



Application of the Lempel-Ziv complexity to the alignment-free sequence comparison of protein families

Sofiène Bacha and Denis Baurain
Université de Liège, Belgium



■ **Background**

- **Issues with multiple alignment**
- **Alignment-free sequence comparison**
- **Lempel-Ziv complexity**
- **LZ distance metrics**

■ **Methods**

- **Encoding schemes**
- **Benchmarking strategy**

■ **Results**

- **Benchmark graphs**
- **Addendum: decision trees**
- **Application to phylogenetics**

■ **Discussion**

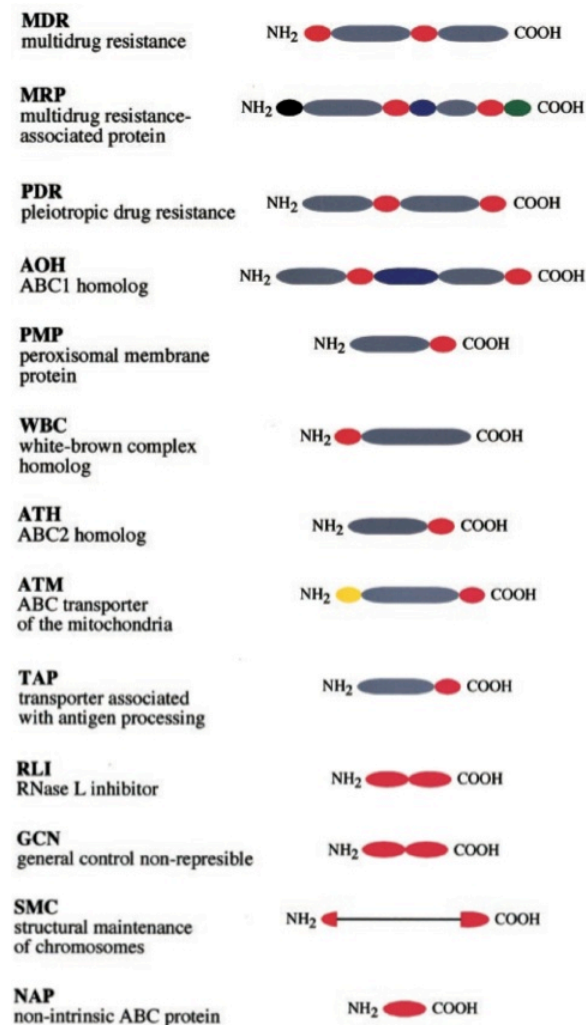
- **Performance considerations**
- **Conclusions and perspectives**



Background

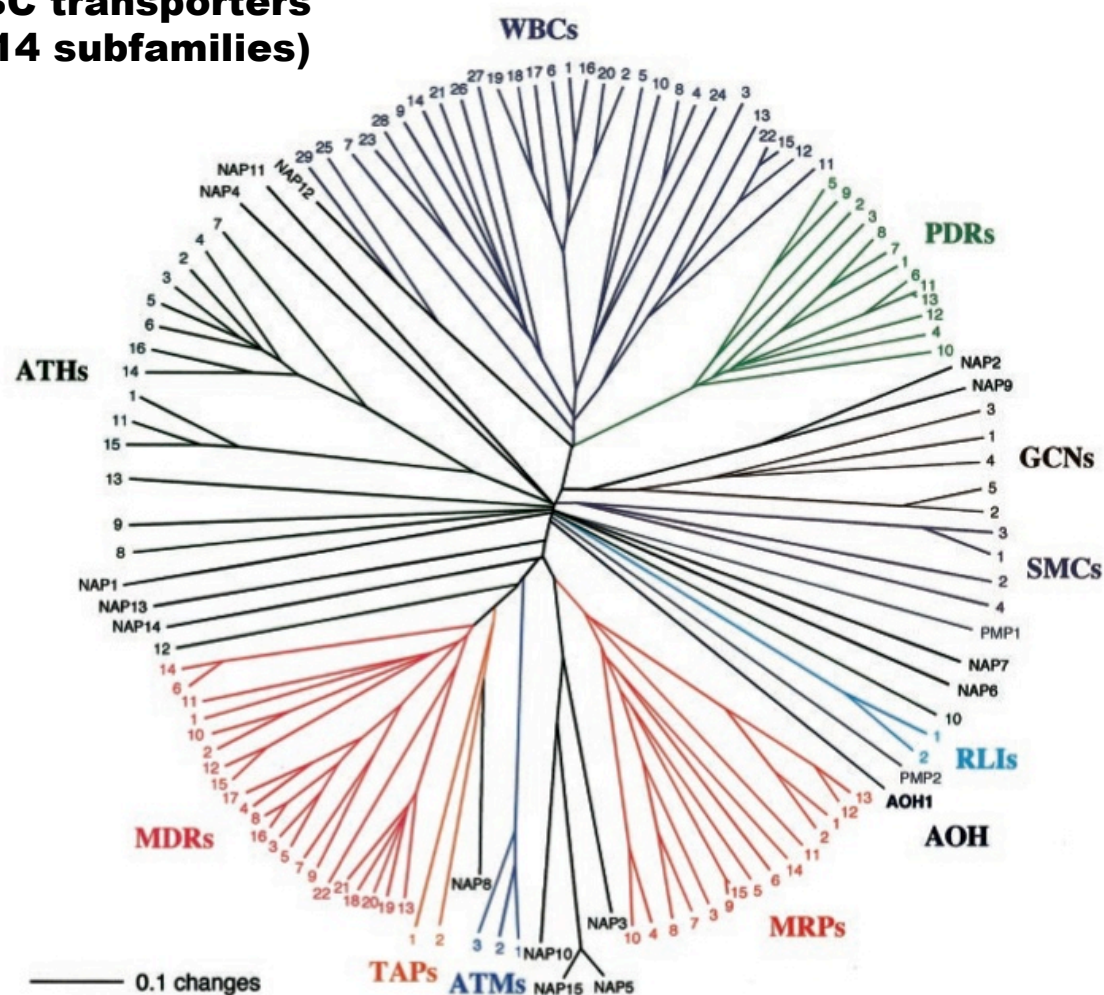
Issues with multiple alignment

- multiprotein family phylogenies are ***NOT*** organism phylogenies
- paralogues vs orthologues
- paralogues often exhibit domain shuffling and/or domain duplication
- rapid diversification may follow gene duplication and then give rise to numerous subfamilies (e.g. ABC transporters)

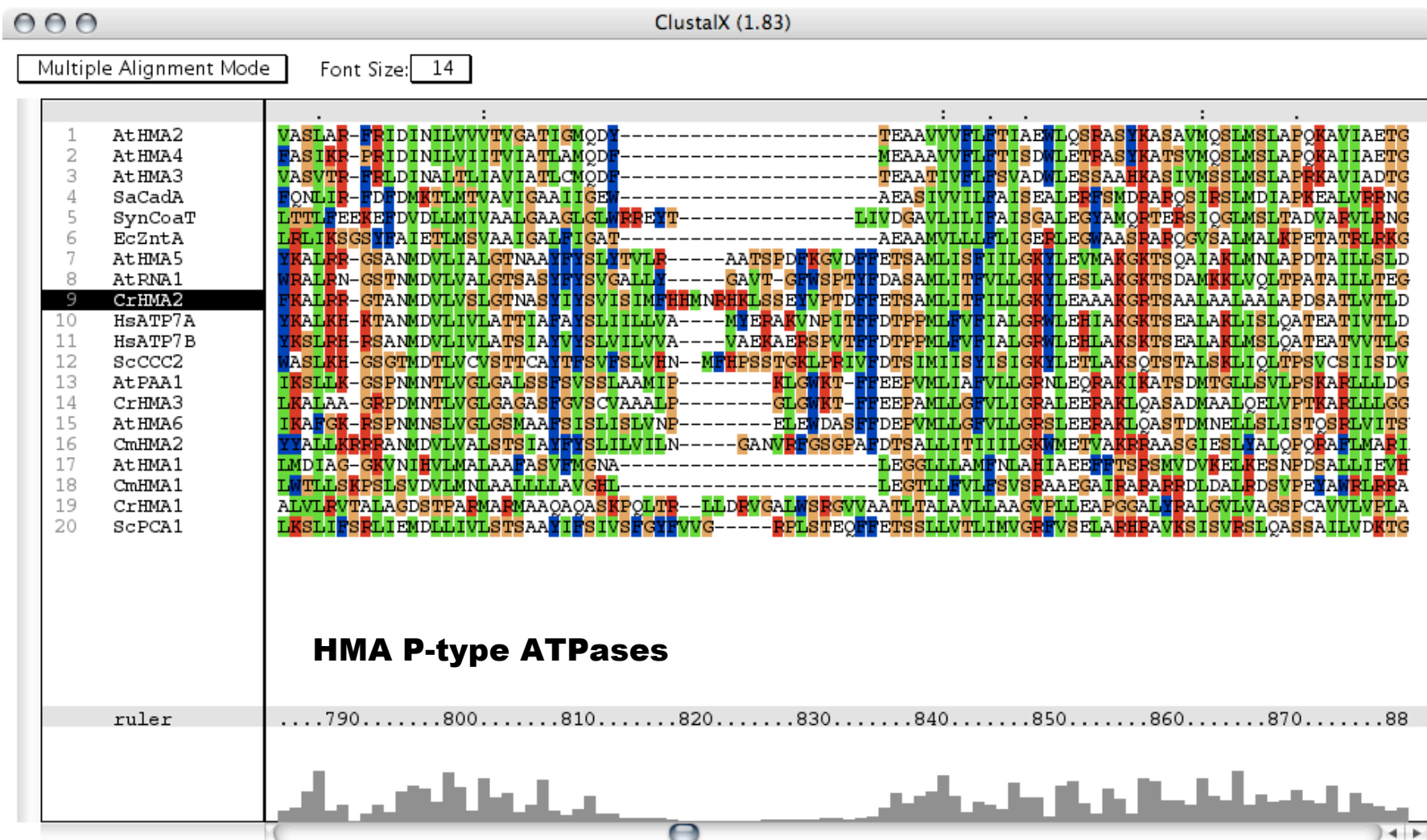


Issues with multiple alignment

***Arabidopsis* ABC transporters
(129 proteins, 14 subfamilies)**



Issues with multiple alignment



Alignment-free sequence comparison

■ word statistics

- vectors of counts or frequencies of k -mers
- e.g. squared Euclidean distance:

$$d_k^E(S, Q) = \sum_{i=1}^K (c_{k,i}^S - c_{k,i}^Q)^2$$

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
seq S	3	2	1	2	3	7	3	3	5	8	5	3	5	2	9	4	3	4	2	3
seq Q	4	2	1	2	2	6	3	2	6	8	5	3	7	2	6	4	4	3	2	4

■ information theory

- algorithmic complexity
- estimation through sequence compression

k=1; K=20

Lempel-Ziv complexity

■ Exhaustive history H_E

- $S = \text{AACGTACCATTG}$
- $H_E(S) = A \cdot AC \cdot G \cdot T \cdot ACC \cdot AT \cdot TG$
- $Q = \text{ACGGTCACCAA}$
- $H_E(Q) = A \cdot C \cdot G \cdot GT \cdot CA \cdot CC \cdot AA$

■ The LZ complexity is the number of components in H_E

- $c(S) = c_E(S) = 7$
- $c(Q) = 7$

■ Subadditivity of the LZ complexity

- $SQ = \text{AACGTACCATTGACGGTCACCAA}$
- $H_E(SQ) = A \cdot AC \cdot G \cdot T \cdot ACC \cdot AT \cdot TG \cdot ACGG \cdot TC \cdot ACCAA; c(SQ) = 10$
- $c(SQ) \leq c(S) + c(Q)$ [indeed, $10 < 7 + 7$]

■ LZ complexity and sequence similarity

- $R = \text{CTAGGGACTTAT}$
- $H_E(R) = C \cdot T \cdot A \cdot G \cdot GGA \cdot CTT \cdot AT; c(R) = 7$
- $H_E(RQ) = C \cdot T \cdot A \cdot G \cdot GGA \cdot CTT \cdot AT \cdot ACG \cdot GT \cdot CA \cdot CC \cdot AA; c(RQ) = 12$
- $c(SQ) < c(RQ)$ because S is closer [more similar] to Q than to R [e.g. ACG and $ACCA$]

LZ distance metrics

$$d(S, Q) = \max\{c(SQ) - c(S), c(QS) - c(Q)\} \quad (1)$$

$$d^*(S, Q) = \frac{\max\{c(SQ) - c(S), c(QS) - c(Q)\}}{\max\{c(S), c(Q)\}} \quad (2)$$

$$d_1(S, Q) = c(SQ) - c(S) + c(QS) - c(Q) \quad (3)$$

$$d_1^*(S, Q) = \frac{c(SQ) - c(S) + c(QS) - c(Q)}{c(SQ)} \quad (4)$$

$$d_1^{**}(S, Q) = \frac{c(SQ) - c(S) + c(QS) - c(Q)}{\frac{1}{2}[c(SQ) + c(QS)]} \quad (5)$$



Methods



Encoding schemes

■ DNA sequences

- the *exact match* approach of the LZ distance metrics works well with the small and simple DNA alphabet

■ AA sequences

- LZ dm are expected to miss the subtle and overlapping *similarities* characterizing the larger and more complex AA-alphabet
- *substitution matrices* (e.g. PAM, BLOSUM...) are not applicable since LZ dm does not compare residues on a *pairwise* basis

■ strategy

- we propose several variants of a simple approach where AA sequences are *encoded* to different alphabets *prior to the computation* of the LZ complexity


■ key idea

- *to capture as much information as possible* in order to enhance the alignment-free sequence comparison of proteins

Encoding schemes


1. binary codons

	T	C	A	G	T	C	A	G	T	C	A	G
Met (M)	0	0	1	0	1	0	0	0	0	0	0	1
Cys (C)	1	0	0	0	0	0	0	1	1	1	0	0
Ile (I)	0	0	1	0	1	0	0	0	1	1	1	0
Pro (P)	0	1	0	0	0	1	0	0	1	1	1	1
Arg (R)	0	1	1	0	0	0	0	1	1	1	1	1

 12 bytes


2. binary codons; 3rd position discarded

	T	C	A	G	T	C	A	G
Met (M)	0	0	1	0	1	0	0	0
Cys (C)	1	0	0	0	0	0	0	1
Ile (I)	0	0	1	0	1	0	0	0
Pro (P)	0	1	0	0	0	1	0	0
Arg (R)	0	1	1	0	0	0	0	1

 8 bytes


4. alphanumeric codons

Met (M)	A	T	G															
Cys (C)	T	G	T	T	G	C												
Ile (I)	A	T	T	A	T	C	A	T	A									
Pro (P)	C	C	T	C	C	C	C	C	A	C	C	G						
Arg (R)	C	G	T	C	G	C	C	G	A	C	G	G	A	G	A	A	G	G

 3 to 18 bytes

5. alphanumeric codons; compressed

Met (M)	A	T	G					
Cys (C)	T	G	T	C				
Ile (I)	A	T	T	C	A			
Pro (P)	C	C	T	C	A	G		
Arg (R)	C	A	G	T	C	A	G	
Ser (S)	T	A	C	G	T	C	A	G

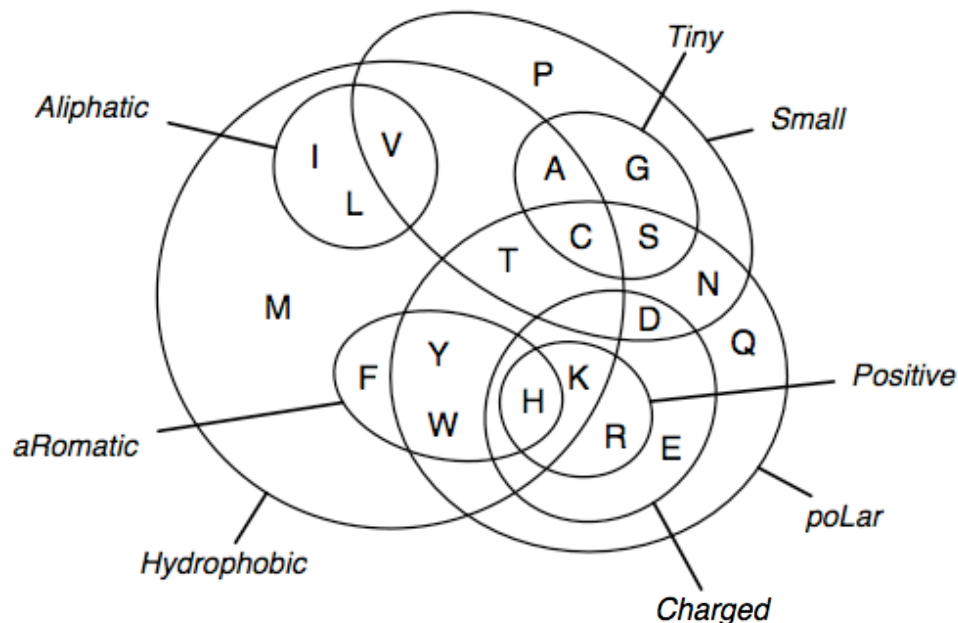
 3 to 8 bytes

Encoding schemes

3. biochemical groupings; binary sets

	T	S	P	L	C	H	R	A
Met (M)	0	0	0	0	0	1	0	0
Cys (C)	1	1	0	1	0	1	0	0
Ile (I)	0	0	0	0	0	1	0	1
Pro (P)	0	1	0	0	0	0	0	0
Arg (R)	0	0	1	1	1	0	0	0

+ 8 bytes



6. biochemical groupings; single byte

Ala (A)	Phe (F)	His (H)	Asp (D)	Asn (N)	Ser (S)	Cys (C)	Pro (P)
Gly (G)	Trp (W)	Lys (K)	Glu (E)	Gln (Q)	Thr (T)	Met (M)	
Ile (I)	Tyr (Y)	Arg (R)					
Leu (L)							
Val (V)							
<i>aLiphatic</i>	<i>aRomatic</i>	<i>Basic</i>	<i>Acidic</i>	<i>aMide</i>	<i>aLCohol</i>	<i>Sulfur</i>	<i>Proline</i>
L	R	B	A	M	C	S	P

X 1 byte

Benchmarking strategy

■ dataset

- 1,683 protein domains

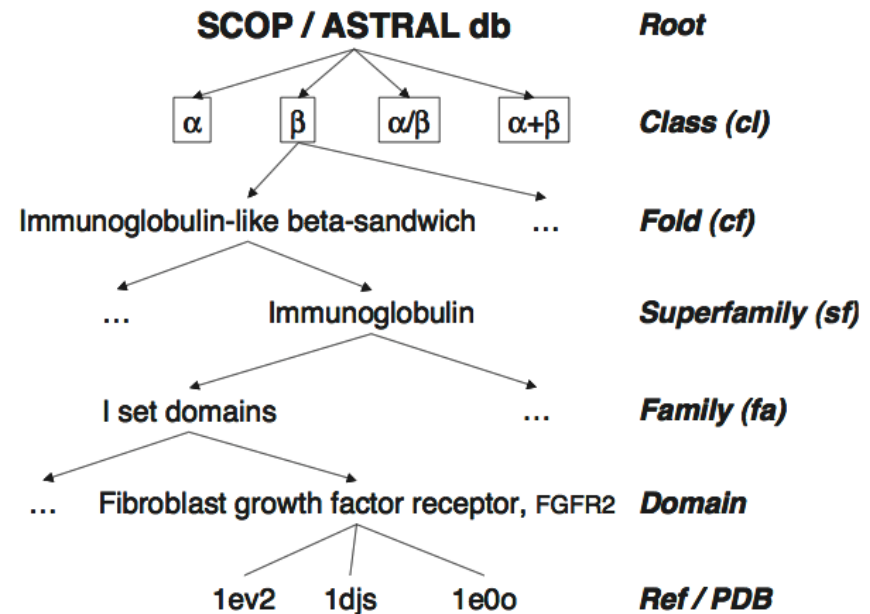
■ SCOP / ASTRAL db

- hierarchical organization based on 3D-structures
- *family* and *superfamily* levels reflect phylogeny
- *fold* and *class* levels reflect broad structure similarity

■ distance methods

- squared Euclidean distance
- W-metric (BLOSUM50)
- SW local alignment score
- LZ dm (raw and encoded)

■ ROC curves and AUCs

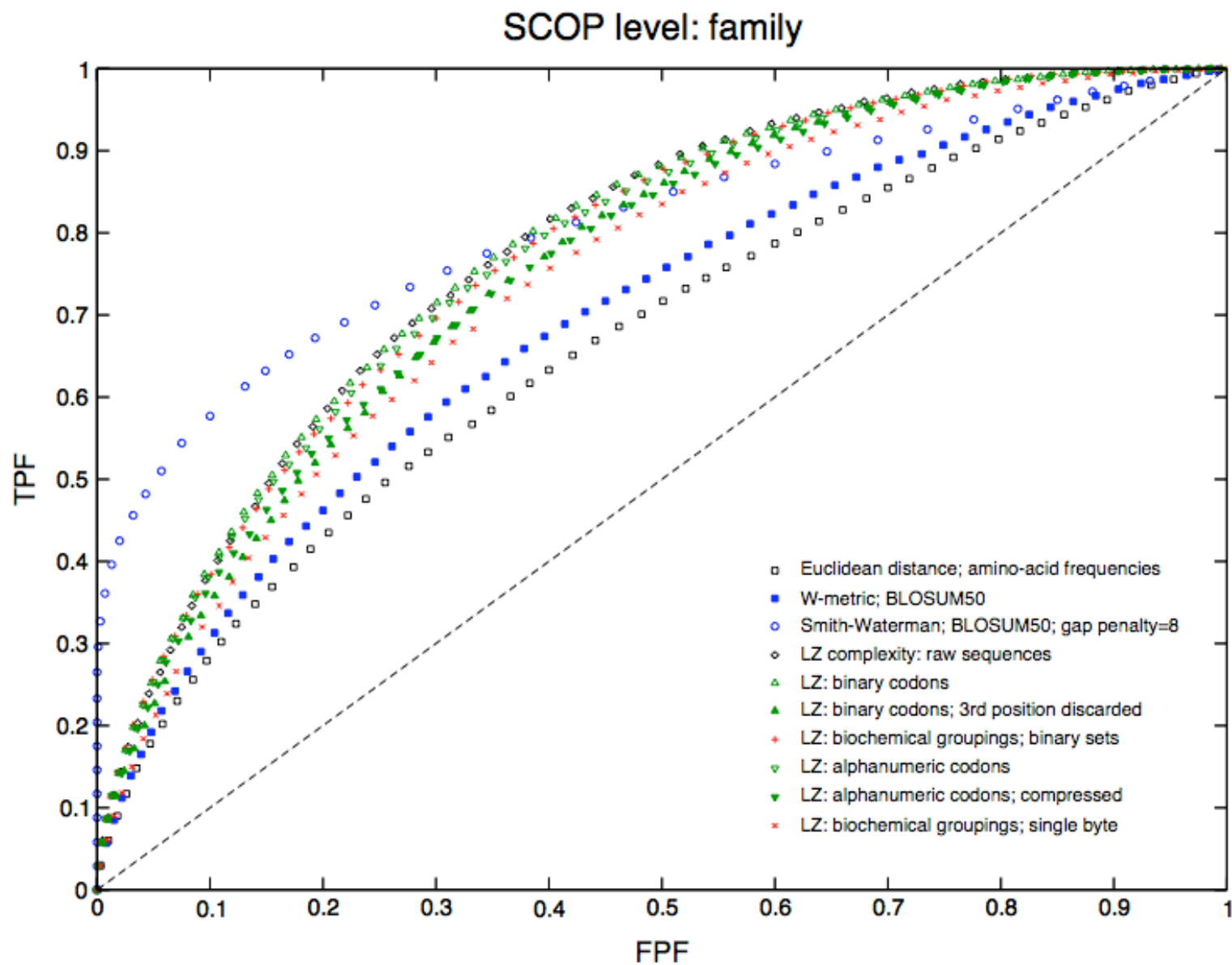


$$d^W(S, Q) = \sum_{i=1}^K \sum_{j=1}^K (f_i^S - f_i^Q) \cdot (f_j^S - f_j^Q) \cdot w_{ij}$$

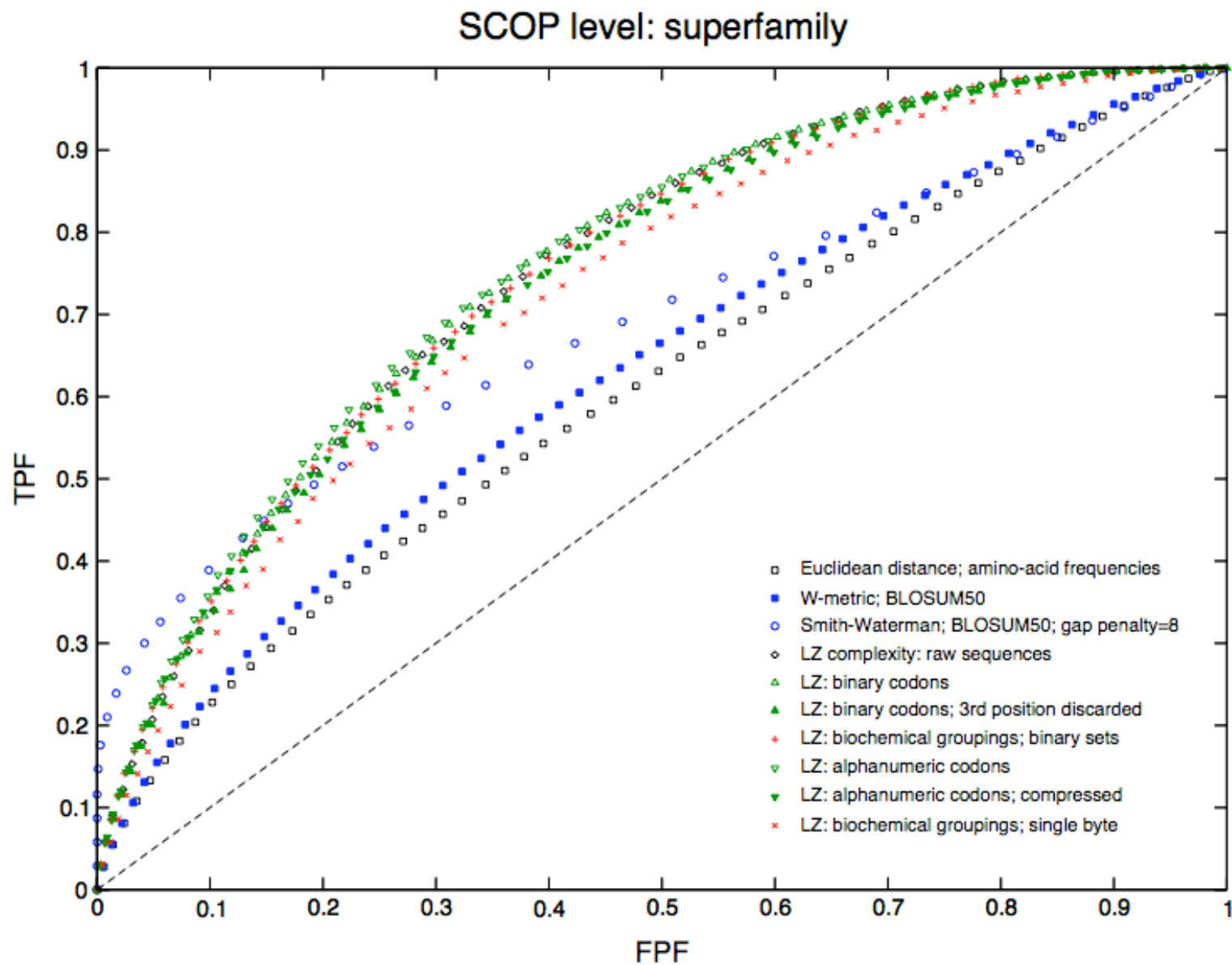


Results

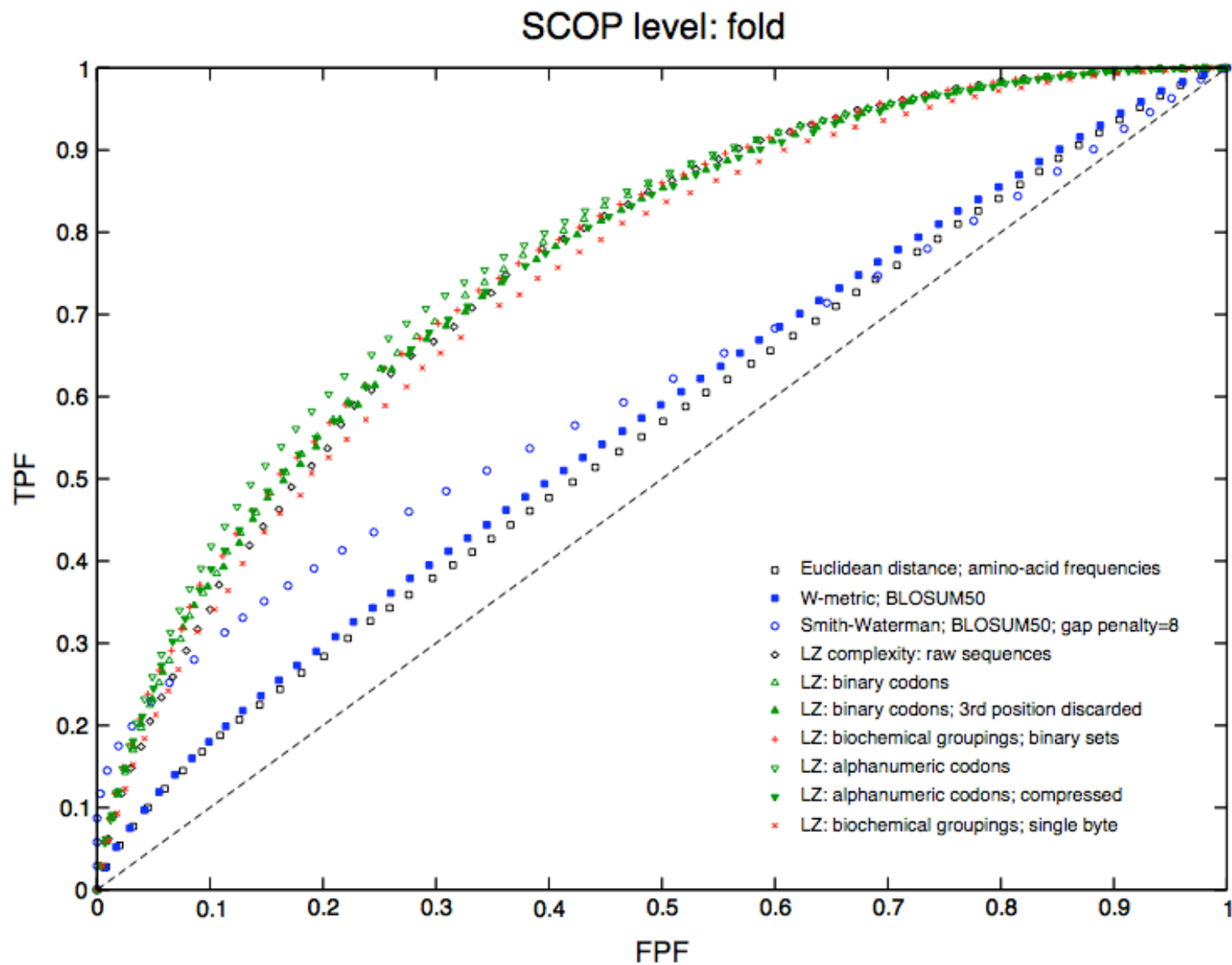
Benchmark graphs



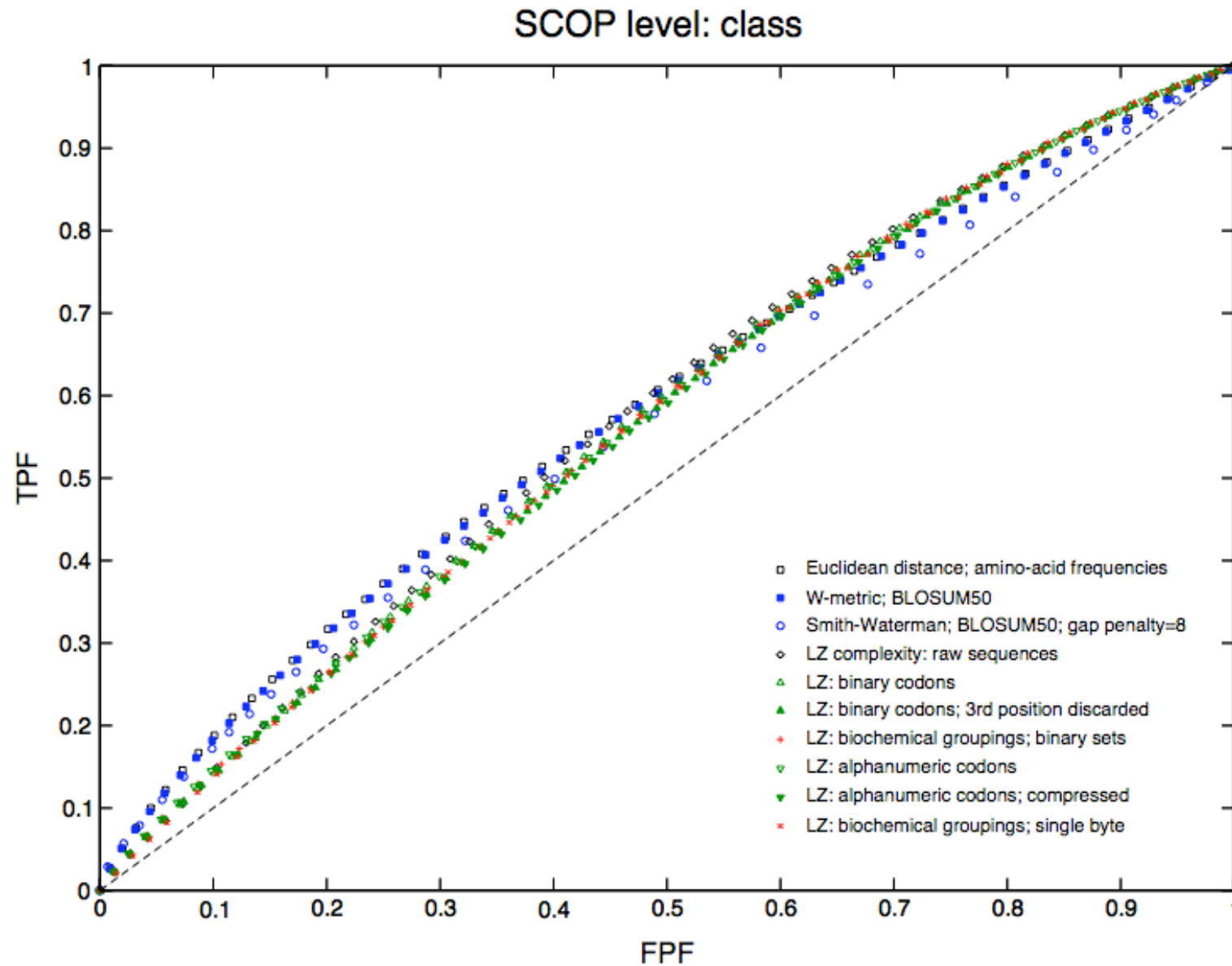
Benchmark graphs



Benchmark graphs

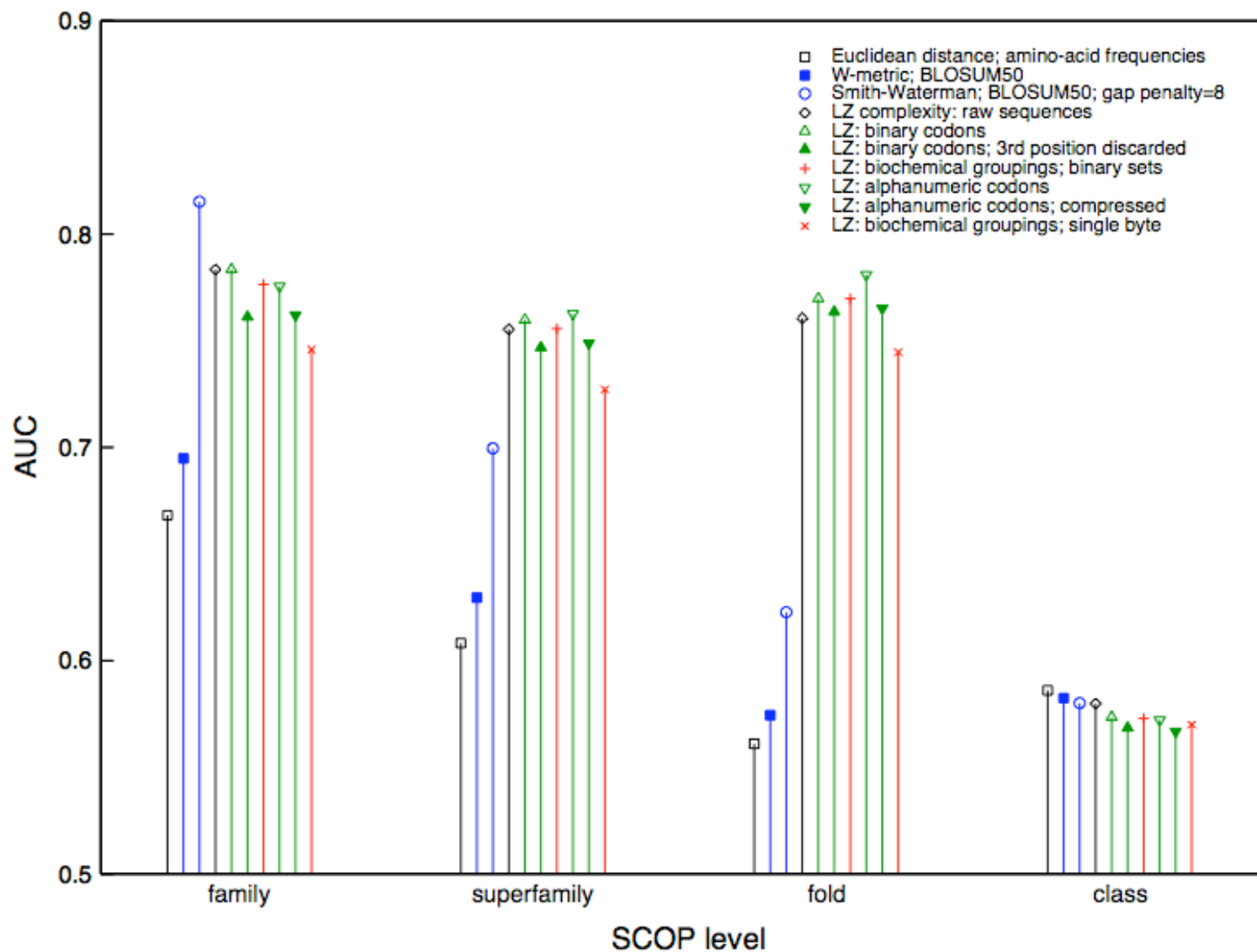


Benchmark graphs



Benchmark graphs

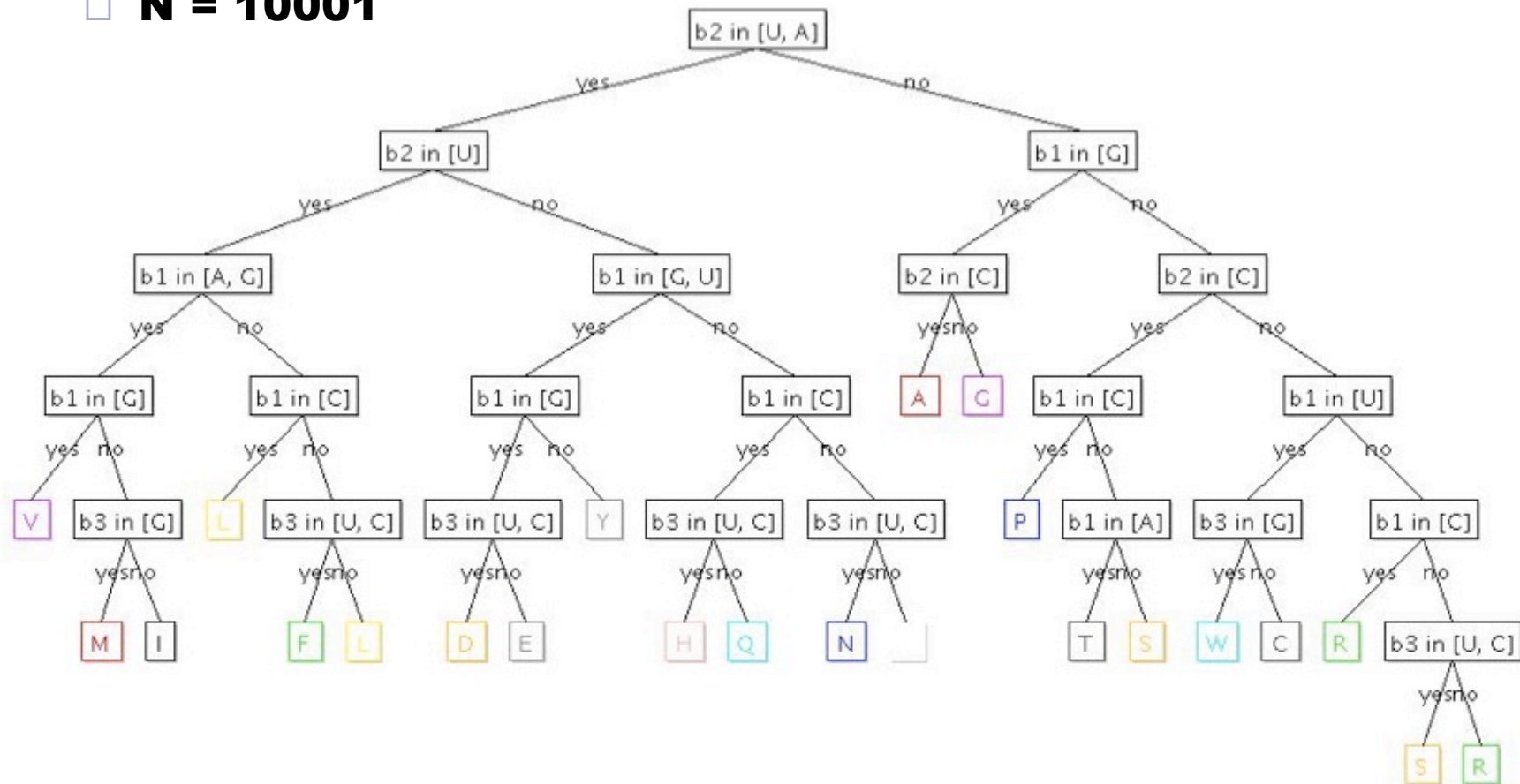
Overview: AUC values



Addendum: decision trees

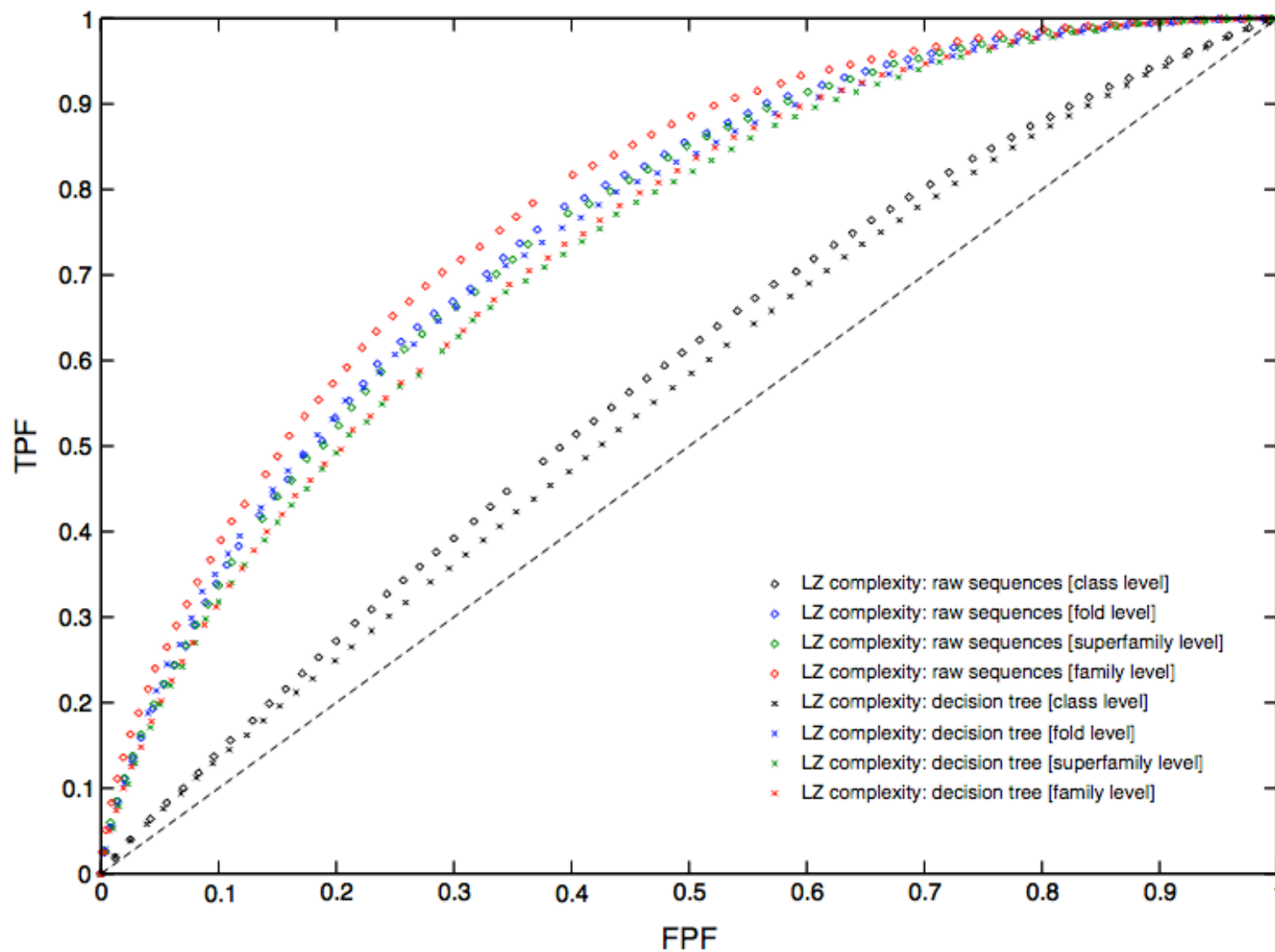
■ codon-based

- **R = 00001 (shortest path)**
- **M = 11101**
- **N = 10001**

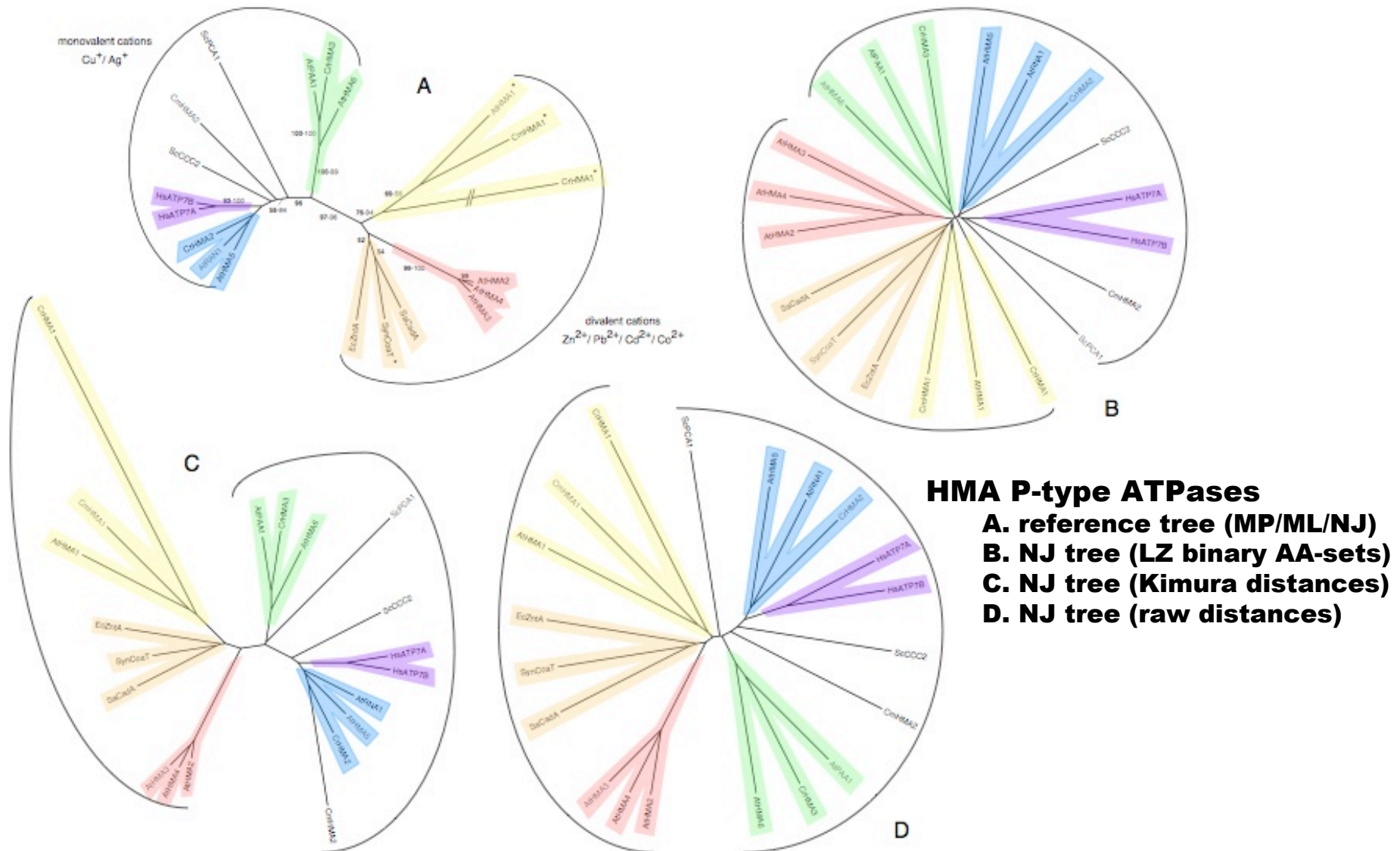


Addendum: decision trees

LZ complexity: raw sequences vs decision tree



Application to phylogenetics





Discussion



Performance considerations

	Eu	W _m	SW	LZ	LZ1	LZ2	LZ3	LZ4	LZ5	LZ6	LZ7
CPU-time	0.2	8	264	1	75	37	37	48	16	1	16
language	perl		perl/C	C							

■ one relative unit

- 7 min 45 on a PowerPC G4 at 1.25 GHz (Mac OS X)
- 6 min 10 on a Pentium 4 at 2.4 GHz (SuSE Linux)

■ implementation

- perl/c is a perl wrapper for the *water* program (written in C) of the EMBOSS software package
- our software is algorithmically optimized but not technically optimized (further optimizations on the way)



Conclusions

■ benchmarks

- while computationally affordable, the LZ complexity *outperforms all other methods* at the three lower levels of the SCOP classification, except the very slow SW local alignment at the family level
- at superfamily and fold levels, our sequence encodings show slightly *better results* than the default complexity, but at the expense of considerable *computational burden*

■ phylogenetics

- the LZ complexity is able to retrieve *most clades* found through alignment-based methods but would need some kind of *distance correction* to be really useful



Perspectives

- **refine the code in order to publish it somewhere as an *Applications Note***
 - probably not novel enough for a research paper
- **improvement of the method for its use in phylogenetic inference**
 - change reconstruction algorithm (Fitch instead of NJ)
 - fix suboptimal folding of AA biochemical groups
 - e.g. Dayhoff groups (C-ILMV-FWY-AGPST-HKR-DENQ)
 - modify the LZ complexity itself to favor large patterns
 - what about the effect on distance properties?
- **other suggestions?**