# Correlated evolutionary scenarios of metabolic functions

Murray Patterson[1], Thomas Bernard[1] and Daniel Kahn[1,2]

[1]LBBE, Université de Lyon 1 and [2]Département MIA, INRA
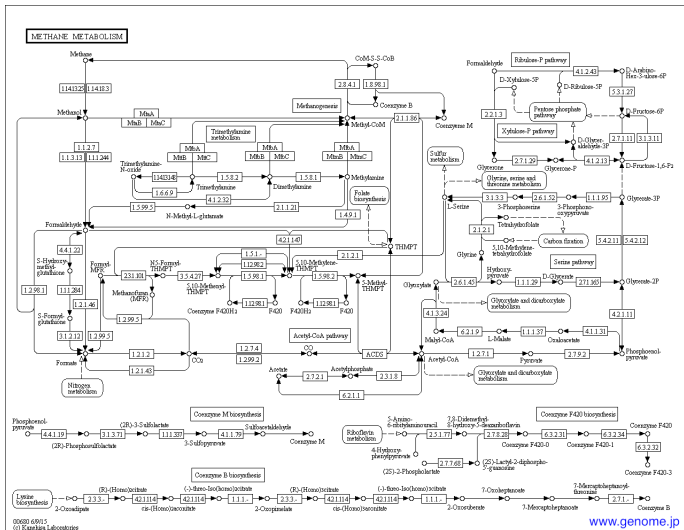
MCEB, L'île de Porquerolles, Jun 2015

# Metabolic functions

Determined by the enzymes (ECs) present – enzymes catalyze reactions

- many enzymes can be present at the same time, creating a complex metabolic network

# Metabolic functions

Determined by the enzymes (ECs) present – enzymes catalyze reactions

- many enzymes can be present at the same time, creating a complex metabolic network

# Metabolic functions, cont'd

### Inferring metabolic functions

In principle, this can be inferred from the <span style="color:red">nucleic acid sequence</span> content, e.g., as encoded by a gene or a genome

---

[1]PRIAM: PRofils pour l'Identification Automatique du Métabolisme (priam.prabi.fr)
[2]pbil.univ-lyon1.fr/databases/hogenom

# Metabolic functions, cont'd

## Inferring metabolic functions

In principle, this can be inferred from the nucleic acid sequence content, e.g., as encoded by a gene or a genome

- ▶ we do this systematically using PRIAM[1] – a set of profiles (PSSMs) for protein modules covering the Swiss-Prot database

---

[1]PRIAM: PRofils pour l'Identification Automatique du Métabolisme (priam.prabi.fr)
[2]pbil.univ-lyon1.fr/databases/hogenom

# Metabolic functions, cont'd

### Inferring metabolic functions

In principle, this can be inferred from the nucleic acid sequence content, e.g., as encoded by a gene or a genome

- we do this systematically using PRIAM[1] – a set of profiles (PSSMs) for protein modules covering the Swiss-Prot database
- idea : an rps-blast of the profiles against a protein sequence delivers a set of enzymes (ECs) with associated (presence) probabilities

---

[1]PRIAM: PRofils pour l'Identification Automatique du Métabolisme (priam.prabi.fr)
[2]pbil.univ-lyon1.fr/databases/hogenom

# Metabolic functions, cont'd

### Inferring metabolic functions

In principle, this can be inferred from the nucleic acid sequence content, e.g., as encoded by a gene or a genome

- we do this systematically using PRIAM[1] – a set of profiles (PSSMs) for protein modules covering the Swiss-Prot database
- idea : an rps-blast of the profiles against a protein sequence delivers a set of enzymes (ECs) with associated (presence) probabilities

### The evolution of metabolic functions

We then want to explore evolutionary scenarios of these functions in order to understand the dependencies between them

- we do a study on HOGENOM 6[2] – a database of homologous gene families

---

[1]PRIAM: PRofils pour l'Identification Automatique du Métabolisme (priam.prabi.fr)
[2]pbil.univ-lyon1.fr/databases/hogenom

# Hogenom Families

There are 296 917 families in HOGENOM 6

---

# Hogenom Families

There are 296 917 families in HOGENOM 6

- of which 10 699 (3.60% of the 296 917) families[3] have an EC assignment ($P(\text{EC}) > 0.5$) , i.e., at least one of its sequences has an EC assigned

---

[3]note that we restrict to families with a tree on at least 3 sequences

# Hogenom Families

There are 296 917 families in HOGENOM 6
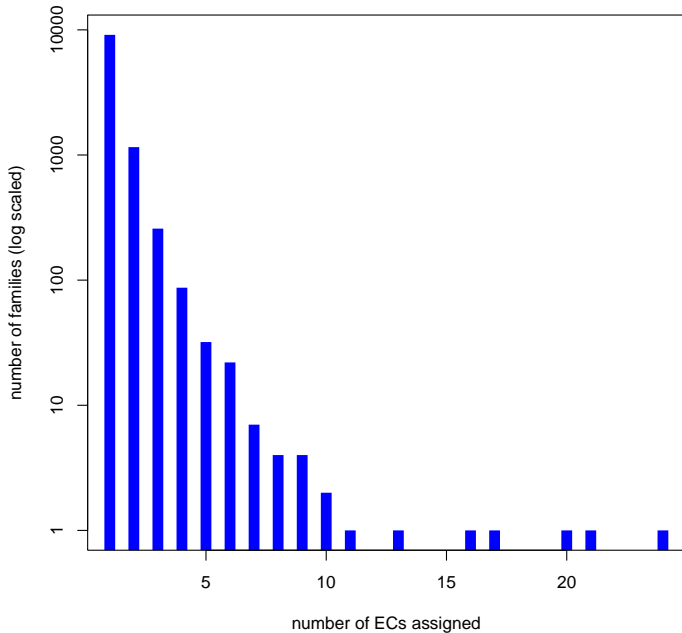
- of which 10 699 (3.60% of the 296 917) families[3] have an EC assignment ($P(EC) > 0.5$) , i.e., at least one of its sequences has an EC assigned

## First question :

How is EC assignment distributed amongst these 10 699 families?

---

[3]note that we restrict to families with a tree on at least 3 sequences

**Number of ECs assigned to a family (histogram)**

number of families (log scaled)

number of ECs assigned

# Families with one and two EC assignements

Of the 10 699 families with at least one EC assignment, no family has more than 24 ECs assigned to it, while

# Families with one and two EC assignements

Of the 10 699 families with at least one EC assignment, no family has more than 24 ECs assigned to it, while

- 9120 : 85.2% of the families have a single EC assignment, and

# Families with one and two EC assignements

Of the 10 699 families with at least one EC assignment, no family has more than 24 ECs assigned to it, while

- 9120 : 85.2% of the families have a single EC assignment, and
- 1156 : 10.8% of the families have exactly two EC assignments

# Families with one and two EC assignements

Of the 10 699 families with at least one EC assignment, no family has more than 24 ECs assigned to it, while

- 9120 : 85.2% of the families have a single EC assignment, and
- 1156 : 10.8% of the families have exactly two EC assignments

that is, 10 278 : 96.0% of the (10 699) families with any EC assignment have either one or two EC assignments

# Families with one and two EC assignements

Of the 10 699 families with at least one EC assignment, no family has more than 24 ECs assigned to it, while
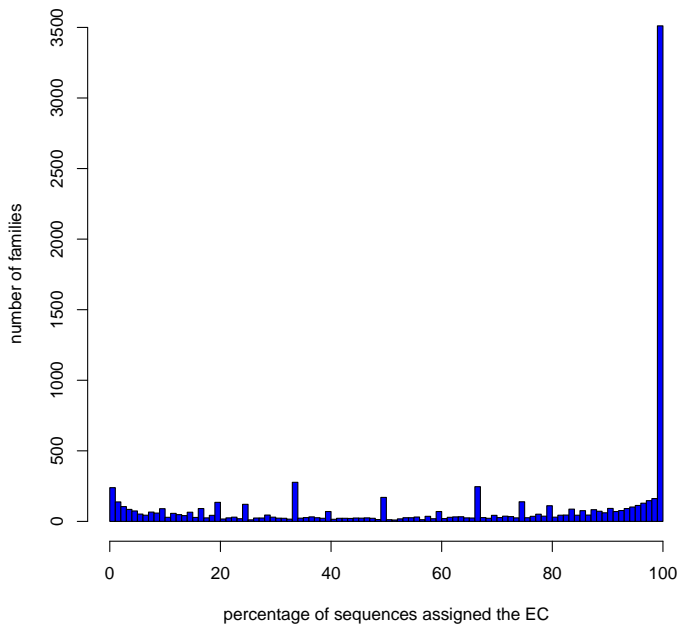
- 9120 : 85.2% of the families have a single EC assignment, and
- 1156 : 10.8% of the families have exactly two EC assignments

that is, 10 278 : 96.0% of the (10 699) families with any EC assignment have either one or two EC assignments

## Families with one EC assignment :

What is the distribution of the percentage of sequences assigned the EC?

**Families with one EC assignment (histogram)**

number of families vs. percentage of sequences assigned the EC

# Families with one and two EC assignements, cont'd

### Families with one EC assignment :

In fact, of the 9120 families with one EC assignment, 3249 : 35.6% of these families are completely assigned with its EC

# Families with one and two EC assignements, cont'd

### Families with one EC assignment :

In fact, of the 9120 families with one EC assignment, 3249 : 35.6% of these families are completely assigned with its EC

- this is 30.4% of the (10 699) families with any EC assignment

# Families with one and two EC assignements, cont'd

### Families with one EC assignment :

In fact, of the 9120 families with one EC assignment, 3249 : <span style="color:red">35.6%</span> of these families are completely assigned with its EC

- this is 30.4% of the (10 699) families with any EC assignment

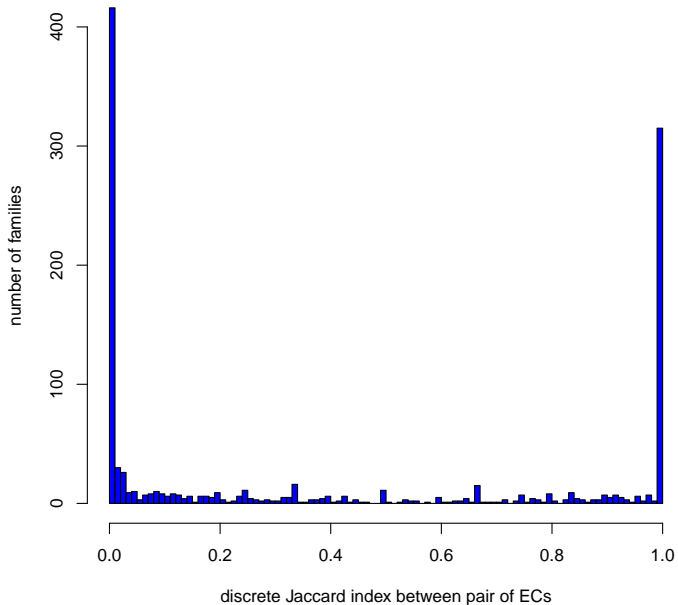### Families with two EC assignments :

What is the distribution of the <span style="color:red">Jaccard index</span>,

$$\frac{|A \cap B|}{|A \cup B|}$$

of the pair of EC assignments for a family, where $A$ (resp., $B$) is the set of sequences assigned with one (resp., the other) EC?

**Families with two EC assignments (histogram)**

number of families

discrete Jaccard index between pair of ECs

# Families with two EC assignments

In fact, of the 1156 families with two EC assignments, 318 : 27.5% (resp., 309 : 26.7%) have Jaccard index 0 (resp., 1)

# Families with two EC assignments

In fact, of the 1156 families with two EC assignments, 318 : 27.5% (resp., 309 : 26.7%) have Jaccard index 0 (resp., 1)

- ▶ What about the families in the middle?

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- ▶ for which there is a (an ultrametric) species tree on 1460 of these organisms

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- ▶ for which there is a (an ultrametric) species tree on 1460 of these organisms
- ▶ from which we consider a subtree on 1452 organisms, due to incomplete genome sequences, or non-species (i.e., mitochondria, plasmids or nucleomorphs)
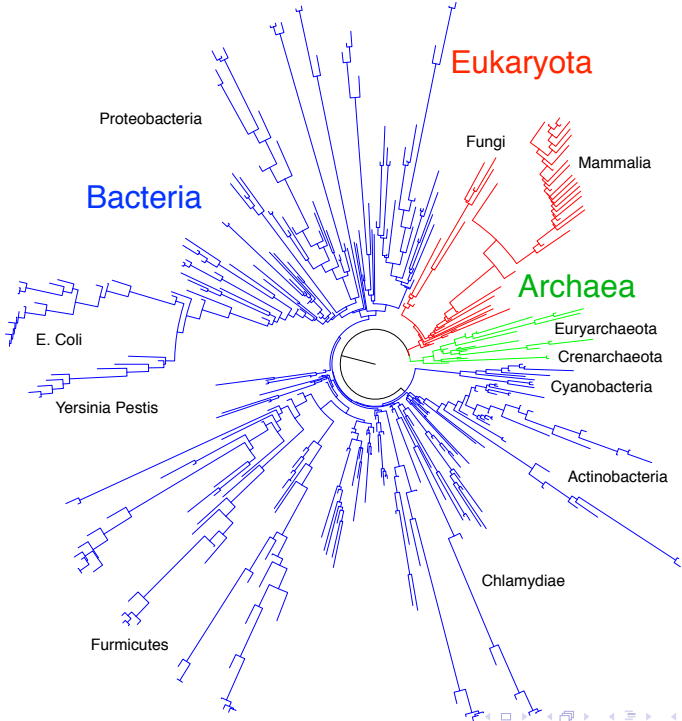
# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- ▶ for which there is a (an ultrametric) species tree on 1460 of these organisms
- ▶ from which we consider a subtree on 1452 organisms, due to incomplete genome sequences, or non-species (i.e., mitochondria, plasmids or nucleomorphs)

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- for which there is a (an ultrametric) species tree on 1460 of these organisms
- from which we consider a subtree on 1452 organisms, due to incomplete genome sequences, or non-species (i.e., mitochondria, plasmids or nucleomorphs)

For a collection of protein sequences (as encoded by the genome of an organism), the PRIAM search tool

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- for which there is a (an ultrametric) species tree on 1460 of these organisms
- from which we consider a subtree on 1452 organisms, due to incomplete genome sequences, or non-species (i.e., mitochondria, plasmids or nucleomorphs)

For a collection of protein sequences (as encoded by the genome of an organism), the PRIAM search tool

- screens the sequences with the PRIAM profiles (using rps-blast) to get a set of hits, with associated proabilities, for each profile

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- for which there is a (an ultrametric) species tree on 1460 of these organisms
- from which we consider a subtree on 1452 organisms, due to incomplete genome sequences, or non-species (i.e., mitochondria, plasmids or nucleomorphs)

For a collection of protein sequences (as encoded by the genome of an organism), the PRIAM search tool

- screens the sequences with the PRIAM profiles (using rps-blast) to get a set of hits, with associated proabilities, for each profile
- applies a set of EC-specific logical rules to the (sets of) hits of all profiles that concern the given EC to infer an overall probability for each enzyme in the collection (i.e., organism)

# Hogenom Organisms

The gene families of HOGENOM 6 fall into 1471 organisms (compose 1471 genome sequences)

- for which there is a (an ultrametric) species tree on 1460 of these organisms
- from which we consider a subtree on 1452 organisms, due to incomplete genome sequences, or non-species (i.e., mitochondria, plasmids or nucleomorphs)

For a collection of protein sequences (as encoded by the genome of an organism), the PRIAM search tool

- screens the sequences with the PRIAM profiles (using rps-blast) to get a set of hits, with associated proabilities, for each profile
- applies a set of EC-specific logical rules to the (sets of) hits of all profiles that concern the given EC to infer an overall probability for each enzyme in the collection (i.e., organism)

Here we apply PRIAM search to the organisms of HOGENOM 6

# Evolutionary scenarios of metabolic functions

## Construction
A given EC $x$ is at each leaf of the species tree with a certain probability

# Evolutionary scenarios of metabolic functions

## Construction

A given EC $x$ is at each leaf of the species tree with a certain probability

- ▶ use MapNH[4] to infer gain (resp., loss) probabilities $p_b(x), b \in \mathcal{B}$ of this EC $x$ on the branches $\mathcal{B}$ of the species tree

[4]biopp.univ-montp2.fr/forge/testnh

# Evolutionary scenarios of metabolic functions

## Construction

A given EC $x$ is at each leaf of the species tree with a certain probability

- ▶ use MapNH[4] to infer gain (resp., loss) probabilities $p_b(x), b \in \mathcal{B}$ of this EC $x$ on the branches $\mathcal{B}$ of the species tree

## Analysis

We combine the probabilities on the branches of the tree to obtain

---

[4]biopp.univ-montp2.fr/forge/testnh

# Evolutionary scenarios of metabolic functions

### Construction

A given EC $x$ is at each leaf of the species tree with a certain probability

- use MapNH[4] to infer gain (resp., loss) probabilities $p_b(x), b \in \mathcal{B}$ of this EC $x$ on the branches $\mathcal{B}$ of the species tree

### Analysis

We combine the probabilities on the branches of the tree to obtain

- $p(x) = \sum_{b \in \mathcal{B}} p_b(x)$ : the overall gain (resp., loss) probability of EC $x$ on the tree,

---

[4]biopp.univ-montp2.fr/forge/testnh

# Evolutionary scenarios of metabolic functions

## Construction

A given EC $x$ is at each leaf of the species tree with a certain probability

- use MapNH[4] to infer gain (resp., loss) probabilities $p_b(x), b \in \mathcal{B}$ of this EC $x$ on the branches $\mathcal{B}$ of the species tree

## Analysis

We combine the probabilities on the branches of the tree to obtain

- $p(x) = \sum_{b \in \mathcal{B}} p_b(x)$ : the overall gain (resp., loss) probability of EC $x$ on the tree,

- $p(x, y) = \sum_{b \in \mathcal{B}} p_b(x) \cdot p_b(y)$ : the overall *joint* gain (resp., loss) of ECs $x$ and $y$ on the tree

---

[4]biopp.univ-montp2.fr/forge/testnh

# Evolutionary scenarios of metabolic functions

## Construction

A given EC $x$ is at each leaf of the species tree with a certain probability

- use MapNH[4] to infer gain (resp., loss) probabilities $p_b(x), b \in \mathcal{B}$ of this EC $x$ on the branches $\mathcal{B}$ of the species tree
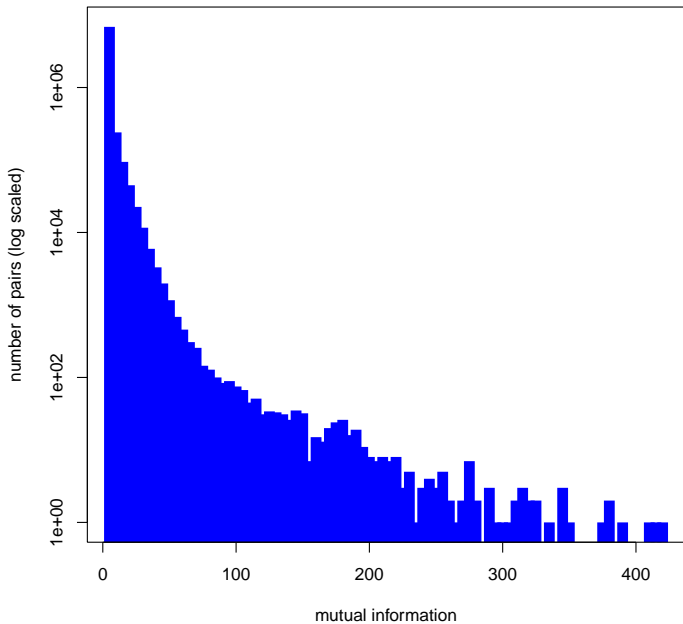
## Analysis

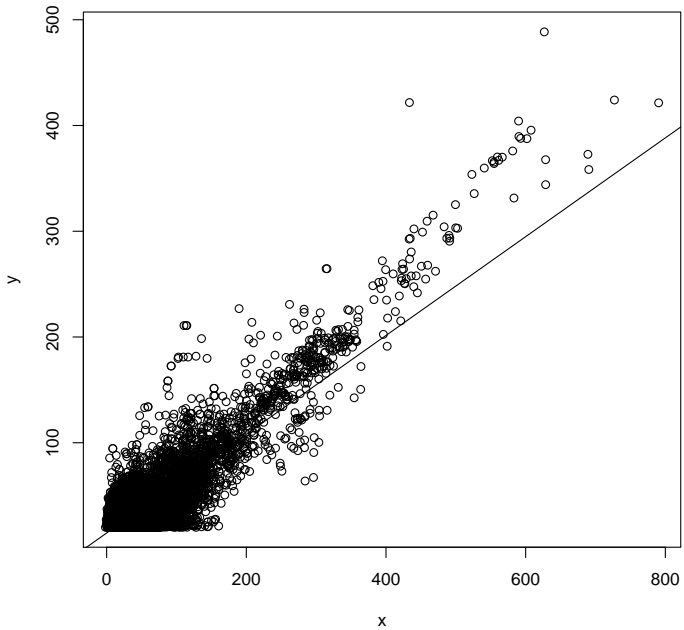We combine the probabilities on the branches of the tree to obtain

- $p(x) = \sum_{b \in \mathcal{B}} p_b(x)$ : the overall gain (resp., loss) probability of EC $x$ on the tree,

- $p(x, y) = \sum_{b \in \mathcal{B}} p_b(x) \cdot p_b(y)$ : the overall *joint* gain (resp., loss) of ECs $x$ and $y$ on the tree

Compute the mutual information (MI) for each pair of ECs $x$ and $y$ :

$$\mathsf{MI}(x, y) = \sum_{x' \in \{x, \bar{x}\}} \sum_{y' \in \{y, \bar{y}\}} p(x', y') \log \left( \frac{p(x', y')}{p(x')p(y')} \right)$$

---

[4]biopp.univ-montp2.fr/forge/testnh

**Mutual information for a pair of ECs (histogram)**

number of pairs (log scaled)

mutual information

# Community Analysis of Mutual Information

## The Graph

| threshold | nodes (isolated) | edges | avg. degree |
|-----------|------------------|-------|-------------|
| 5.0E-3 | 1124 (+ 1608 = 2732) | 22005 | 39.2 |
| 5.0E-4 | 974 (+ 1758 = 2732) | 11792 | 24.2 |
| 5.0E-5 | 836 (+ 1896 = 2732) | 6962 | 16.7 |

# Community Analysis of Mutual Information

## The Graph

| threshold | nodes (isolated) | edges | avg. degree |
|-----------|------------------|-------|-------------|
| 5.0E-3 | 1124 (+ 1608 = 2732) | 22005 | 39.2 |
| 5.0E-4 | 974 (+ 1758 = 2732) | 11792 | 24.2 |
| 5.0E-5 | 836 (+ 1896 = 2732) | 6962 | 16.7 |

## Node Communities                              (Blondel et al., 2008)

| threshold | number | # /w 2 ECs | # /w 3 ECs | mean size | max size |
|-----------|--------|------------|------------|-----------|----------|
| 5.0E-3 | 40 | 22 (55%) | 6 (15%) | 28.1 | 335 |
| 5.0E-4 | 43 | 26 (60%) | 7 (16%) | 22.7 | 268 |
| 5.0E-5 | 53 | 29 (55%) | 6 (11%) | 15.8 | 223 |

# Community Analysis of Mutual Information

## The Graph

| threshold | nodes (isolated) | edges | avg. degree |
|---|---|---|---|
| 5.0E-3 | 1124 (+ 1608 = 2732) | 22005 | 39.2 |
| 5.0E-4 | 974 (+ 1758 = 2732) | 11792 | 24.2 |
| 5.0E-5 | 836 (+ 1896 = 2732) | 6962 | 16.7 |

## Node Communities                      (Blondel et al., 2008)

| threshold | number | # /w 2 ECs | # /w 3 ECs | mean size | max size |
|---|---|---|---|---|---|
| 5.0E-3 | 40 | 22 (55%) | 6 (15%) | 28.1 | 335 |
| 5.0E-4 | 43 | 26 (60%) | 7 (16%) | 22.7 | 268 |
| 5.0E-5 | 53 | 29 (55%) | 6 (11%) | 15.8 | 223 |

## Link Communities                      (Ahn et al., 2010)

| threshold | number | # /w 2 ECs | # /w 3 ECs | mean size | max size |
|---|---|---|---|---|---|
| 5.0E-3 | 5571 | 4480 (80%) | 473 (8%) | 3.1 | 129 |
| 5.0E-4 | 3084 | 2381 (77%) | 318 (10%) | 3.0 | 107 |
| 5.0E-5 | 2043 | 1554 (76%) | 265 (13%) | 2.7 | 93 |

METHANE METABOLISM

Chistoserdova et al., 2003

# So, what's next?

### Reconciled Trees

Perform the same study on a Cyanobacteria dataset reconciled by the method of Szöllősi et al., (2013)

# So, what's next?

### Reconciled Trees
Perform the same study on a Cyanobacteria dataset reconciled by the method of Szöllősi et al., (2013)

- ▶ we will have (more detailed) evolutionary scenarios for ECs,

# So, what's next?

### Reconciled Trees
Perform the same study on a Cyanobacteria dataset reconciled by the method of Szöllősi et al., (2013)

- ▶ we will have (more detailed) evolutionary scenarios for ECs,
- ▶ but also ancestral sequences (idea : run PRIAM on these and compare to the ECs scenarios)

# So, what's next?

### Reconciled Trees

Perform the same study on a Cyanobacteria dataset reconciled by the method of Szöllősi et al., (2013)

- ▶ we will have (more detailed) evolutionary scenarios for ECs,
- ▶ but also ancestral sequences (idea : run PRIAM on these and compare to the ECs scenarios)

### Stoichiometry Analysis

Given a collection of enzymes (ECs), PRIAM can construct a draft metabolic network, using Kegg pathway data

# So, what's next?

### Reconciled Trees

Perform the same study on a Cyanobacteria dataset reconciled by the method of Szöllősi et al., (2013)

- ▶ we will have (more detailed) evolutionary scenarios for ECs,
- ▶ but also ancestral sequences (idea : run PRIAM on these and compare to the ECs scenarios)

### Stoichiometry Analysis

Given a collection of enzymes (ECs), PRIAM can construct a draft metabolic network, using Kegg pathway data

- ▶ construct extant, and ancestral networks

# So, what's next?

### Reconciled Trees
Perform the same study on a Cyanobacteria dataset reconciled by the method of Szöllősi et al., (2013)

- ▶ we will have (more detailed) evolutionary scenarios for ECs,
- ▶ but also ancestral sequences (idea : run PRIAM on these and compare to the ECs scenarios)

### Stoichiometry Analysis
Given a collection of enzymes (ECs), PRIAM can construct a draft metabolic network, using Kegg pathway data

- ▶ construct extant, and ancestral networks
- ▶ perform a comparison from a reaction stoichiometry point of view (Poolman et al., 2007)

# Acknowledgements

- Ancestrome Project and UCBL (funding, etc.)
- Simon Penel and Vincent Daubin (hogenom)
- Laurent Guéguen and Julien Dutheil (bio++)
- Dominique Guyot (paraload)
- Gergely Szöllősi (reconciled trees)

# Thank you!

Any questions?

# The goals of this project

We outline a few of the main goals of this project :

# The goals of this project

We outline a few of the main goals of this project :

1. to investigate principles of metabolic network evolution : how is function related to evolution? How is evolution determined by function?

# The goals of this project

We outline a few of the main goals of this project :

1. to investigate principles of metabolic network evolution : how is function related to evolution? How is evolution determined by function?

2. in the context of the Ancestrome project : to introduce functional dependencies into the likelihood calculations

# The goals of this project

We outline a few of the main goals of this project :

1. to investigate principles of metabolic network evolution : how is function related to evolution? How is evolution determined by function?

2. in the context of the Ancestrome project : to introduce functional dependencies into the likelihood calculations

3. evaluating hypotheses about ancestral environments : metabolism sheds light on environmental factors, which could provide clues on the events associated with the emergence of ancestral phyla

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

## From the evolutionary tree

Determine the (pairwise) relationships between ECs using mutual information

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

## From the evolutionary tree

Determine the (pairwise) relationships between ECs using mutual information

- the mutual information between two (discrete) random variables $X$ and $Y$ is defined as :

$$I(X;Y) = \sum_{y \in Y} \sum_{x \in X} p(x,y) \log \left( \frac{p(x,y)}{p(x)p(y)} \right)$$

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

## From the evolutionary tree

Determine the (pairwise) relationships between ECs using mutual information

- the mutual information between two (discrete) random variables $X$ and $Y$ is defined as :

$$I(X;Y) = \sum_{y \in Y} \sum_{x \in X} p(x,y) \log \left( \frac{p(x,y)}{p(x)p(y)} \right)$$

- from here we can build a large graph that represents these pairwise dependencies, and try to find modules within this graph

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

## From the structure of the networks

Conversely, we may detect functional dependencies purely from the structure of a metabolic network (i.e., at a given node of the species tree)

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

## From the structure of the networks

Conversely, we may detect functional dependencies purely from the structure of a metabolic network (i.e., at a given node of the species tree)

- from a reaction stoichiometry there exist methods for finding correlations among subsets of reactions (Poolman et al., 2007)

# Dependencies between metabolic functions

Our goal is to find (and to understand) the dependencies between different metabolic functions within these evolutionary scenarios

## From the structure of the networks
Conversely, we may detect functional dependencies purely from the structure of a metabolic network (i.e., at a given node of the species tree)

- from a reaction stoichiometry there exist methods for finding correlations among subsets of reactions (Poolman et al., 2007)

- from this, we get dependencies between sets of reactions, functions (their ECs)

# Reconstructing evolutionary scenarios

For each of the individual family and whole genome viewpoints, there are again two viewpoints on the reconstruction of evolutionary scenarios – in terms of either

# Reconstructing evolutionary scenarios

For each of the individual family and whole genome viewpoints, there are again two viewpoints on the reconstruction of evolutionary scenarios – in terms of either

## ECs
Apply a method (parsimony or ML) to propagate the ECs to the ancestral nodes of the tree

- the collection of ECs at the ancestral nodes will then determine the functions active at these nodes

# Reconstructing evolutionary scenarios

For each of the individual family and whole genome viewpoints, there are again two viewpoints on the reconstruction of evolutionary scenarios – in terms of either

## ECs

Apply a method (parsimony or ML) to propagate the ECs to the ancestral nodes of the tree

- the collection of ECs at the ancestral nodes will then determine the functions active at these nodes

## Nucelic acid sequences

Apply a method (parsimony or ML) to propagate the sequences to the ancestral nodes of the tree

- we can apply PRIAM search on these ancestral sequences (just like we did for the extant ones) to get collections of ancestral ECs at the ancestral nodes
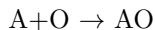
# reaction stoichimetry

|       | $s_1$ | $s_2$ |
|-------|-------|-------|
| $e_a$ | 0.76  | 0.45  |
| $e_b$ | 0.39  | 0.87  |
| $e_c$ | 0.68  | 0.23  |

## reaction stoichimetry

|       | $s_1$ | $s_2$ |
|-------|-------|-------|
| $e_a$ | 0.76  | 0.45  |
| $e_b$ | 0.39  | 0.87  |
| $e_c$ | 0.68  | 0.23  |

$\longrightarrow$

| EC 1. | $A+O \rightarrow AO$ |
|-------|----------------------|
| EC 2. | $AB+C \rightarrow A+BC$ |
| EC 3. | $AB+H_2O \rightarrow AOH+BH$ |
| EC 6. | $X+Y+ATP \rightarrow XY+ADP+Pi$ |

# reaction stoichimetry

|       | $s_1$ | $s_2$ |
|-------|-------|-------|
| $e_a$ | 0.76  | 0.45  |
| $e_b$ | 0.39  | 0.87  |
| $e_c$ | 0.68  | 0.23  |

| EC 1. | $A + O \rightarrow AO$ |
|-------|------------------------|
| EC 2. | $AB + C \rightarrow A + BC$ |
| EC 3. | $AB + H_2O \rightarrow AOH + BH$ |
| EC 6. | $X + Y + ATP \rightarrow XY + ADP + Pi$ |