

## FOR THE RECORD

# A diverse superfamily of enzymes with ATP-dependent carboxylate–amine/thiol ligase activity

MICHAEL Y. GALPERIN AND EUGENE V. KOONIN

National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland 20894

(RECEIVED June 12, 1997; ACCEPTED September 8, 1997)

**Abstract:** The recently developed PSI-BLAST method for sequence database search and methods for motif analysis were used to define and expand a superfamily of enzymes with an unusual nucleotide-binding fold, referred to as palmate, or ATP-grasp fold. In addition to D-alanine-D-alanine ligase, glutathione synthetase, biotin carboxylase, and carbamoyl phosphate synthetase, enzymes with known three-dimensional structures, the ATP-grasp domain is predicted in the ribosomal protein S6 modification enzyme (RimK), urea amidolyase, tubulin-tyrosine ligase, and three enzymes of purine biosynthesis. All these enzymes possess ATP-dependent carboxylate-amine ligase activity, and their catalytic mechanisms are likely to include acylphosphate intermediates. The ATP-grasp superfamily also includes succinate–CoA ligase (both ADP-forming and GDP-forming variants), malate–CoA ligase, and ATP–citrate lyase, enzymes with a carboxylate–thiol ligase activity, and several uncharacterized proteins. These findings significantly extend the variety of the substrates of ATP-grasp enzymes and the range of biochemical pathways in which they are involved, and demonstrate the complementarity between structural comparison and powerful methods for sequence analysis.

**Keywords:** ATP binding site; ATP-grasp fold; biotin carboxylase; glutathione synthetase; purine biosynthesis; succinate thiokinase; tubulin-tyrosine ligase

With the rapid accumulation of three-dimensional (3D) protein structures and the complementary development of structure-to-structure comparison methods, there has been lately a remarkable growth in the number of protein superfamilies delineated through structural conservation alone, in the absence of detectable sequence similarity (Holm & Sander, 1996). One of such structural superfamilies unites two groups of peptide synthetases, namely D-alanine:D-alanine ligase (DD-ligase) and glutathione synthetase (GSHase), with biotin carboxylases (BCases) and carbamoyl phosphate synthase (Fan et al., 1995; Artymiuk et al., 1996; Thoden

et al., 1997). All these enzymes catalyze a reaction that involves an ATP-dependent ligation of a carboxyl group carbon of one substrate with an amino or imino group nitrogen of the second one and includes, in each case, the formation of acylphosphate intermediates (Gushima et al., 1983; Ogita & Knowles, 1988; Meister, 1989; Fan et al., 1994). Structural alignment of DD-ligase, GSHase, and BCCase revealed three conserved motifs, corresponding to the phosphate-binding loop and the  $Mg^{2+}$ -binding site of the ATP-binding domain (Artymiuk et al., 1996). In each of these enzymes, ATP binds in a cleft formed by two structural elements, each containing two antiparallel  $\beta$ -strands and a loop (Hibi et al., 1996). A similar ATP-binding fold, referred to as GSHase fold, palmate ( $\beta$ -sheet) fold (Yamaguchi et al., 1993), or ATP-grasp fold (Murzin, 1996) has been detected in succinyl–CoA synthetase (SCS) (Wolodko et al., 1994; Matsuda et al., 1996). Sequence similarity with BCases indicates that this superfamily additionally includes the biotin-dependent carboxylase domains of pyruvate carboxylase and propionyl–CoA carboxylase (Artymiuk et al., 1996). Here, using recently developed sensitive methods for sequence database search and sequence motif analysis, we further expand the ATP-grasp superfamily to include the enzyme involved in ribosomal protein S6 modification, urea amidolyase, tubulin-tyrosine ligase, three enzymes of purine biosynthesis, and several uncharacterized proteins. These findings significantly extend the range of the biochemical pathways, in which ATP-grasp enzymes are involved, the variety of their substrates, and emphasize the complementarity between structural comparison and powerful methods for sequence analysis.

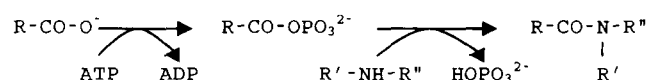
In the course of detailed comparative analysis of the protein sequences encoded in complete bacterial and archaeal genomes (Koonin et al., 1997), we observed that the *Escherichia coli* RimK protein, which is involved in post-translational modification of the ribosomal protein S6 (Reeh & Pedersen, 1979; Kang et al., 1989), had highly conserved homologs in all completely sequenced bacterial and archaeal genomes and also showed a significant similarity to GSHases. When the non-redundant protein sequence database at the National Center for Biotechnology Information was searched using BLASTGP program, which is an extension of the BLAST method (Altschul et al., 1990) incorporating statistical analysis of local alignments with gaps (Altschul & Gish, 1996;

Reprint requests to: Michael Y. Galperin, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland 20894; e-mail: galperin@ncbi.nlm.nih.gov.

Altschul et al., 1997), alignments of the RimK sequence with GSHases were detected with a probability of occurring by chance,  $P < 10^{-8}$ . We further iterated this search using the recently developed PSI-BLAST (Position-Specific Iterative BLAST) program, which converts local alignment produced by BLASTGP into position-specific weight matrices that are then used for iterative database scanning (Altschul et al., 1997). This search detected the known proteins with the ATP-grasp fold, GSHases, DD-ligases, and BCases, with a high statistical significance ( $P < 10^{-4}$ ). It also revealed, at the same significance level, a similar domain in urea amidolyase, phosphoribosylamine-glycine ligase, phosphoribosylglycinamide formyltransferase, and phosphoribosylaminoimidazole carboxylase. In addition, marginal similarity was detected with the sequences of SCS and tubulin-tyrosine ligase (TTL). Finally, when the alignment block containing the phosphate-binding site with a flexible glycine-rich loop flanked by two anti-parallel  $\beta$ -strands from these enzymes (residues 137–157 in DD-ligase) was used in a motif search using MoST program (Tatusov et al., 1994), a total of 125 different sequences were retrieved, of which 122 were considered members of the ATP-grasp superfamily.

The extended ATP-grasp superfamily currently includes 15 groups of enzymes, catalyzing ATP-dependent ligation of a carboxylate-

containing molecule to an amino or thiol group-containing molecule (Table 1). The list of reactions catalyzed by these enzymes demonstrates their flexibility with respect to both carboxyl and amino/thiol group-containing substrates. Thus, phosphoribosylglycinamide formyltransferase uses formic acid as a substrate, showing that the moiety at the carboxyl group can be as simple as H atom. On the other hand, in case of RimK and TTL, the carboxyl-containing substrates are proteins. In carbamoyl phosphate synthetase, the amino group containing substrate is simply ammonia (derived from glutamine), while in biotin carboxylases this substrate is N' atom of enzyme-bound biotin molecule. This shows that primary and secondary amines can both be used by enzymes of this family. The reaction catalyzed by ATP-dependent carboxylate-amine ligases can be summarized as follows:



In this scheme, R can be a hydrogen atom, hydroxyl group, an organic molecule, or even a protein; R' can be either a hydrogen atom or a part of a biotin ring, and R'' can be an amino-group

**Table 1.** Carboxylate-amine/thiol ligases containing ATP-grasp domains

Enzyme	Function or pathway	SWISS-PROT symbol	Active form	Carboxylate substrate	Amine or thiol substrate
Ribosomal protein S6 modification protein <sup>a</sup>	Ribosome biogenesis	RIMK_ECOLI	Monomer?	Ribosomal protein S6	Glutamate
Glutathione synthetase (EC 6.3.2.3)	Glutathione biosynthesis	GSHB_ECOLI	Tetramer	$\gamma$ -Glutamyl-cysteine	Glycine
D-Alanine-D-alanine ligase (EC 6.3.2.4)	Peptidoglycan biosynthesis	DDLA_ECOLI DDLB_ECOLI	Dimer	D-Alanine	D-Alanine
Phosphoribosylamine-glycine ligase <sup>a</sup> (EC 6.3.4.13)	Purine biosynthesis	PUR2_ECOLI	Monomer	Glycine	5-Phosphoribosylamine
Phosphoribosylglycinamide formyltransferase <sup>a</sup> (EC 2.1.2.-)	Purine biosynthesis	PURT_ECOLI	Monomer	HCOO <sup>-</sup>	5'-Phosphoribosylglycinamide
Phosphoribosylaminoimidazole carboxylase <sup>a</sup> (EC 4.1.1.21)	Purine biosynthesis	PURK_ECOLI PUR6_YEAST	Dimer	HCO <sub>3</sub> <sup>-</sup>	5'-Phosphoribosyl-5-aminoimidazole
Acetyl-CoA carboxylase, biotin carboxylase subunit (EC 6.3.4.14)	Fatty acid biosynthesis	ACCC_ECOLI COAC_YEAST	Heterohexamer Tetramer	HCO <sub>3</sub> <sup>-</sup>	Biotin-enzyme
Propionyl-CoA carboxylase (EC 6.4.1.3)	Amino acid catabolism	PCCA_HUMAN	Heterodimer	HCO <sub>3</sub> <sup>-</sup>	Biotin-enzyme
Pyruvate carboxylase (EC 6.4.1.1)	Gluconeogenesis	PYC_HUMAN	Tetramer	HCO <sub>3</sub> <sup>-</sup>	Biotin-enzyme
Urea amidolyase <sup>a</sup> (EC 6.3.4.6)	Urea hydrolysis	DURI_YEAST	Monomer	HCO <sub>3</sub> <sup>-</sup>	Biotin-enzyme
Carbamoyl-phosphate synthetase, large chain (EC 6.3.5.5)	Arginine biosynthesis pyrimidine biosynthesis	CARB_ECOLI PYR1_HUMAN	Heterodimer Hexamer	HCO <sub>3</sub> <sup>-</sup> NH <sub>2</sub> COO <sup>-</sup>	NH <sub>3</sub> —
Tubulin-tyrosine ligase <sup>a</sup> (EC 6.3.2.25)	Microtubules assembly	TTL_PIG	Monomer	$\alpha$ -Tubulin	Tyrosine
Succinyl-CoA synthetase, $\beta$ subunit (EC 6.2.1.5, 6.2.1.4)	Citric acid cycle	SUCC_ECOLI SUCB_PIG	Heterotetramer Heterodimer	Succinate	Coenzyme A
Malate-CoA ligase, $\beta$ subunit <sup>a</sup> (EC 6.2.1.9)	Growth on C-1 compounds	MTKB_METEX	Heterodimer	Malate, succinate	Coenzyme A
ATP-citrate lyase <sup>a</sup> (EC 4.1.3.8)	Lipid biosynthesis	ACLY_HUMAN	Monomer	Citrate	Coenzyme A

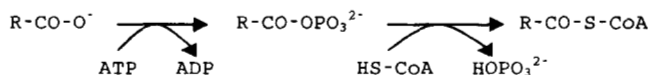
<sup>a</sup>Newly identified members of the superfamily.

containing a molecule, or a part of a biotin ring (Table 1). In several cases, substrates of ATP-grasp enzymes do not have an amino group. One of such enzymes is the enterococcal DD-ligase (vancomycin-resistance protein), which ligates D-alanine with D-lactate (Fan et al., 1994; Evers et al., 1996). Another variant of

the same catalytic mechanism works in succinyl-CoA synthetase, where it is the thiol group of HS-CoA that performs the nucleophilic attack on the succinyl phosphate intermediate. This probably also occurs in related enzymes, such as malate-CoA ligase and ATP-citrate lyase (Wells, 1991):

RIMK_ECOLI	88	arDKRr	31	gGAPLVVklveG-tGciGVVLA	19	ILVQBYIkeagqCDI	22	rsnlhrGGaasv	20	dVAGVDLlr	4	PLVBEVNASPGlegi	30
RIMK_HAEIN	104	arDKWk	30	ISSPPIIKktlng-sGciGVILA	19	MLGQDFIeeagnADI	22	ranchrGGktek	20	dVAGVDLlr	4	LLVLBEVNASPGlemi	25
Y011_MYCGE	93	anDKYe	23	ksFPVIVKkrns-hGckdVHLV	16	WIVQPFLLs-igtVEY	22	kanfsqGAevsl	20	GYAIDPFL	5	VIVBEVBDaAgaral	31
Y012_MYCGE	27	adNKgl	43	MeFPVIVKsvfG-sfcdyVFLC	16	AIVQKYIt-cskgEs	22	rsnlngGakaer	20	FYCGIDPLF	5	LIFCEVnpvqltrs	16
MJ0620	93	asNKrl	29	kfeeAVLkPifG-CGGeGIVrV	22	FYIQEFIk-pvrnEh	24	knmvsqGGrvek	20	FYAGVDLle	4	LKVLBEVNSTPswiGl	22
MJ1001	95	tsDRKf	32	LrFPVIVKnsfs-kCclkvFMA	17	KLIQEFId-fkenLD	24	rtnllylGnvvek	20	VILGVDLlp	4	YYVIBENSSPgtkGf	21
			125	160	167	198	206	273					
GSBH_ECOLI	121	cnEKlf	26	hs-DIILKPlDG-mGCASIFrV	21	CMANVTYPaiKGDK	23	rgnlaaGGrgep	23	IFVGLDIIg	-	drLTEENvtSptCir	25
GSBH_ANASP	129	anEKMy	26	kg-ATVLRPlGn-kACEGLFL	21	VMVOTYLPearKGDK	22	rnmmatGGTvak	23	IFVGLDVIg	-	gYLTBEVntSptGir	26
PUR2_ECOLI	101	egsKaf	30	kgAPIVIKAdG1-AAGkGVIVA	25	IVIEBFLD-GBEASF	24	dtGpntGGmGay	35	AGLMDkqg	1	pkVIBENCRFGdpet	133
PUR2_YEAST	105	easKaf	30	tdkAFVIVKAdG1-AAGkGVIP	26	VVIEQFLE-GDEISL	24	dklntGGmGay	40	MILVKdskt	4	peVLEBNVRFGdpet	494
PUR2_ARATH	212	egsKnf	30	qGAPIVIKAdG1-AAGkGVIVA	25	VVIEBFLD-GBEASF	24	dtGpntGGmGay	38	AGLMEIEkks	2	LSFIBENVRFGdpec	120
PUR2_HUMAN	103	essKrf	31	dFPALVVKAsG1-AAGkGVIVA	25	IVIEBLLD-GBEVS	24	dgGpntGGmGay	38	AGIML-tnk	1	pkVLEBNCRFGdpec	712
PURT_ECOLI	109	tmNRGf	31	IGYPCIVKpVms-SSckGqTPI	23	VIVGQVVKfdFETL	16	hrqgdGdyreSw	27	GLFGVBLFV	4	VIFSEVspRPhdtGm	103
PURT_BACSU	102	tmDRGf	31	IGFPVIVKpLms-SSckGqSVC	23	VVIEBFLP-fesEITL	16	hvgkdGdyieSw	27	GLFGVBLFV	4	VYFSEVspRPhdtG1	103
PURK_ECOLI	86	iadRlt	30	LGELAIIVKrrtGGyGdGrGqWrL	15	CIVGQGINfsgEVS	19	hqgdilrtsvAf	25	GVMAMECVF	4	LLINBEVlspRVhmsGh	107
PURK_BACSU	104	tQNRt	30	LrLPAVLKtcrGGyGdGqGqVFI	17	CILBSWVsfkmEVS	19	hknmlfqsivP	24	GpLAVEMFL	5	LLVnBEVlspRPhnsGh	102
PUR6_YEAST	102	iqDRYi	32	LGFPFVIVKrsrtlaydGrGnFVV	17	VLYBKWAPfktEGLV	19	hkdnicdlCyAp	25	GTFGVEMFY	5	LLINBEVlspRPhnsGh	292
			97	144	150	180	187	215	270	273			
DDLB_ECOLI	93	smDKRl	37	LGLPVIVKPSre-GSSvGMSkV	19	VLIIEKVLs-GPEEIV	19	ydyeaKRYLsdet	30	GNRGRDVML	5	FYLLBEVntsPgmtsh	26
DDLA_ECOLI	137	cmDKdv	34	LGLPLFVKPAnq-GSSvGVskV	12	AIVEQGIK-GRBIEC	22	yaydtKVIDedg	31	GMARVDVFL	5	VVINBEVntlPgftni	39
VANA_ENTFC	129	cmDKsl	28	FtYVVFVKPars-GSSfGvKkV	19	ALIEQAVs-GCEVCG	27	qevepekgSena	30	GLARVDMFL	5	IVLBEVntlPgftsy	28
			116	159	165	201	209	236	241	288	294		
ACCC_ECOLI	112	mgDKVs	33	IGYPVIVKASGG-GGGrGMRVV	25	VYMEKYLBNPRHIEI	20	QRRHQRVVEAP	26	GAGTPEFLF	4	FYFBEVntRlQVSH	152
ACCC_ANASP	113	mgDKst	32	IGYPVIVKATAG-GGGrGMRVL	25	VYIEKPIERPRHIEI	20	QRRNQKLIBEAP	26	GAGTPEFLF	5	FYFBEVntRlQVSH	148
ACCC_METJA	112	mgSKfn	32	IGFPVIVKASAG-GGGMMSVA	25	VPIEKYLENPRHIEI	20	QRRHQKLIBEAP	26	SAGTVEFLY	4	FYFBEVntRlQVSH	204
COAC_YEAST	182	lqDKIs	56	IGFPVIVKAsEG-GGCKGIRqV	21	IPIMKLAGrARHIEV	20	QRRHQKLIBEAP	26	SAGTVEFLY	6	FYFBEVntRlQVSH	1844
COAC_RAT	238	lqDKIA	58	IGYVIVKAsEG-GGCKGIRqV	21	IPVnMLAKQSRHIEV	20	QRRHQKLIBEAP	26	SAGTVEFLY	5	FYFBEVntRlQVSH	1899
PCCA_HUMAN	148	mgDKIE	32	IGYVIVKAsAG-GGCKGMRIA	25	LLIEKFIIDNPRHIEI	20	QRRNQRVVEAP	26	SAGTVEFLV	5	FYFBEVntRlQVSH	368
PYC1_YEAST	132	vgDKVs	32	IGYVIVKAAfG-GGGrGMRVV	25	CFVBRFLDKPKHIEV	20	QRRHQRVVEAP	26	nAGTAEFLV	5	HYFBEVntRlQVSH	860
PYC_HUMAN	148	mgDKVE	32	YGFPVIVKAAfG-GGGrGMRVV	25	LFVBEKFIKPKHIEV	20	QRRHQRVVEAP	26	nAGTAEFLV	5	HYFBEVntRlQVSH	844
DURI_YEAST	743	lglKhs	31	LeYVIVKStAG-GGCIgLqkV	25	VFLERFIENARHIEV	20	QRRNQRVVEAP	26	CAGTVEFLY	6	FYFBEVntRlQVSH	906
			129	169	215	208	215	243	285	299			
CARB_ECOLI	125	aeDRrr	30	VGFPCIVRPSft-MGSGGGIA	21	LLIDRSILigwKEVEM	22	GHTGDSITVAP	27	GGSNVQFV	6	LIVBEVnDpVSRSSA	764
CARB_ECOLI	671	aeDRer	30	IGYPLVVRPSYv-LGcraMEIV	21	VLLIDRFLDdaveVDD	21	GVHSGDSACSLP	26	GLMNVOFV	4	VYLLBEVnDpAARTvP	222
MJ1378	127	aeDRl	30	IGYPLVVRPAft-LGctGGGIA	21	VLLIDRSVLIgweFEL	22	GHTGDSITVAP	26	GGCNIOFV	4	YrVBEVnDpVSRSSA	175
MJ1381	200	aeDRee	30	IGYPLVVRPSYv-LGcraMqIV	21	VLLIDKFLDdaieLDV	21	GVHSGDSATVIP	26	GLLNVOYAV	4	VYVBEVnDpRASrtvP	238
CARB_YEAST	146	aeDRdl	30	VkYVIVRAYSAY-LGclGSGFA	19	VLMVRKSLKgwKEVY	22	GVHSGDSMVFP	26	GeCNVOYAL	6	YrVBEVnDpVSRSSA	791
CARB_YEAST	690	aeNRhk	30	VnYVIVRAYSAY-LSCAAMSvV	21	VVMSRFILGgaQELDV	21	GVHSGDSMLVLP	26	GpFNQIILK	5	LkVBEVnDpVSRSSA	247
PYR1_YEAST	554	teDRel	30	IGFPVIVRAAY-LGclGSGFA	19	VLMVRKSMKgwKEVY	22	GHTGDSITVAP	22	lGVVGCeNI	10	YCIBEVnDpVSRSSA	1479
PYR1_YEAST	1091	aeNRyk	30	IGYVIVRAYSAY-LSCAAMntV	21	VVITKYIEEnaKEIEM	21	GVHSGDATLIVP	26	GpYNQOFTA	4	TkVBEVnDpVSRSSA	943
PYR1_HUMAN	511	teDRra	30	VGYVIVRAAFA-VGclGSGFA	19	VLLDKSLKgwKEVY	22	GHTGDSITVAP	22	lGVVGCeNI	10	YCIBEVnDpVSRSSA	1533
PYR1_HUMAN	1044	aeNRfk	30	VGYVIVRAYSAY-LSCAAMnVA	21	VVISKFIQeaKEIDV	21	GVHSGDATLIVP	26	GpFNQIILK	4	LkVBEVnDpVSRSSA	1001
			9	46	52	99	107	176					
TTL_PIG	66	gaDKLc	70	egnVVIARSSAG-AKcEGILIS	16	hVIOKYLRLPRLLEP	23	eGVLRTasepYh	83	qLFGdFMV	5	VWLEBEVNGAPACAqk	38
TTL_BOVIN	66	gaDKLc	70	egnVVIARSSAG-AKcEGILIS	16	hVIOKYLRLPRLLEP	23	eGVLRTasepYh	83	qLFGdFMV	5	VWLEBEVNGAPACAqk	36
ZK1128_6*	242	rkDRlw	7	asrhVIVRPPAS-ARctGISVT	11	LVAQHYIERPLTInr	23	qGLVRFASvYs	88	eLFGDIIIL	5	pWLEBEVnDpVSRSSA	156
C55A6_2*	660	kkDRly	49	fpgeFIVKPtNS-rQCKGIFFA	10	LLVSRYLKDPYLVnn	23	eGLARlAsrYd	94	eLFGdDVLV	5	pWLEBEVnDpVSRSSA	283
KIAA0173	669	rkDRlw	38	srqkWiVIRPPAS-ARctGIQVI	11	LLVORYLEKPKYLIsg	23	dGLVRFASckYs	92	eLFGdDML	5	pWLEBEVnDpVSRSSA	283
YBU4_YEAST	390	dmDKlm	75	qdkWIVRPPAS-dKcGgGIRVF	50	FIIQBYLLENPLLLas	24	rMLALFAakpFv	82	eTYGVDLIL	5	VKLEBEVnDpVSRSSA	50
D2013_9*	390	vkDlLa	37	qhnVVIKRPwnl-ARcmDMTvt	14	kIVcEYIPRPLLFpr	29	rFWIRfaineFs	77	aMYGVDIIL	9	STLBEVnDpVSRSSA	28
KIAA0153	360	vkDcLa	39	ednhWICRPwnl-ARSlDTHVT	14	kVVSRYIEsPVLFR	28	vFWRlRfsnraFa	87	wdnGpDgrr	3	pqIIEVnDpVSRSSA	29
			9	46	52	99	107	176					
SUCC_ECOLI	5	YQAKQL	27	GAGPVVVkrqcvh-AGGRGKAG	36	ILVBEaatDiakELYL	19	GVEIEKVAeetp	29	GkLVOcftk	13	LALBEVnDpVSRSSA	181
SUCB_PIG	27	YQSKKL	27	nakeIVLkraqil-ACGRGKGVf	44	MVAbal-DisrETYL	19	GVDIEEVAAsnp	29	GpLqNCAad	13	ATqVBEVnDpVSRSSA	181
MTKA_METEX	5	YQAKEL	27	GGsfWVRaqih-ACGRGKAG	36	VYVBEtadpferELYL	19	GMDIEEIAAkep	26	GlnIKQVSA	16	GTMPEVnDpVSRSSA	183
ACL_Y_HUMAN	8	qtCKEL	36	lsqnLVVkpqdl-ikrRGKLG1	33	FLLEp-5-gaeEFVY	17	GVDVGDVAKaq	23	pdkkKILA	14	FTYIEVnDpVSRSSA	894
			9	46	52	99	107	176					
MJ0776	123	isNKYk	29	efktCILRPIyGSGGsilKtEL	14	IIA0EYIrgksfsn	17	kGMyaGnltpyi	20	GMSGIDFLI	4	pYIVBEVnDpVSRSSA	99
MJ0815	135	agnKY1	15	ppkkYVVKkids-CGCKfnLfd	2	FLIQEYIDgenfsn	19	rgfvGgevninh	21	GYVGDVIV	4	YIIEBEVnDpVSRSSA	34
slr1616	112	csNKkr	26	iglPCIVRPsstGtGGSvVFFA	19	pIA0EYID-ineGef	30	rgGIVeissys	24	GpInICarV	4	LMPEBEVnDpVSRSSA	48
YSC19931	309	lrakLR	28	mtkPFVVRPvdGed--hnIYIU	35	UYEQFMdtdnfeVd	24	rrnthGkevyri	20	MICGDELLR	4	sYVBEVnDpVSRSSA	649
Consensus		..\$+u.		uGupUUU+pooc.OGG.guUUU		UUU-\$UU\$. \$EU-U		....\$g.....		ouUUU-UUU		UUUUUEUNuru.....	
Structure				ββββββ		ββββ		βββ		ββββ βββ		ββββ	

**Fig. 1.** Multiple alignment of the conserved regions in the ATP-grasp proteins. The proteins are listed under their unique SWISS-PROT or GenBank identifiers; proteins from *Caenorhabditis elegans* are marked with an asterisk. The numbers indicate distances to the ends of each protein and the sizes of the gaps between aligned segments. The names of proteins with known 3D structures are shown in bold; the positions of the conserved residues are indicated above such residues. Red shading indicates the amino acid residues that were shown to be involved in ATP binding by X-ray analysis; corresponding amino acid residues in other sequences are shaded blue. Magenta shading indicates the residues identified by site-specific mutagenesis. Conserved amino acid residues of the active center of biotin carboxylases are shaded green. Yellow shading indicates uncharged amino acid residues (A, I, L, V, M, F, Y, or W) with a propensity to form  $\beta$ -strands. Conserved small residues (G, A, S, or C) are shown in green, the residues conserved within a protein family are in caps, the ones conserved among several protein families are in bold. The consensus includes amino acid residues conserved in all sequences (upper case) and those conserved in the majority of the sequences (lower case). U stands for a bulky hydrophobic residue (I, L, V, M, F, Y, W), O stands for a small residue (G, A, S, C), + stands for K or R, - stands for D or E, \$ indicates any charged residue (D, E, K, R, N, Q), and dot stands for any residue.



However, these exceptions appear to be the only ones where an ATP-grasp enzyme works without an amino group-containing substrate. Remarkably, SCS of *E. coli* can use GTP and ITP in addition to ATP (Murakami et al., 1972; Kelly & Cha, 1977), while the enzymes from pig hearts and *Dictyostelium discoideum* are GTP specific (discussed by Nishimura, 1986; Anschutz et al., 1993).

The multiple alignment of the ATP-grasp superfamily generated from the PSI-BLAST data using CLUSTALW (Thompson et al., 1994) and MACAW (Schuler et al., 1991) programs (Fig. 1) shows a clear pattern of conservation in the three motifs described by Artymiuk et al. (1996). These motifs are also detectable in the sequences of SCS, in accordance with the structural classification (Murzin, 1996), and TTL, confirming that TTL contains an ATP-grasp domain (Fig. 1). The alignment also includes several uncharacterized proteins, for which no function could be deduced from the sequence data.

This alignment demonstrates that most of the amino acid residues that interact with ATP in DD-ligase (Fan et al., 1994, 1995) are conserved in all ATP-grasp domains, and the few allowed substitutions are consistent with their predicted role in ATP binding. Thus, Lys-97 and Lys-144 of DD-ligase, interacting electrostatically with  $\alpha$ - and  $\beta$ -phosphates of ATP, are conserved in a majority of sequences and are only substituted by Arg in carbamoyl phosphate synthetases. Glu-180, hydrogen bonded to the amino

group of adenine, can be substituted by Gln or Asp, while Glu-187, hydrogen bonded to ribose OH groups, can be substituted by His, Asp, or Asn. Amino acid residues Asp-257, Glu-270, and Asn-272, participating in coordination of  $\text{Mg}^{2+}$ , and Trp-182, Leu-183, and Leu-269, providing hydrophobic interactions for adenine and ribose rings, are also highly conserved. In accordance with the site mutagenesis data (Shi & Walsh, 1995), Tyr-216 can be substituted by other hydrophobic residues. On the other hand, most of the amino acid residues forming the active center of BCCase (Waldrop et al., 1994), such as Tyr-82, His-236, Lys-238, Glu-241, Gln-294, and Glu-296, are conserved only among BCases. These residues are not conserved even in phosphoribosylaminoimidazole carboxylase (Fig. 1), which also uses bicarbonate as a substrate, but is likely biotin-independent (Mueller et al., 1994). Finally, we found no sequence conservation around Glu-15, which in DD-ligase binds amino groups of the substrate. This may reflect the variety of the amine-containing substrates of the enzymes of this family.

The ATP-grasp domain is ubiquitous, with multiple representatives of the superfamily encoded in each of the completely sequenced genomes (Table 2). It is of interest, however that, unlike metabolic enzymes, the two groups of ATP-grasp proteins involved in protein modification, while highly conserved, are limited in their phylogenetic distribution to prokaryotes (RimK) or eukaryotes (TTL) (Table 2). Conceivably, these enzymes could have evolved from ancestral metabolic enzymes. Phylogenetic distribution of the ATP-grasp enzymes suggests some interesting functional clues. *Methanococcus jannaschii* encodes two paralogous proteins that are both orthologous to RimK, yet there is no gene for

**Table 2.** ATP-grasp domains encoded in prokaryotic and eukaryotic genomes

Organism	RimK	Glutathione synthetase	D-Alanine-D-alanine ligase	Purine biosynthesis enzymes			Biotin carboxylases	Carbamoyl phosphate synthetase	Tubulin-tyrosin ligase	Unknown
<i>Escherichia coli</i>	RimK	GshB	DdIA DdIB	PurD	PurT	PurK	AccC	CarB	—	—
<i>Hemophilus influenzae</i>	HI1531	—	HI1140	HI0888	—	HI1616	—	—	—	—
<i>Mycoplasma genitalium</i>	MG011 MG012	—	—	—	—	—	—	—	—	—
<i>Synechocystis</i> sp.	—	slr1238 slr2002	slr1874	slr1159	slr0861	slr0578	slr0053	slr0370	—	slr1616
<i>Methanococcus jannaschii</i>	MJ0620 MJ1001	—	—	MJ0937	MJ1486	—	MJ1229	MJ1378 MJ1381	—	MJ0776 MJ0815
<i>Saccharomyces cerevisiae</i>	—	— <sup>b</sup>	—	YGL234w	—	YOR128c	YBR208c, YBR218c, YGL062w, YNR016c, YM8261.01c	YJL130c YJR109c	YBR094w	YSCL9931
<i>Caenorhabditis elegans</i>	—	— <sup>b</sup>	—	F38B6.4	—	—	D2023.2, F26D10.2, F27D9.5, F32B6.e	D2085.1	C55A6.2 ZK1128.6	F46F11.1
<i>Homo sapiens</i>	—	— <sup>b</sup>	—	P22102	—	—	P05165, P11498, S41121	P27708 P31327	KIAA0153 KIAA0173	—

<sup>a</sup>Proteins are indicated by their original authors' designations; *E. coli* proteins are from SWISS-PROT database, *H. influenzae*—from Fleischmann et al. (1995), *M. genitalium*—from Fraser et al. (1995), *Synechocystis* sp.—from Kaneko et al. (1996), *M. jannaschii*—from Bult et al. (1996), *S. cerevisiae*—from the *Saccharomyces* genome database (Stanford University). Yeast, human, and worm proteins are listed under their GenBank identifiers.

<sup>b</sup>Eukaryotic glutathione synthetase has no detectable sequence similarity with the bacterial enzyme.

the bacterial ribosomal protein S6 in the *M. jannaschii* genome. A similar pattern can be seen in the recently sequenced genomes of two other archaea, *Methanobacterium thermoautotrophicum* and *Archaeoglobus fulgidus*. Thus, it seems likely that RimK actually has another, conserved, but not yet discovered activity. Of further interest is the finding of a group of uncharacterized eukaryotic proteins containing an ATP-grasp domain and showing the highest, though limited, similarity to RimK. These putative enzymes may be involved in an as yet unknown protein modification mechanism.

A substantial number of enzymes that are similar in function to ATP-dependent carboxylate-amine ligases do not show any detectable sequence similarity to the ATP-grasp superfamily. These include such peptide synthetases as  $\gamma$ -glutamyl-cysteine synthetase, eukaryotic GSHase, and bacterial peptidoglycan biosynthesis proteins MurC, MurD, and MurE. MurD was shown recently to have a typical Rossmann fold (Bertrand et al., 1997), rather than ATP-grasp. Pyruvate phosphate dikinase, which reportedly has the ATP-grasp fold (Herzberg et al., 1996; Murzin, 1996), shows no apparent sequence similarity to DD-ligase, GSHase, BCCase, or SCS. Finally, in spite of certain similarities in the reaction mechanism, there is no indication that glutamine synthetases belong to the ATP-grasp superfamily.

On a general note, the findings presented here show that with refinement of sequence comparison methods their sensitivity may match that of methods based on structure comparison. Such developments are particularly important, given the parallel rapid growth of sequence and structure databases as sequence analysis methods complement structure analysis by assigning subtly similar sequences to superfamilies with known folds.

**Acknowledgments:** We thank Stephen Altschul, David Lipman, Tom Madden, and Alejandro Schaffer for providing us with the pre-release version of the program PSI-BLAST, and James Knox and Alex Kuzin for helpful discussion.

## References

- Altschul SF, Gish W. 1996. Local alignment statistics. *Methods Enzymol* 266:460–480.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* 215:403–410.
- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucl Acids Res* 25:3389–3402.
- Anschutz AL, Um HD, Siegel NR, Veron M, Klein C. 1993. P36, a *Dicystostelium discoideum* protein whose phosphorylation is stimulated by GDP, is homologous to the alpha-subunit of succinyl-CoA synthetase. *Biochim Biophys Acta* 1162:40–46.
- Artymiuk PJ, Poirrette AR, Rice DW, Willett P. 1996. Biotin carboxylase comes into the fold. *Nat Struct Biol* 3:128–132.
- Bertrand JA, Auger G, Fanchon E, Martin L, Blanot D, van Heijenoort J, Dideberg O. 1997. Crystal structure of UDP-N-acetylmuramoyl-L-alanine:D-glutamate ligase from *Escherichia coli*. *EMBO J* 16:3416–3425.
- Bult CJ, White O, Olsen GJ, Zhou L, Fleischmann RD, Sutton GG, Blake JA, FitzGerald LM, Clayton RA, Gocayne JD, et al. 1996. Complete genome sequence of the methanogenic archaeon, *Methanococcus jannaschii*. *Science* 273:1058–1073.
- Evers S, Casadewall B, Charles M, Dutka-Malen S, Galimand M, Courvalin P. 1996. Evolution of structure and substrate specificity in D-alanine:D-alanine ligases and related enzymes. *J Mol Evol* 42:706–712.
- Fan C, Moews PC, Walsh CT, Knox JR. 1994. Vancomycin resistance: Structure of D-alanine:D-alanine ligase at 2.3 Å resolution. *Science* 266:439–443.
- Fan C, Moews PC, Shi Y, Walsh CT, Knox JR. 1995. A common fold for peptide synthetases cleaving ATP to ADP: glutathione synthetase and D-alanine:D-alanine ligase of *Escherichia coli*. *Proc Natl Acad Sci USA* 92:1172–1176.
- Fleischmann RD, Adams MD, White O, Clayton RA, Kirkness EF, Kerlavage AR, Bult CJ, Tomb JF, Dougherty BA, Merrick JM, et al. 1995. Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd. *Science* 269:496–512.
- Fraser CM, Gocayne JD, White O, Adams MD, Clayton RA, Fleischmann RD, Bult CJ, Kerlavage AR, Sutton G, Kelley JM, et al. 1995. The minimal gene complement of *Mycoplasma genitalium*. *Science* 270:397–403.
- Gushima H, Miya T, Murata K, Kimura A. 1983. Purification and characterization of glutathione synthetase from *Escherichia coli* B. *J Appl Biochem* 5:210–218.
- Herzberg O, Chen CC, Kapadia G, McGuire M, Carroll LJ, Noh SJ, Dunaway-Mariano D. 1996. Swiveling-domain mechanism for enzymatic phosphotransfer between remote reaction sites. *Proc Natl Acad Sci USA* 93:2652–2657.
- Hibi T, Nishioka T, Kato H, Tanizawa K, Fukui T, Katsube Y, Oda J. 1996. Structure of the multifunctional loops in the nonclassical ATP-binding fold of glutathione synthetase. *Nat Struct Biol* 3:16–18.
- Holm L, Sander C. 1996. Alignment of three-dimensional protein structures: Network server for database searching. *Methods Enzymol* 266:653–662.
- Kaneko T, Sato S, Kotani H, Tanaka A, Asamizu E, Nakamura Y, Miyajima N, Hirose M, Sugiura M, Sasamoto S, et al. 1996. Sequence analysis of the genome of the unicellular cyanobacterium *Synechocystis* sp. strain PCC6803. II. Sequence determination of the entire genome and assignment of potential protein-coding regions. *DNA Res* 3:109–136.
- Kang WK, Icho T, Isono S, Kitakawa M, Isono K. 1989. Characterization of the gene rimK responsible for the addition of glutamic acid residues to the C-terminus of ribosomal protein S6 in *Escherichia coli* K12. *Mol Gen Genet* 217:281–288.
- Kelly CJ, Cha S. 1977. Nucleotide specificity of succinate thiokinases from bacteria. *Arch Biochem Biophys* 178:208–217.
- Koonin EV, Mushegian AR, Galperin MY, Walker DR. 1997. Comparison of archaeal and bacterial genomes: Computer analysis of protein sequences predicts novel functions and suggests a chimeric origin for archaea. *Mol Microbiol* 25:619–637.
- Matsuda K, Mizuguchi K, Nishioka T, Kato H, Go N, Oda J. 1996. Crystal structure of glutathione synthetase at optimal pH: Domain architecture and structural similarity with other proteins. *Protein Eng* 9:1083–1092.
- Meister A. 1989. Mechanism and regulation of the glutamine-dependent carbamyl phosphate synthetase of *Escherichia coli*. *Adv Enzymol Relat Areas Mol Biol* 62:315–374.
- Mueller EJ, Meyer E, Rudolph J, Davisson VJ, Stubbe J. 1994. N5-carboxyaminoimidazole ribonucleotide: Evidence for a new intermediate and two new enzymatic activities in the de novo purine biosynthetic pathway of *Escherichia coli*. *Biochemistry* 33:2269–2278.
- Murakami K, Mitchell T, Nishimura JS. 1972. Nucleotide specificity of *Escherichia coli* succinic thiokinase. Succinyl coenzyme A-stimulated nucleoside diphosphate kinase activity of the enzyme. *J Biol Chem* 247:6247–6252.
- Murzin AG. 1996. Structural classification of proteins: New superfamilies. *Curr Opin Struct Biol* 6:386–394.
- Nishimura JS. 1986. Succinyl-CoA synthetase structure-function relationships and other considerations. *Adv Enzymol Relat Areas Mol Biol* 58:141–172.
- Ogita T, Knowles JR. 1988. On the intermediacy of carboxyphosphate in biotin-dependent carboxylations. *Biochemistry* 27:8028–8033.
- Reeh S, Pedersen S. 1979. Post-translational modification of *Escherichia coli* ribosomal protein S6. *Mol Gen Genet* 173:183–187.
- Schuler GD, Altschul SF, Lipman DJ. 1991. A workbench for multiple alignment construction and analysis. *Proteins* 9:180–190.
- Shi Y, Walsh CT. 1995. Active site mapping of *Escherichia coli* D-Ala-D-Ala ligase by structure-based mutagenesis. *Biochemistry* 34:2768–2776.
- Tatusov RL, Altschul SF, Koonin EV. 1994. Detection of conserved segments in proteins: Iterative scanning of sequence databases with alignment blocks. *Proc Natl Acad Sci USA* 91:12091–12095.
- Thoden JB, Holden HM, Wesenberg G, Rauschel FM, Rayment I. 1997. Structure of carbamoyl phosphate synthetase: A journey of 96 Å from substrate to product. *Biochemistry* 36:6305–6316.
- Thompson JD, Higgins DG, Gibson TJ. 1994. CLUSTAL W: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22:4673–4680.
- Waldrop GL, Rayment I, Holden HM. 1994. Three-dimensional structure of the biotin carboxylase subunit of acetyl-CoA carboxylase. *Biochemistry* 33:10249–10256.
- Wells TN. 1991. ATP-citrate lyase from rat liver. Characterisation of the citryl-enzyme complexes. *Eur J Biochem* 199:163–168.
- Wlodko WT, Fraser ME, James MN, Bridger WA. 1994. The crystal structure of succinyl-CoA synthetase from *Escherichia coli* at 2.5-Å resolution. *J Biol Chem* 269:10883–10890.
- Yamaguchi H, Kato H, Hata Y, Nishioka T, Kimura A, Oda J, Katsube Y. 1993. Three-dimensional structure of the glutathione synthetase from *Escherichia coli* B at 2.0 Å resolution. *J Mol Biol* 229:1083–1100.