., Chou J., Drumm M.L., Iannuzzi C., Collins F.S. and Tsui L. (1989) *Science* **245**, 1066–1073.
- Rich D.P., Anderson M.P., Gregory R.J., Cheng S.H., Paul S., Jefferson D.M., McCann J.D., Klinger K.W., Smith A. and Welsh M.J. (1991) *Nature* **347**, 358–363.
- Rossmann M.G., Moras D. and Olsen K.W. (1974) *Nature* **250**, 194–199.
- Rost B. (1996) *Meth. in Enzym.* **266**, 525–539.
- Saraste M., Sibbald P.R. and Wittinghofer A. (1990) *Trends in Biochem. Sci.* **15**, 430–434.
- Shirakihara Y., Leslie A., Abrahams J.P., Walker J., Ueda T., Sekimoto Y., Kambara M., Saiga K., Odaka M., Yoshida M. and Kagawa Y. (1997) *Structure* **5**, 825–836.
- Story R.M., Weber I.T., Steitz T.A. (1992) *Nature* **355**, 567.
- Schulz G.E., Elzinga M., Marx F., Schrimmer R.H. (1974) *Nature* **250**, 120–123.
- Taylor W.R. and Green N.M. (1989) *FEBS* **179**, 241–248.
- Teem J.L., Berger L.S., Ostedgaard D.P., Rich D.P., Tsui L.-C. and Welsh M.J. (1993) *Cell* **73**, 335–346.
- Thomas P.J., Pedersen P.L. (1993) *J. Bioener. & Biomem.* **25**:11–20.
- Tsui L.-C. (1992) *Trends of Genetics* **8**, 392–398.
- Walker J.E., Saraste M., Runswick M.J. and Gay N.J. (1982) *EMBO J.* **1**, 945–951.
- Weber J. and Senior A.E. (1997) *Biochem. Biophys. Acta* **1319**, 19–58.
- White S. (1994) Ed. S. White, Oxford University Press, New York, 97–124.
- Yu X., Egelman E.H. (1997) *Nat. Struct. Biol.* **4**, 101.