

Integrating the Universal Metabolism into a Phylogenetic Analysis

Chomin Cunchillos* and Guillaume Lecointre†

*Institut Charles Darwin International, Romainville, France; and †UMR7138 Systématique, Adaptation, Evolution, Département Systématique et Evolution, Muséum National d'Histoire Naturelle, Paris, France

The darwinian concept of “descent with modification” applies to metabolic pathways: pathways sharing similarities must have inherited them from an exclusive, hypothetical ancestral pathway. Comparative anatomy of biochemical pathways is performed using five criteria of homology. Primary homologies of “type I” were defined as several pathways sharing the same enzyme with high specificity for its substrate. Primary homologies of “type II” were defined as the sharing of similar enzymatic functions, cofactors, functional family, or recurrence of a set of reactions. Standard cladistic analysis is used to infer the evolutionary history of metabolic development and the relative ordering of biochemical reactions through time, from a single matrix integrating the whole basic universal metabolism. The cladogram shows that the earliest pathways to emerge are metabolism of amino acids of groups I and II (Asp, Asn, Glu, and Gln). The earliest enzymatic functions are mostly linked to amino acid catabolism: deamination, transamination, and decarboxylation. For some amino acids, catabolism and biosynthesis occur at the same time (Asp, Glu, Lys, and Met). Catabolism precedes anabolism for Asn, Gln, Arg, Trp, His, Tyr, and Phe, and anabolism precedes catabolism for Pro, Ala, Leu, Val, Ile, Cys, Gly, Ser, and Thr. The urea cycle evolves from arginine synthesis. Metabolism of fatty acids and sugars develops after the full development of metabolism of amino acids of groups I and II, and they are associated with the anabolism of amino acids of groups III and IV. Syntheses of aromatic amino acids are branched within sugar metabolism. The Krebs cycle occurs relatively late after the setting of metabolism of amino acids of groups I and II. One portion of the Krebs cycle has a catabolic origin, whereas the other portion has an anabolic origin in pathways of amino acids of groups III and IV. It is not possible to order glycolysis and gluconeogenesis with regard to the Krebs cycle, as they all belong to “period 6.” Pentose-phosphate and Calvin cycles are later (periods 7 and 8, respectively). Cladistic analysis of the structure of biochemical pathways makes hypotheses in biochemical evolution explicit and parsimonious.

Introduction

Metabolic pathways are series of successive biochemical reactions catalyzed by enzymes. As the product of enzymatic activity, they can be seen as a part of the phenotype. Most of the phenotypic attributes are inherited and classically analyzed through standard phylogenetic methods to depict the distribution of character states among the living world. Classically, attributes of organisms are coded into a matrix to express primary homologies (de Pinna 1991), to infer synapomorphies (i.e., secondary homologies) through tree reconstruction, and to obtain a relative ordering through time of branches and character changes. Attributes of metabolic pathways can be coded the same way to order enzymatic innovations through time. From these premises, comparative anatomy of metabolic pathways has been developed by Cunchillos and Lecointre (2000, 2002). The darwinian concept of descent with modification was shown to apply to metabolic pathways, legitimizing the use of cladistic analysis to infer the history of metabolic development. Phylogenetic analysis of amino acid metabolism was performed by defining a new kind of taxon, new characters, and four criteria of primary homology (Cunchillos and Lecointre 2002, 2003). To achieve this goal, a taxon was defined as a “pathway”; that is, the series of reactions from the tip of the pathway to its point of connection into the Krebs cycle. This definition is valid either for catabolism, where the tip of the pathway is the nutrient, or for anabolism, where the tip is the final

product. Characters were defined as presence of enzymes, type of reactions, or use of cofactors, according to four criteria of primary homology. Primary homologies of “type I” were defined as several pathways sharing the same enzyme with high specificity for its substrate. In such case, the name of the enzyme itself is the name of the character (a column in the matrix). Primary homologies of “type II” were defined as the sharing of similar enzymatic functions by different pathways (although the enzymes recognized can have different specificities), sharing the same cofactors or sharing the same functional family (see Cunchillos and Lecointre [2003] for details). Because the parts of metabolism analyzed are universal, the inferred history is located in time before the phylogenetic development of the species “tree of life” as we know it today.

Using these new concepts, the historical development of aliphatic amino acid catabolism was first inferred in relation to the development of the Krebs cycle, and then aliphatic amino acid anabolism was incorporated into the matrix. Obviously, such a methodological framework now has to incorporate the most widely shared pathways through life, called here “universal metabolism”; that is, anabolism and catabolism of the three fundamental kinds of biomolecules: amino acids, fatty acids, and saccharides. It is of interest because that framework offers new and transparent means to answer questions about metabolism evolution. One of them is the relative timing of the rise of glycolysis, the Krebs cycle, and amino acid biosynthesis. Meléndez-Hevia, Waddell, and Cascante (1996) considered glycolysis as earlier than both the amino acid biosynthesis and the Krebs cycle because glucose is implicitly considered as more profitable as the earliest nutrient than are amino acids. Amino acid catabolism was, therefore, thought by these

Key words: metabolism, metabolic pathways, biochemical evolution.

E-mail: lecointre@mnhn.fr.

Mol. Biol. Evol. 22(1):1–11, 2005

doi:10.1093/molbev/msh253

Advance Access publication September 8, 2004

authors to be secondary. Conversely, Cunchillos and Lecointre (2003) found the Krebs cycle as the product of metabolism of amino acids of groups I and II (in the sense of Cordón [1990], namely, Asp, Asn, Glu, Gln, Arg, and Pro). However, their matrix did not include glycolysis. For this reason, glycolysis and gluconeogenesis have been included here. Is the metabolism of monosaccharides (also including Calvin and pentose-phosphate cycles) the first? In the same framework, when did fatty acid metabolism rise, compared with other pathways? A number of other metabolic pathways have also been included in the present work to submit a greater diversity of widely shared reactions to a transparent procedure for inferring a more complete view of their temporal development. For instance, the development of the urea cycle is traditionally thought to be linked to the metabolism of arginine and might have been possible as early as arginine metabolism, and the metabolism of aromatic amino acids is classically thought as arising later than the appearance of aliphatic amino acid metabolism.

Earlier authors (e.g., Horowitz 1945; Cordón 1990) have speculated on the evolutionary timing of the development of pathways from a theoretical point of view. For example, because free aliphatic amino acids were considered one of the very first sources of molecules in abiotic environments, upstream reactions of amino acid catabolic pathways must have occurred before downstream reactions, whereas downstream reactions of amino acid anabolic pathways must have occurred before upstream reactions (see Cunchillos and Lecointre [2000]). The present methods have the power to test such hypotheses. By using those methods, those hypotheses were corroborated for some aliphatic amino acids only (Cunchillos and Lecointre 2003). The evolutionary timing of development of other pathways can now be tested in the same way. In the present work, new data increasing complexity of characters required to face new methodological challenges.

Materials and Methods

Taxonomic Sampling

The present work focuses on the metabolism of the three main kinds of universal biochemical compounds, namely, amino acids, fatty acids, and monosaccharides. Analysis of the metabolism of these compounds requires the coding of the following assemblages of pathways: Krebs cycle, Calvin cycle, pentose-phosphate cycle, urea cycle, fatty acid anabolism and catabolism, amino acid anabolism and catabolism, glycolysis, and gluconeogenesis. These pathways are shared by all living things, at least primitively, with few exceptions (Calvin cycle), and discussion of secondary losses in some species or species groups is beyond the scope of the present work.

Enzymes acting on complex molecules have been excluded from that universal core of enzymatic activities. Complex molecules are assemblages of those elementary molecules whose metabolic evolution is being studied. Complex molecules include, for instance, coenzymes, triglycerids, phospholipids, nucleosids, and polymers (e.g., glycogen, starch, proteins, DNA, and RNA). They all are secondary products of the core metabolism studied here. Purines and pyrimidines are considered as complex mol-

ecules themselves because they are never synthesized as such *in vivo*. Indeed, their precursors are already attached to other compounds (ribose or ribose-phosphates), so their recognition as isolated compounds have no biological basis. This delineation leads us to sample the following pathways for phylogenetic analysis.

As in Cunchillos and Lecointre (2000), taxa are defined from the tip of the pathway to its point of contact into the Krebs cycle. To name pathways, prefixes “d” and “s” are used to refer to degradation and synthesis, respectively. For example, dGLN is the set of enzymatic activities involved in converting glutamine to oxoglutarate, whereas sGLN is the synthetic pathway from oxoglutarate to glutamine. When degradation or synthesis of a compound can occur in different ways, they are numbered in the name of the pathway. For instance, cysteine can be degraded via mercaptopyruvate (dCYS2) or directly through pyruvate and acetyl-CoA (dCYS1). Taxons are listed in tables 1 and 2. Fatty acid catabolic and anabolic pathways stop at acetyl-CoA. Monosaccharide anabolic pathways (gluconeogenesis and pentose-phosphate cycle) stop at oxaloacetate, and monosaccharide catabolic pathways (glycolysis and Calvin cycle) stop at acetyl-CoA. Amino acid anabolism and catabolism stop at different points of the Krebs cycle or at acetyl-CoA, depending on the amino acid (oxoglutarate, succinyl-CoA, oxaloacetate, or acetyl-CoA). The Krebs cycle is divided into two parts, designated using the two main points of entrance: KC1 from oxaloacetate to oxoglutarate and KC2 from oxoglutarate to oxaloacetate. The urea cycle is not delineated in reference to the Krebs cycle, but independently, including the reactions of its own cycle.

Characters and Homologies

Characters are defined according to four criteria of homology (see Cunchillos and Lecointre [2003] for details). Primary homologies of “type I” were defined as several pathways sharing the same enzyme with high specificity for its substrate. The name of the enzyme itself is then given to the character. Primary homologies of “type II” were defined as the sharing of similar enzymatic functions by different pathways (although the enzymes recognized can have different specificities, “IIa”), sharing the same cofactors (“IIb”), or being of the same functional family (“IIc”). For the coding of homologies of type II, in a number of taxa, the coding involves reactions embedded within the Krebs cycle. The anabolism of alanine (sALA), valine (sVAL), serine (sSER1 and sSER2), glycine (sGLY1 and sGLY2), and cysteine (sCYS1 and sCYS2), gluconeogenesis, and the Calvin cycle all start with β -decarboxylation of oxaloacetate. However, there cannot be any β -decarboxylation without a preceding β -dehydrogenation of an alcohol to obtain the carbonyl precursor. For the pathways considered here, this is done within the Krebs cycle, from malate to oxaloacetate. So these pathways (taxons) have to incorporate β -dehydrogenation as a type II homology (coding “1” at character 27: alcohol NAD-dehydrogenation), as if they were starting from malate. For symmetrical reasons, the second portion of the Krebs cycle (from oxoglutarate to oxaloacetate, “KC2”) does include the presence of a β -decarboxylation from

Table 1
Names of Characters, with the Corresponding Number of
International Nomenclature

[1] Pyridoxalphosphate enzymes; IIb
[2] Deamination (NAD); IIa
[3] Transamination (PLP); IIa
[4] Amide deamination; IIa
[5] Aldehyde dehydrogenation (NAD); IIa
[6] α -Decarboxylation; IIa
[7] Deamination (PLP); IIa
[8] Carboxylation (biotine); IIa
[9] No Alcohol, no aldehyde NAD dehydrogenases; IIa
[10] β -Decarboxylation; IIa
[11] Decarboxylation; IIc
[12] Deamination; IIc
[13] Acid-ammonia ligases; IIa
[14] Phosphorylation; IIc
[15] Acyl-CoA synthases; IIa
[16] Thiamine pyrophosphate enzymes; IIb
[17] Transketolases; IIa
[18] β -Oxidation/reduction sequence repeat; IIc
[19] β -Oxidation sequence; IIa
[20] No niacin, no flavin Oxidases; IIa
[21] Tetrahydrofolate enzymes; IIb
[22] Transaldolases; IIa
[23] Acetyltransferases; IIa
[24] Succinyltransferases; IIa
[25] Malonyltransferases; IIa
[26] Ammonia lyases; IIa
[27] Alcohol dehydrogenation (NAD); IIa
[28] FAD dehydrogenation; IIa
[29] Hydratases; IIa
[30] Isomerases; IIa
[31] Hydrolases; IIa
[32] Phosphoribosyltransferase; IIa
[33] <i>Amino-acid dehydrogenase</i> : 1.4.1.5; I
[34] <i>Aspartate aminotransferase</i> : 2.6.1.1; I
[35] <i>Asparaginase</i> : 3.5.1.1; I
[36] <i>Glutamate dehydrogenase</i> : 1.4.1.3; I
[37] <i>Pyruvate dehydrogenase</i> : 1.2.4.1; I
[38] <i>Serine deaminase</i> : 4.3.1.19; I
[39] <i>Serine hydroxymethyltransferase</i> : 2.1.2.1; I
[40] <i>Cysteine desulphhydrase</i> : 4.4.1.1; I
[41] <i>Propionyl-CoA carboxylase</i> : 6.4.1.3; I
[42] <i>Branched-chain-amino-acid aminotransferase</i> : 2.6.1.42; I
[43] <i>Branched-chain ketoacid dehydrogenase</i> : 1.2.4.4; I
[44] <i>Acetyl-CoA acetyltransferase</i> : 2.3.1.9; I
[45] <i>Amino adipate-semialdehyde dehydrogenase</i> : 1.2.1.31; I
[46] <i>Amino adipate aminotransferase</i> : 2.6.1.39; I
[47] <i>Saccharopine dehydrogenase I</i> : 1.5.1.9; I
[48] <i>Alanine aminotransferase</i> : 2.6.1.2; I
[49] <i>Glutaminase</i> : 3.5.1.2; I
[50] <i>Glutamine-oxo-acid aminotransferase</i> : 2.6.1.15; I
[51] ω -Amidase: 3.5.1.3; I
[52] <i>Arginase</i> : 3.5.3.1; I
[53] <i>Ornithine-oxo-acid aminotransferase</i> : 2.6.1.13; I
[54] <i>Pyrroline-carboxylate reductase</i> : 1.5.1.2; I
[55] <i>Cysteine aminotransferase</i> : 2.6.1.3; I
[56] <i>Mercaptopyruvate sulphurtransferase</i> : 2.8.1.2; I
[57] <i>Cysteine dioxygenase</i> : 1.13.11.20; I
[58] <i>Aspartate decarboxylase</i> : 4.1.1.12; I
[59] <i>Homocysteine methyltransferase</i> : 2.1.1.10; I
[60] <i>Cystathionine synthase</i> : 4.2.1.22; I
[61] <i>Hydroxymethylbutyryl-CoA dehydrogenase</i> : 1.1.1.178; I
[62] <i>Methylcrotonoyl-CoA carboxylase</i> : 6.4.1.4; I
[63] <i>Methylglutaconyl-CoA hydratase</i> : 4.2.1.18; I
[64] <i>Hydroxymethylglutaryl-CoA lyase</i> : 4.1.3.4; I
[65] <i>Ketoacid CoA-transferase</i> : 2.8.3.5; I
[66] <i>Methylmalonate-semialdehyde dehydrogenase</i> : 1.2.1.27; I
[67] <i>Saccharopine dehydrogenase 2</i> : 1.5.1.7; I
[68] <i>Citrate synthase</i> : 2.3.3.1; I
[69] <i>Citrate dehydratase</i> : 4.2.1.3; I

Table 1
continued

[70] <i>Isocitrate dehydrogenase</i> : 1.1.1.41; I
[71] <i>α-Oxoglutarate dehydrogenase</i> : 1.2.4.2; I
[72] <i>Succinyl-CoA hydrolase</i> : 3.1.2.3; I
[73] <i>Succinate dehydrogenase</i> : 1.3.99.1 I
[74] <i>Fumarate hydratase</i> : 4.2.1.2; I
[75] <i>Malate dehydrogenase</i> : 1.1.1.37; I
[76] <i>Asparagine synthase (ADP-forming)</i> : 6.3.1.4; I
[77] <i>Asparagine synthase (glutamine-hydrolysing)</i> : 6.3.5.4; I
[78] <i>Glutamine synthetase</i> : 6.3.1.2; I
[79] <i>Amino-acid acetyltransferase</i> : 2.3.1.1; I
[80] <i>Acetylglutamate kinase</i> : 2.7.2.8; I
[81] <i>N-acetyl-glutamyl-P reductase</i> : 1.2.1.38; I
[82] <i>Acetylorithine aminotransferase</i> : 2.6.1.11; I
[83] <i>Acetylorithine deacetylase</i> : 3.5.1.16; I
[84] <i>Ornithine carbamoyltransferase</i> : 2.1.3.3; I
[85] <i>Argininosuccinate synthetase</i> : 6.3.4.5; I
[86] <i>Argininosuccinate lyase</i> : 4.3.2.1; I
[87] <i>Glutamate kinase</i> : 2.7.2.11; I
[88] <i>Glutamate semialdehyde dehydrogenase</i> : 1.5.1.12; I
[89] <i>Oxaloacetate decarboxylase</i> : 4.1.1.3; I
[90] <i>Phosphoenolpyruvate carboxykinase (GTP)</i> : 4.1.1.32; I
[91] <i>Enolase</i> : 4.2.1.11; I
[92] <i>Phosphoglycerate phosphomutase</i> : 5.4.2.1; I
[93] <i>Glycerate kinase</i> : 2.7.1.31; I
[94] <i>Hydroxypyruvate reductase</i> : 1.1.1.81; I
[95] <i>Serine-pyruvate aminotransferase</i> : 2.6.1.51; I
[96] <i>Phosphoglycerate dehydrogenase</i> : 1.1.1.95; I
[97] <i>Phosphoserine aminotransferase</i> : 2.6.1.52; I
[98] <i>Phosphoserine phosphatase</i> : 3.1.3.3; I
[99] <i>Acyl-CoA dehydrogenase</i> : 1.3.99.3; I
[100] <i>Serine acetyltransferase</i> : 2.3.1.30; I
[101] <i>Cysteine synthase</i> : 2.5.1.47; I
[102] <i>Aspartate kinase</i> : 2.7.2.4; I
[103] <i>Aspartate-semialdehyde dehydrogenase</i> : 1.2.1.11; I
[104] <i>Homoserine dehydrogenase</i> : 1.1.1.3; I
[105] <i>Homoserine succinyltransferase</i> : 3.1.46; I
[106] <i>Cystathionine synthase</i> : 2.5.1.48; I
[107] <i>Cystathionine lyase</i> : 4.4.1.8; I
[108] <i>Homoserine kinase</i> : 2.7.1.39; I
[109] <i>Threonine synthase</i> : 4.2.3.1; I
[110] <i>Acetolactate synthase</i> : 2.2.1.6; I
[111] <i>Keto-acid reductoisomerase</i> : 1.1.1.86; I
[112] <i>Dihydroxyacid dehydratase</i> : 4.2.1.9; I
[113] <i>Isopropylmalate synthase</i> : 2.3.3.13; I
[114] <i>Isopropylmalate isomerase</i> : 4.2.1.33; I
[115] <i>Isopropylmalate dehydrogenase</i> : 1.1.1.85; I
[116] <i>Dihydrodipicolinate synthase</i> : 4.2.1.52; I
[117] <i>Dihydrodipicolinate reductase</i> : 1.3.1.26; I
[118] <i>Piperideindicarboxylate succinyltransferase</i> : 2.3.1.117; I
[119] <i>Succinyl-diaminopimelate aminotransferase</i> : 2.6.1.17; I
[120] <i>Succinyl-diaminopimelate desuccinylase</i> : 3.5.1.18; I
[121] <i>Diaminopimelate epimerase</i> : 5.1.1.7; I
[122] <i>Diaminopimelate decarboxylase</i> : 4.1.1.20; I
[123] <i>Homocitrate synthase</i> : 2.3.3.14; I
[124] <i>Homoaconitate hydratase</i> : 4.2.1.36; I
[125] <i>Homoisocitrate dehydrogenase</i> : 1.1.1.155; I
[126] <i>Phosphoglycerate kinase</i> : 2.7.2.3; I
[127] <i>Triosephosphate dehydrogenase</i> : 1.2.1.12; I
[128] <i>Triose-phosphate isomerase</i> : 5.3.1.1; I
[129] <i>Fructose-bisphosphate aldolase</i> : 4.1.2.13; I
[130] <i>Fructose-bisphosphatase</i> : 3.1.3.11; I
[131] <i>Glucose-6-phosphate isomerase</i> : 5.3.1.9; I
[132] <i>6-Phosphofructokinase</i> : 2.7.1.11; I
[133] <i>Pyruvate kinase</i> : 2.7.1.40; I
[134] <i>Glycolaldehydetransferase</i> : 2.2.1.1; I
[135] <i>Sedoheptulose diphosphatase</i> : 3.1.3.37; I
[136] <i>Ribulose-phosphate 3-epimerase</i> : 5.1.3.1; I
[137] <i>Ribulose-5-phosphate isomerase</i> : 5.3.1.6; I
[138] <i>Phosphoribulokinase</i> : 2.7.1.19; I
[139] <i>Ribulose diphosphate carboxylase</i> : 4.1.1.39; I
[140] <i>Glucose-6-phosphate 1-dehydrogenase</i> : 1.1.1.49; I

Table 1
continued

[141] <i>Phosphogluconolactonase</i> : 3.1.1.31; I
[142] <i>Phosphogluconic acid dehydrogenase</i> : 1.1.1.44; I
[143] <i>Dihydroxyacetone transferase</i> : 2.2.1.2; I
[144] <i>Dihydrolipoamide dehydrogenase</i> : 1.8.1.4; I
[145] <i>Acetate thiokinase</i> : 6.2.1.13; I
[146] <i>Acyl-CoA oxidase</i> : 1.3.3.6; I
[147] <i>Enoyl-CoA hydratase</i> : 4.2.1.17; I
[148] β - <i>Hydroxyacyl dehydrogenase</i> : 1.1.1.35; I
[149] <i>Acetyl-CoA C-acyltransferase</i> : 2.3.1.16; I
[150] <i>Trans-2-enoyl-CoA reductase (NAD)</i> : 1.3.1.44; I
[151] <i>Acyl-ACP-hydrolase</i> : 3.1.2.14; I
[152] <i>Acetyl-CoA carboxylase</i> : 6.4.1.2; I
[153] <i>Acyl-carrier-protein</i>] <i>S-malonyltransferase</i> : 2.3.1.39; I
[154] <i>Acyl-carrier-protein</i>] <i>S-acetyltransferase</i> : 2.3.1.38; I
[155] β - <i>Ketoacyl synthetase</i> : 2.3.1.41; I
[156] β - <i>Ketoacyl reductase</i> : 1.1.1.100; I
[157] <i>Enoyl acyl carrier protein hydratase</i> : 4.2.1.58; I
[158] <i>Acyl-ACP dehydrogenase</i> : 1.3.1.10; I
[159] β - <i>Hydroxydecanoate dehydratase</i> : 4.2.1.60; I
[160] <i>Glutaryl-CoA dehydrogenase</i> : 1.3.99.7; I
[161] <i>Phenylalanine 4-monoxygenase</i> : 1.14.16.1; I
[162] <i>Hydroxyphenylpyruvate dioxygenase</i> : 1.13.11.27; I
[163] <i>Dihydroxyphenylacetate 2,3-dioxygenase</i> : 1.13.11.15; I
[164] <i>Maleylacetoacetate isomerase</i> : 5.2.1.2; I
[165] <i>Fumarylacetoacetate</i> : 3.7.1.2; I
[166] <i>Tryptophan 2,3-dioxygenase</i> : 1.13.11.11; I
[167] <i>Kynurenine formamidase</i> : 3.5.1.9; I
[168] <i>Kynurenine hydroxylase</i> : 1.14.13.9; I
[169] <i>Kynureninase</i> : 3.7.1.3; I
[170] <i>Hydroxyanthranilate oxygenase</i> : 1.13.11.6; I
[171] <i>Picolinic acid carboxylase</i> : 4.1.1.45; I
[172] <i>Aminomuconate-semialdehyde dehydrogenase</i> : 1.2.1.32; I
[173] <i>Histidine ammonia-lyase</i> : 4.3.1.3; I
[174] <i>Urocanate hydratase</i> : 4.2.1.49; I
[175] <i>Imidazolone propionic acid hydrolase</i> : 3.5.2.7; I
[176] <i>Glutamate formyltransferase</i> : 2.1.2.5; I
[177] <i>Aromatic amino acid aminotransferase</i> : 2.6.1.57; I
[178] <i>Deoxy-7-phosphoheptulonate synthase</i> : 2.5.1.54; I
[179] <i>Dehydroquininate synthase</i> : 4.2.3.4; I
[180] <i>Dehydroquininate dehydratase</i> : 4.2.1.10; I
[181] <i>Shikimate 5-dehydrogenase</i> : 1.1.1.25; I
[182] <i>Shikimate kinase</i> : 2.7.1.71; I
[183] <i>Riboflavin synthase</i> : 2.5.1.9; I
[184] <i>Anthranilate synthase</i> : 4.1.3.27; I
[185] <i>Anthranilate phosphoribosyltransferase</i> : 2.4.2.18; I
[186] <i>Phosphoribosylanthranilate isomerase</i> : 5.3.1.24; I
[187] <i>Indoleglycerol phosphate synthetase</i> : 4.1.1.48; I
[188] <i>Tryptophan synthase</i> : 4.2.1.20; I
[189] <i>Prephenate dehydrogenase</i> : 1.3.1.12; I
[190] <i>Chorismate mutase</i> : 5.4.99.5; I
[191] <i>Prephenate dehydratase</i> : 4.2.1.51; I
[192] <i>Methylmalonyl-CoA epimerase</i> : 5.1.99.1; I
[193] <i>Methylmalonyl-CoA mutase</i> : 5.4.99.2; I
[194] <i>Carbamate kinase</i> : 2.7.2.2; I
[195] <i>ATP phosphoribosyltransferase</i> : 2.4.2.17; I
[196] <i>Phosphoribosyl-ATP diphosphatase</i> : 3.6.1.31; I
[197] <i>Phosphoribosyl-AMP cyclohydrolase</i> : 3.5.4.19; I
[198] <i>Phosphoribosylformiminoaminophosphoribosylimidazole-carboxamide isomerase</i> : 5.3.1.16; I
[199] <i>Imidazoleglycerol-phosphate dehydratase</i> : 4.2.1.19; I
[200] <i>Histidinol-phosphate transaminase</i> : 2.6.1.9; I
[201] <i>Histidinol-phosphatase</i> : 3.1.3.15; I
[202] <i>Histidinol dehydrogenase</i> : 1.1.1.23; I

NOTE.—Homology types defined in the text. For characters 1 to 32, homologies are of type II. For characters 33 to 204, homologies are of type I (Enzyme Nomenclature 1973).

oxaloacetate to pyruvate or phosphoenolpyruvate in the coding of character 10.

A new kind of homology of type II was used in the present work. Pathways can be similar in the recurrence of a set of reactions made by the same enzymes from different substrates. This homology is used in the metabolism of fatty acids (character 18) and called "IId." All characters are named in table 1, with the number from enzyme nomenclature and the type of homology involved. In figures 1 and 2, characters involving homologies of type II are numbered from 1 to 32. It must be stressed that, if homologies of type I are named strictly following the international enzymatic nomenclature, homologies of type II are not. For example, hydratase (character 29) is used here in a wider meaning than in the international nomenclature. A reaction involving the pyridoxal-phosphate is coded as involving a hydratase because a molecule of water is implied, whereas it is not the case for international nomenclature. Conversely, our delineation of homologies of type II is more precise for dehydrogenases (or reductases) involving NAD. Alcohol-NAD-dehydrogenases, aldehyde-NAD-dehydrogenases, deaminase-NAD-dehydrogenases, and FAD-dehydrogenases have been separated into different characters.

A number of characters referring to homologies of type II contain question marks. These are included when the taxon-pathway does not exhibit the appropriate type of substrate for the enzymatic function, the coenzyme, or the functional family at hand. A question mark is also used when the enzymatic function to be coded is meaningless for the taxon. For instance, the coding of functions of degradation into anabolic pathways is meaningless: there is no way for an α -decarboxylation into a biosynthetic pathway.

Phylogenetic Reconstruction Tree Search

The matrix contains 75 taxa and 202 characters (table 2). Characters are treated as unordered and unweighted. Heuristic searches were conducted with NONA (Goloboff 1998) as implemented into WINCLADA (Nixon 1999a), using TBR branch swapping. For a better exploration of trees, the Parsimony Ratchet (Hopper Islands [Nixon 1999b]) was used. The proportion of data to be reweighted was set between 25% and 50% and the number of iterations progressively increased from 25,000 to 150,000 (option amb- poly=). This increase in iterations was used to check that the number of supplementary MP trees gained each time was decreasing or null. Each time the number of trees was recorded after having collapsed, all unsupported nodes in all trees ("hard collapse").

Rooting

The tree was rooted using an all-zero hypothetical ancestor (HYPANC). This is justified by the fact that, in the coding of character states, zero was given to the absence of enzymes or to the absence of performance of particular functions (even in presence of a putative suitable substrate) or to absence of use of a cofactor. Such a rooting option will automatically put the simplest pathways closer to the root. However, this does not make any assumption of the nature of the corresponding enzymatic reactions.

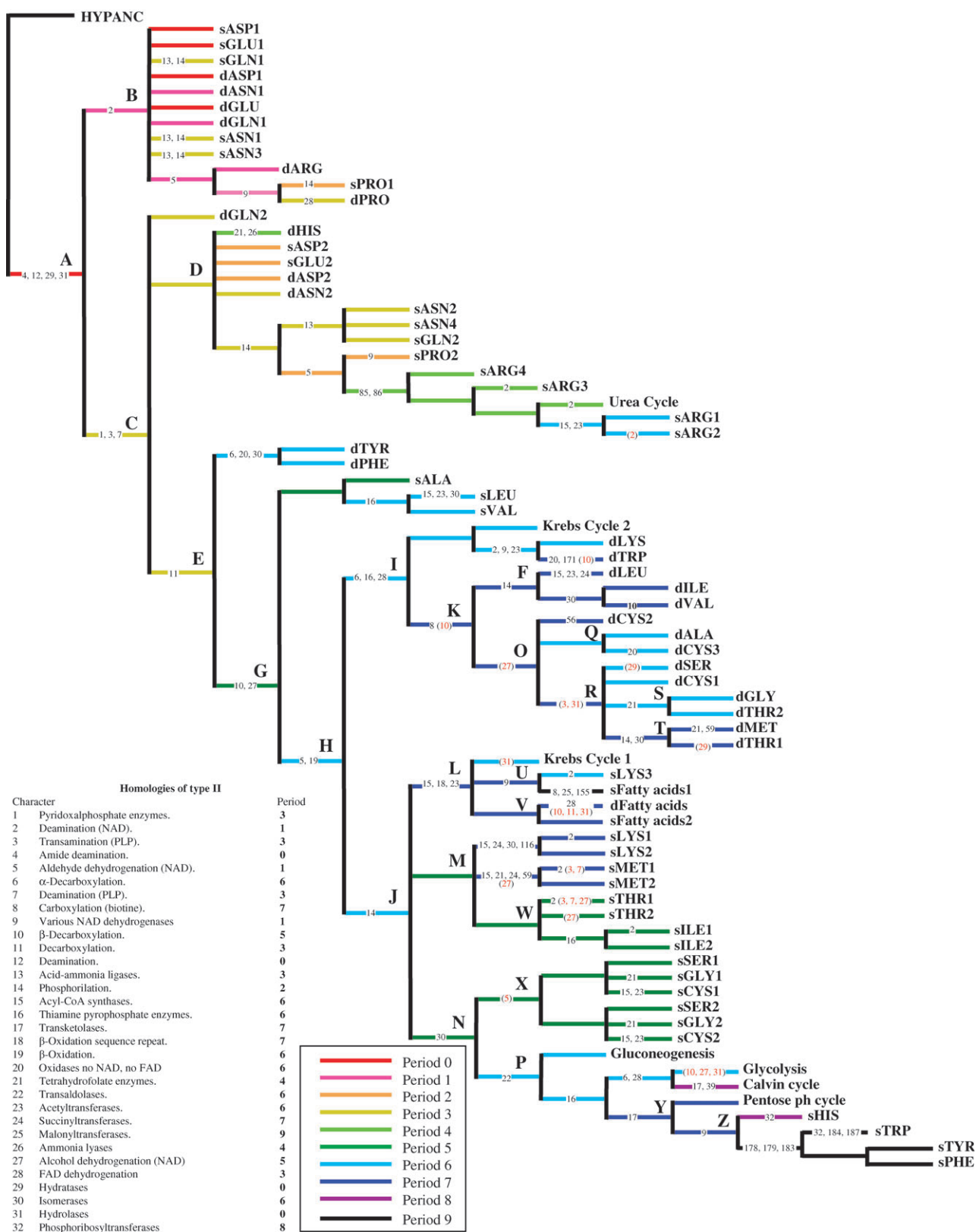


FIG. 1.—Strict consensus of 20 equiparsimonious trees obtained through the Parsimony Ratchet of WINCLADA calculated on the matrix. Each tree is of 347 steps (CI = 0.58, RI = 0.79). Character reversions in red.

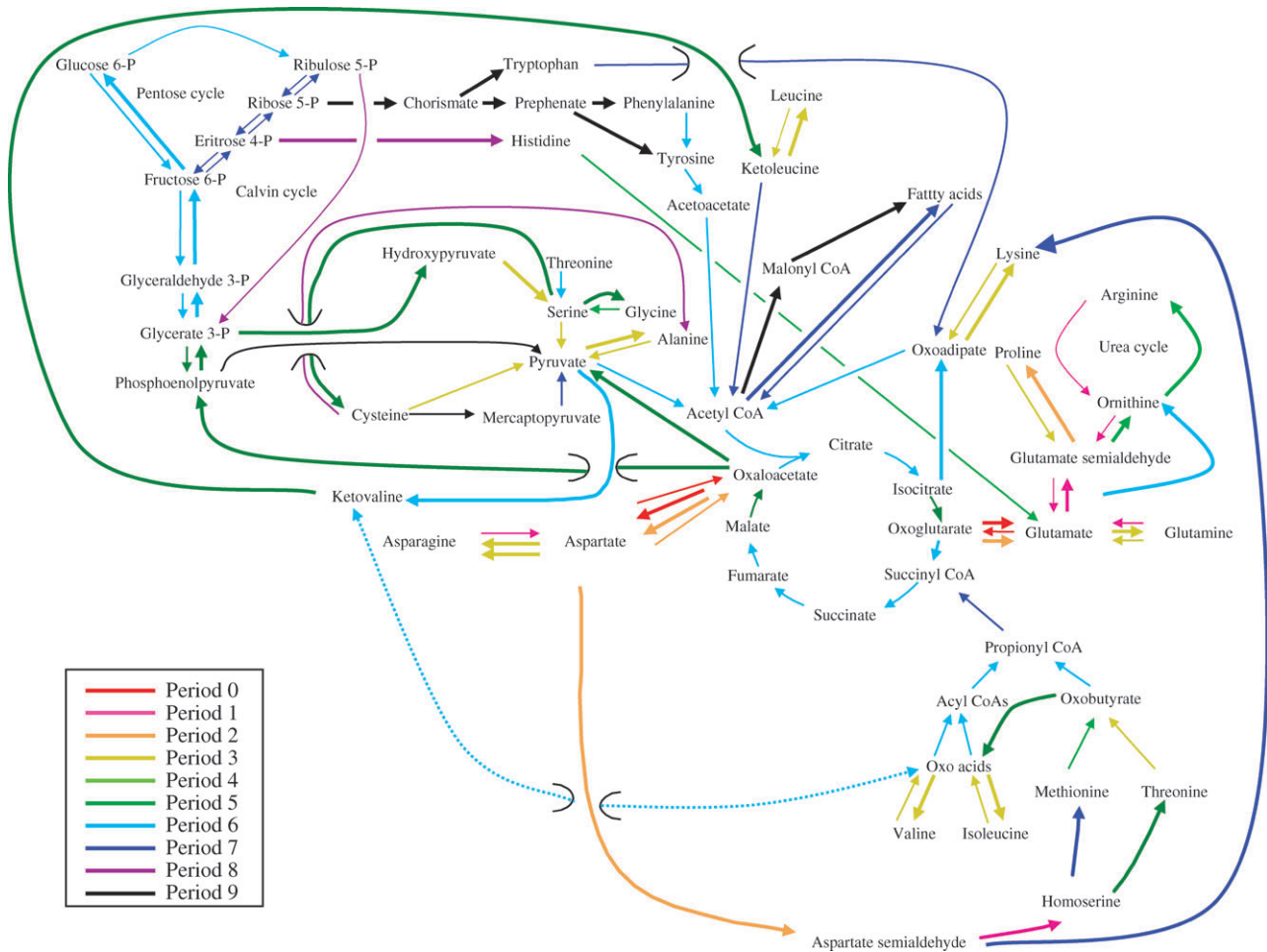


FIG. 2.—General view of metabolic pathways and their connecting points to the Krebs cycle. Colors are successive time spans (or “periods”) as inferred from the tree. Reactions in black (period 9) do not mean that all these reactions appeared at the same time. These are just reactions later than the purple period (period 8).

in homologies of type II, the downstream branches are of the same period. When a new homology of type II occurs on a branch, it defines the next period. When the homology of type II occurs in the next branch but is already present in earlier periods (homoplasy), it does not define a new period. For example, the node H is light blue (fig. 1, period 6), and the downstream node J exhibits a change in homology of type II (homology II n°14) while it remains in light blue. There is no new period here because the phosphorylation is already available in a previous period (period 2, within the clade B: sPRO1).

No homoplasy. When a character exhibits homoplasy, classically one has a choice in localizing changes in the tree (convergences, reversions, or more complex optimizations). Only transformations with unambiguous localization (i.e., when place on the tree does not depend on any choice) will be taken into account in the use of the second criterion.

Classically, homologies of type II involve several enzymes with different specificities. However, for 14 cases, single enzymes with the same high specificity

innovate a new enzymatic mechanism. They are recorded as homologies of type I, but because they also correspond to enzymatic changes used to define homologies of type II, they are taken into account for defining periods. In figure 1, they are shown with asterisks.

Polytomies

It must be stressed that both the tree topology and the new nonhomoplastic changes in homologies of type II are involved in the definition of periods. How this can be used when parts of the tree are unresolved? A polytomy is just an absence of resolution. A polytomy normally contains branches of different periods, simply because there are different numbers of new homologies of type II occurring on them. For example, within the clade B (fig. 1), a polytomy involves branches of the same period as the origin of the clade B (period 1, pink) and branches of later periods 2 (orange) and 3 (yellow). Some of these branches are of period 3 because they exhibit two additional homologies of type II (13 and 14). Homologies of type II are, therefore, cumulative in defining periods.

The Problem of Question Marks

Question marks bring complications in the definition of time spans. As a consequence of using the parsimony criterion, some homologies of type II are optimized onto a node deeper than the ones where the derived state is really observed, because of question marks in the rest of terminals of the clade. For example, character 8 (carboxylations using biotin) is optimized to occur on the node K. Among the 12 pathways included within K, only five (dLEU, dILE, dVAL, dMET, and dTHR1) really exhibit this function; all others are coded “?” for that function. Thus, the early position of the gain of character 8 is because of the parsimony criterion used to manage question marks in character optimization, whereas real observation of this function would optimize two gains of character 8 onto nodes T and F only. This causes no problem, except for the ordering the time spans. K is of period 7 because of the gain of this homology of type II. However, seven downstream pathways (dALA, dCYS2, dCYS3, dSER, dCYS1, dGLY, and dTHR2) do not really have to wait for that period to be achieved, because they do not need to use this function, for which they are coded “?”. Therefore, assigning their complete development to period 7 is an artifact of optimization affecting time-span assignment. Because these pathways are actually possible as early as the previous period (period 6), they have been assigned to that period. This is the reason why there can be an apparent contradiction in observing a node or a branch to which is assigned a period earlier than an upstream (more inclusive) branch (dCYS1 is possible in period 6 while node K is in period 7). Such situations are met four times in the analysis and are always caused by question marks.

Results

From our data matrix, the heuristic search using the Parsimony Ratchet stabilized the number of equiparsimonious trees to 20 (347 steps, CI = 0.58, RI = 0.79). The strict consensus is shown figure 1. The first enzymes to occur (node A) are hydrolases, hydratases, and those involving deaminations of amino acids. Therefore, degradation of amino acids must have preceded all other kind of reactions.

Cordón (1990) defined four groups of amino acids, according to common reactions in their metabolism and the point of entry into the Krebs cycle: group I enters by oxaloacetate, group II by ketoglutarate, group III by pyruvate and acetyl-CoA, and group IV by succinyl-CoA. In our tree (fig.1), amino acids of Cordón's groups I and II develop first. At period 6 (light blue), it is possible to synthesize and degrade all these amino acids (Asp, Asn, Arg, Pro, Gln, and Glu). Based on the full development of the synthesis of arginine (sARG3 and sARG4), the urea cycle is fully developed by the period 4 (light green). The complete metabolism of Cordón's groups III (Ser, Gly, Cys, and Ala) and IV (Thr, Val, Ile, and Met) starts at period 5 (deep green) and mostly develops during periods 5 (for syntheses) and 6 to 7 (for degradations). Aromatic amino acids are degraded as early as period 4 for histidine and period 6 or later for others, whereas their synthesis is

so sophisticated that it is only possible very late in periods 8 and 9.

The two parts of the Krebs cycle develop within period 6. Period 5 (deep green) is the period of termination of almost all aliphatic amino acid syntheses (except for lysine and methionine, whose syntheses are achieved in period 7). Period 6 (light blue) is very rich; that is, it contains many events that cannot be ordered: the closing of the Krebs cycle, the termination of the degradation of all aliphatic amino acids of groups III (except dCYS2), and the full development of glycolysis and gluconeogenesis. Later in period 7 (deep blue) almost all aliphatic amino acid metabolism (except dCYS2), degradation of fatty acids, and one of the two syntheses of fatty acids (the “intramitochondrial pathway,” here numbered “2”) are all possible. It is noticeable that the relative order of glycolysis and the closure of the Krebs cycle cannot be clarified; both of them appear to be the consequence of amino acid metabolic pathways. Period 8 (purple) closes the complete development of all aliphatic amino acid metabolisms and develops the Calvin cycle. Period 10 is arbitrarily defined as containing all events later than period 9. This includes syntheses of aromatic amino acids and fatty acids along the “extramitochondrial pathway,” here numbered “1”.

Some branches are in a period earlier than one of their upstream branches because of question mark optimization (see above). This is the case for optimizations of characters 4, 8, 18, and 19. Question marks in character 8 (carboxylations using biotin) imply that terminal branches dALA, dCYS3, dSER, dCYS1, dGLY, dTHR2, and branches Q and S are in period 6 while their upstream branches K, O, and R are in period 7. Question marks of character 18 (β -oxidation sequence repeat) in KC1 and sLYS3 place terminals KC1 and sLYS3 in period 6 while their upstream nodes L and U are in period 7. Question marks in character 19 (β -oxidations) for pathways of clade H (except KC2, dLYS, dILE, dVAL, and synthesis and degradation of fatty acids) place the gain of β -oxidation onto the branch H. That branch is assigned to period 6 because of the new gain of β -oxidation (character 5, already available from period 2, is not new). All pathways that exhibit a question mark for character 19 are assigned one period earlier. Branch M, all members of clade W, branch N, and all members of clade X are, therefore, assigned to period 5, which is earlier than period 6 of H. The same reasoning is followed for character 4, which is coded “1” in all pathways of asparagine and glutamine and coded “?” in all others. All others are assigned one period earlier. Question marks explain why one can find new homology of type II in a branch without changing the period. The assignment of the period involves one period forward because of the new type II homology and one period earlier because of a question mark for another character. For example, within the clade B, the node on which character “5” occurs remains in period 1 (pink) because one step forward for the new homology “5” and one step earlier because of question marks on character “4” for dARG, sPRO1, dPRO. Another example is the node E, which remains in period 3 despite the new homology “11”.

Discussion

The rise of different enzymes and enzymatic functions are plotted onto the tree (fig. 1), and therefore each event is associated to a given period. Then, the different parts of a given pathway can belong to different periods according to the enzymes involved. This allows a first view on the temporal development of each pathway, as shown in figure 2. For instance, degradation of valine (dVAL) starts in period 2 and ends in period 7.

Amino acid metabolism is basal in the phylogenetic tree, and all other metabolisms develop from that background. The priority of amino acid degradation over syntheses depends on the amino acid. Full degradations are earlier than full syntheses for amino acids of Cordón's groups I (Asp and Asn) and II (Glu, Gln, Arg, and Pro, excepting proline, for which full synthesis is possible in period 2 while degradation appears in period 3) and for aromatic amino acids (His, Tyr, and Phe), contradicting Cordón's views on priority of aromatic amino acid syntheses over their degradations. For amino acids of Cordón's groups III (Ala, Cys, Gly, and Ser) and IV (Leu, Val, Ile, Lys, Thr, and Met) full syntheses are possible before full degradations (except for Trp), contradicting Cordón's views. In groups I and II, degradation and synthesis tend to be reverse processes, whereas in groups III and IV, they are different pathways. This observation is also true for aromatic amino acids, as the syntheses are different and very late in comparison to degradations. Our results do not fully corroborate predictions of Horowitz (1945) and Cordón (1990) according to which amino acid catabolic pathways develop forwards and which anabolic pathways develop backwards. Indeed, dPRO, dHIS, dALA, dTRP, and dMET develop backwards, whereas none of the amino acid anabolisms are found to develop backwards, except sLYS. Furthermore, in some cases, opportunistic late connections of pathways lead to complex developments neither forwards nor backwards as in sLEU, sVAL, dILE, sGLY, sTHR, sMET, dLEU, and dGLY. Horowitz (1945) and Cordón's (1990) predictions were not based on a data matrix, but (1) on the hypothesis that amino acids of abiotic origin might have been one of the very first sources of energy and structural components available to primitive forms of life (which has been highly corroborated since), (2) on a theoretical model of selective pressure acting on protocells during the course of biochemical evolution, and (3) on consistent comparative reasoning from structures of metabolic pathways and their compounds (see Cunchillos and Lecointre [2000]). No computerized parsimony was used, and it is not surprising that a formalized reconstruction algorithm using parsimony can deal with complex issues of homoplastic late opportunistic pathway connections better than a human mind can.

Metabolism of fatty acids and saccharides develop after the full development of metabolism of amino acids of groups I and II, and they are associated with the anabolism of amino acids of groups III and IV. Sugar metabolisms are within the clade N, with anabolic pathways of the amino acids of group III, as predicted by Cordón (1990). Syntheses of aromatic amino acids are branched within sugar metabolism. From the present data, it is neither

possible to draw conclusions about forward or backwards development of fatty acid and sugar metabolisms nor possible to draw conclusions about priority of either their anabolism or their catabolism. Note that the extramitochondrial synthesis of fatty acids precedes (period 7) the intramitochondrial synthesis (period 9).

Period 6 is very rich, a period during which events are difficult to order and the two portions of the Krebs cycle take place after the setting of metabolisms of amino acids of groups I and II. Interestingly, one portion of the Krebs cycle has a catabolic origin and the other an anabolic origin: KC2 is associated with catabolism of amino acids of groups III and IV, and KC1 is associated with anabolism of the same groups. Despite the late branching of glycolysis in the tree, it is difficult to order the rise of full glycolysis and the rise of full Krebs cycle: they are all embedded within period 6. So the views of Meléndez-Hevia et al. (1996), that glycolysis must have preceded the Krebs cycle because it provided energy, are neither confirmed nor contradicted. However, relative position of a catabolic pathway must not only be assessed in terms of energy but also in the light of the compounds provided. Availability of glucose in early abiotic environments is much more speculative than availability of amino acids (Cunchillos and Lecointre 2000, 2002). For that reason, because the Krebs cycle is derived from amino acid metabolisms of groups III and IV, the hypothesis that the Krebs cycle arose earlier than glycolysis is much more likely.

From the present data, the ordering of gluconeogenesis compared with glycolysis is not possible: both are in period 6. Pentose-phosphate and Calvin cycles arose later (periods 7 and 8, respectively) than glycolysis and gluconeogenesis, as predicted by Cordón (1990). The urea cycle arises as early as the synthesis of arginine in period 4.

The present inferences partially corroborate Cordón's scenario, which, after all, is not surprising, because they are only based on the comparative anatomy of metabolic pathways, whereas Cordón incorporated much biological input. It is very interesting that his scenario was embedded into a consistent and well-documented DNA-free conception of natural selection of enzymatic specificity, leading to the evolution of metabolism through natural selection of proteins without any DNA information storage. This is one of the reasons why the present work has nothing to do with theories of early information storage: to test Cordón's scenario, such storage was not necessary. Concerning information storage in nucleic acids, it should be stressed that the compounds involved in the pathways studied here are basic elements for construction of nucleic acids of metabolic origin. Therefore, those nucleic acids arose later than the pathways. Biological purines and pyrimidines are the products of several already complex molecules; they are never synthesized as such in cells, because their parts are already linked to chemically different compounds such as riboses. From the evolutionary point of view, they are later chemical compounds than all the metabolic pathways compared here. The question of knowing whether there might have been DNAs or RNAs of nonmetabolic origin encoding "information" for enzymatic activities studied cannot be answered here. Cordón's scenario for the rise of protein complexity and enzymatic specificity by natural

selection without invoking DNA/RNA “information” storage is more parsimonious than explanations that require such a storage and is well worth further consideration for positive input in the increasing criticisms of the current uses of notions of genetic control and information (Kupiec and Sonigo 2000; Segal 2003).

Acknowledgments

We thank Patrick Tort for help to C.C., and John Macdonald and Craig Marshall for help in improving the manuscript. The Muséum National d’Histoire Naturelle and the company Saint-Gobain are acknowledged for support.

Literature Cited

- Cordón, F. 1990. *Tratado evolucionista de biología*. Aguilar, Madrid.
- Cunchillos, C. and G. Lecointre. 2000. L’histoire du catabolisme des acides aminés aliphatiques inférée par l’analyse cladistique de deux nouveaux types de caractères : l’enzyme et la réaction enzymatique. Pp. 87–106 *in* V. Barriel and T. Bourgoïn, eds. *Caractères*. Biosystema 18, Société Française de Systématique, Paris.
- . 2002. Early steps of metabolism evolution inferred by cladistic analysis of the structure of amino acid catabolic pathways. *C. R. Biologies* **325**:119–129.
- . 2003. Evolution of amino acid metabolism inferred through cladistic analysis. *J. Biol. Chem.* **278**:47960–47970.
- de Pinna, M. C. C. 1991. Concepts and tests of homology in the cladistic paradigm. *Cladistics* **7**:367–394.
- Enzyme Nomenclature. 1973. Recommendations (1972) of the International Union of pure and applied chemistry and the International Union of biochemistry. Elsevier Scientific Publishing, Amsterdam.
- Goloboff, P. 1998. Nona: computer program and software. Published by the author, Tucuman, Argentina.
- Horowitz, N. H. 1945. On the evolution of biochemical syntheses. *Proc. Natl. Acad. Sci. USA* **31**:153–157.
- Kupiec, J. J. and P. Sonigo. 2000. *Ni Dieu ni gène*. Seuil, Paris.
- Meléndez-Hevia, E., T. G. Waddell, and M. Cascante. 1996. The puzzle of the Krebs citric acid cycle: assembling the pieces of chemically feasible reactions, and opportunism in the design of metabolic pathways during evolution. *J. Mol. Evol.* **43**:293–303.
- Nixon, K. C. 1999*a*. Winclada (BETA). Version 0.9.9. Published by the author, Ithaca, New York.
- . 1999*b*. The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics* **15**:407–414.
- Segal, J. 2003. *Le zéro et le un, histoire de la notion scientifique d’information au 20^{ème} siècle*. Syllepse, Paris.

Geoffrey McFadden, Associate Editor

Accepted August 12, 2004