

# What's in a Likelihood? Simple models of protein evolution and the contribution of structurally viable reconstructions.

## SUPPLEMENTARY MATERIALS

Clemens Lakner<sup>1,2</sup>, Mark T. Holder<sup>3</sup>, Nick Goldman<sup>4</sup>, and Gavin J. P. Naylor<sup>2</sup>

<sup>1</sup> *Department of Biological Science, Section: Ecology and Evolution, Florida State University, Tallahassee, Florida 32306-4120, USA*

<sup>2</sup> *Department of Scientific Computing, Florida State University, Tallahassee, Florida 32306-4120, USA*

<sup>3</sup> *Department of Ecology and Evolution, 6031 Haworth, University of Kansas, 1200 Sunnyside Ave, Lawrence KS 66045*

<sup>4</sup> *European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SD, UK*

# Contents

<b>Reference Data Sets</b>	<b>3</b>
<b>Trees</b>	<b>4</b>
Two taxa . . . . .	4
Three taxa . . . . .	5
Four taxa . . . . .	6
<b>Likelihoods (lnL)</b>	<b>8</b>
Two taxa . . . . .	8
Three taxa . . . . .	9
Four taxa . . . . .	10
<b>Plots</b>	<b>11</b>
Two taxa . . . . .	11
$Z_s$ . . . . .	11
Three taxa . . . . .	12
$PE_s$ . . . . .	12
$Z_s$ . . . . .	17
Four taxa . . . . .	22
$PE_s$ . . . . .	22
$Z_s$ . . . . .	35
Rate variation . . . . .	48

## Reference Data Sets

Sequences used for the transition paths are indicated in boldface font.

1. HPPK (GenBank Accession.version)

BAB96719.1, AAP15678.1, AAL19147.1, **AAO67923.1** (*Salmonella typhi*), CAE13168.1, AAS60560.1, **CAG76218.1** (*Pectobacterium atrosepticum*), BAC95526.1, **AAF93760.1** (*Vibrio cholerae*), CAG21480.1, **YP\_001140682.1** (*Aeromonas salmonicida*);

2. Lysozyme (UniProtKB)

LYZ1\_MOUSE, **LYSC\_HUMAN** (*Homo sapiens*), LYSC\_GORGO, LYSC\_PONPY, LYSC\_HYLLA, LYZ2\_MOUSE, LYSC\_NASLA, LYSC\_SHEEP, LYSC\_SAISC, LYSC\_COLAN, LYSC\_RABIT, LYSC\_TRAOB, LYSC\_SEMEN, LYSC\_PYGNE, LYSC1\_RAT, LYSC\_MACMU, LYSC\_PAPAN, LYSC\_TRISI, LYSC2\_CANFA, LYSC\_SAGOE, LYSC\_ALLNI, LYSC\_CALJA, LYSC\_CERAE, LYSM\_BOVIN, LYSC\_MIOTA, LYSC2\_RAT, LYSCN\_BOVIN, LYSC3\_PIG, LYSC\_AMYCA, **LYS-BUFAN** (*Bufo andrewsi*), LYSC\_PHOVI, LYSC\_LEPWE, **LYSC\_CAMDR** (*Camelus dromedarius*), **LYSC\_CHEMY** (*Chelonia mydas*);

3. Myoglobin (UniProtKB)

MYG\_CARCR, MYG\_CHEMY, MYG\_GRAGE, MYG\_VARVA, MYG\_URILO, MYG\_ANAPO, MYG\_PHAFI, MYG\_CHICK, MYG\_CERMN, MYG\_AETPY, **MYG\_STRCA** (*Struthio camelus*), **MYG\_ALLMI** (*Alligator mississippiensis*), MYG\_PERPO, MYG\_OCHPR, MYG\_LEPMU, MYG\_LUTLU, MYG\_NYCCO, MYG\_TUPGL, MYG\_SAISC, MYG\_CEBAP, MYG\_PROGU, MYG\_OCHCU, MYG\_PIG, MYG\_CALJA, MYG\_AOTTR, MYG\_ROUAE, MYG\_GORBE, MYG\_RABIT, MYG\_CTEGU, MYG\_LAGLA, MYG\_LAGMA, MYG\_MACFA, MYG\_ERYPA, MYG\_GALCR, MYG\_DIDMA, MYG\_HYLAG, MYG\_HUMAN, MYG\_ZALCA, MYG\_PONPY, MYG\_PHOVI, MYG\_PANTR, **MYG\_EQUBU** (*Equus burchellii*), MYG\_ORYAF, MYG\_VULCH, MYG\_ORNAN, MYG\_CANFA, MYG\_LYCPI, **MYG\_RAT** (*Rattus norvegicus*), MYG\_SPAEH, MYG\_ERIEU, MYG\_TACAC, MYG\_INIGE, MYG\_MACRU, MYG\_MOUSE, MYG\_MELME, MYG\_INDPC, MYG\_ZIPCA, MYG\_MEGNO, MYG\_MESST, MYG\_BALBO, MYG2\_STEAT, MYG\_BALAC, MYG\_GLOME, MYG1\_STEAT, MYG\_BALPH, MYG\_ORCOR, MYG\_PENEL, MYG\_SHEEP, MYG\_PHOSI, MYG\_BOVIN, MYG\_PHODA, MYG\_PHYCA, MYG\_KOGBR, MYG\_CEREL, MYG\_BUBBU, MYG\_LOXAF, MYG\_ELEMA;

4. Parvalbumin (UniProtKB)

**PRVA\_RAT** (*Rattus norvegicus*), PRVA\_MOUSE, PRVA\_GERSP, PRVA\_MACFU, PRVA\_HUMAN, PRVA\_BOVIN, PRVA\_RABIT, PRVA\_FELCA, PRVM\_CHICK, PRVA\_RANCA, PRVA\_RANES, **PRVA\_AMPME** (*Amphiuma means*), **PRVA\_LATCH** (*Latimeria chalumnae*), **PRVA\_TRISE** (*Triakis semifasciata*);

## **Trees**

### **Two taxa**

#### **HPPK**

LG: (Salmonella:0.72945,Vibrio:0.0)

WAG: (Salmonella:0.66055,Vibrio:0.0);

Poisson: (Salmonella:0.60856,Vibrio:0.0);

#### **Lysozyme**

LG: (LYSC\_HUMAN:0.54236,LYSC\_CHEMY:0.0);

WAG: (LYSC\_HUMAN:0.51281,LYSC\_CHEMY:0.0);

Poisson: (LYSC\_HUMAN:0.48017,LYSC\_CHEMY:0.0);

#### **Myoglobin**

LG: (MYG\_EQUBU:0.49951,MYG\_ALLMI:0.0);

WAG: (MYG\_EQUBU:0.46228,MYG\_ALLMI:0.0);

Poisson: (MYG\_EQUBU:0.46225,MYG\_ALLMI:0.0);

#### **Parvalbumin**

LG: (PRVA\_RAT:0.58710,PRVA\_TRISE:0.0);

WAG: (PRVA\_RAT:0.54033,PRVA\_TRISE:0.0);

Poisson: (PRVA\_RAT:0.54166,PRVA\_TRISE:0.0);

## Three taxa

### HPPK

LG: (Aeromonas:0.367204,Salmonella:0.514048,Vibrio:0.228230);

WAG: (Aeromonas:0.337242,Salmonella:0.464291,Vibrio:0.218742);

Poisson: (Aeromonas:0.267048,Salmonella:0.412442,Vibrio:0.224312);

### Lysozyme

LG: (LYSC\_HUMAN:0.357404,LYS\_BUFAN:0.340339,LYSC\_CHEMY:0.203214);

WAG: (LYSC\_HUMAN:0.345938,LYS\_BUFAN:0.319364,LYSC\_CHEMY:0.188813);

Poisson: (LYSC\_HUMAN:0.313610,LYS\_BUFAN:0.300220,LYSC\_CHEMY:0.181252);

### Myoglobin

LG: (MYG\_EQUBU:0.229045,MYG\_STRCA:0.100012,MYG\_ALLMI:0.264132);

WAG: (MYG\_EQUBU:0.219730,MYG\_STRCA:0.093948,MYG\_ALLMI:0.245053);

Poisson: (MYG\_EQUBU:0.211233,MYG\_STRCA:0.102493,MYG\_ALLMI:0.248147);

### Parvalbumin

LG: (PRVA\_RAT:0.189325,PRVA\_AMPME:0.316503,PRVA\_TRISE:0.395588);

WAG: (PRVA\_RAT:0.187683,PRVA\_AMPME:0.286560,PRVA\_TRISE:0.361415);

Poisson: (PRVA\_RAT:0.229427,PRVA\_AMPME:0.247373,PRVA\_TRISE:0.330310);

## Four taxa

### HPPK

#### 1. LG

((Vibrio:0.327906,Salmonella:0.400892):0.010265,Pectobacterium:0.302588,Aeromonas:0.445498);  
((Vibrio:0.330949,Pectobacterium:0.307360):0.000882,Salmonella:0.403020,Aeromonas:0.448683);  
((Vibrio:0.230261,Aeromonas:0.363793):0.220977,Salmonella:0.338829,Pectobacterium:0.219824);

#### 2. WAG

((Vibrio:0.313270,Salmonella:0.368740):0.005007,Pectobacterium:0.282739,Aeromonas:0.407168);  
((Vibrio:0.314001,Pectobacterium:0.283188):0.008442,Salmonella:0.364117,Aeromonas:0.403835);  
((Vibrio:0.222122,Aeromonas:0.327729):0.206134,Salmonella:0.314046,Pectobacterium:0.204833);

#### 3. Poisson

((Vibrio:0.284759,Salmonella:0.325622):0.024539,Pectobacterium:0.252609,Aeromonas:0.332781);  
((Vibrio:0.293604,Pectobacterium:0.262784):0.003472,Salmonella:0.333768,Aeromonas:0.340022);  
((Vibrio:0.216645,Aeromonas:0.270011):0.189444,Salmonella:0.273903,Pectobacterium:0.187624);

## Lysozyme

#### 1. LG

((LYSC\_HUMAN:0.226501,LYS\_BUFAN:0.421583):0.000004,LYSC\_CAMDR:0.250204,LYSC\_CHEMY:0.328016);  
((LYSC\_HUMAN:0.162835,LYSC\_CAMDR:0.200649):0.219274,LYS\_BUFAN:0.325866,LYSC\_CHEMY:0.213459);  
((LYSC\_HUMAN:0.226500,LYSC\_CHEMY:0.328016):0.000004,LYS\_BUFAN:0.421583,LYSC\_CAMDR:0.250204);

#### 2. WAG

((LYSC\_HUMAN:0.220983,LYS\_BUFAN:0.395049):0.000004,LYSC\_CAMDR:0.238081,LYSC\_CHEMY:0.314577);  
((LYSC\_HUMAN:0.164085,LYSC\_CAMDR:0.195008):0.208546,LYS\_BUFAN:0.303897,LYSC\_CHEMY:0.203699);  
((LYSC\_HUMAN:0.220982,LYSC\_CHEMY:0.314578):0.000004,LYS\_BUFAN:0.395049,LYSC\_CAMDR:0.238080);

#### 3. Poisson

((LYSC\_HUMAN:0.196094,LYS\_BUFAN:0.375805):0.000004,LYSC\_CAMDR:0.228509,LYSC\_CHEMY:0.293868);  
((LYSC\_HUMAN:0.147103,LYSC\_CAMDR:0.188600):0.189729,LYS\_BUFAN:0.289089,LYSC\_CHEMY:0.194504);  
((LYSC\_HUMAN:0.196094,LYSC\_CHEMY:0.293868):0.000004,LYS\_BUFAN:0.375805,LYSC\_CAMDR:0.228508);

## Myoglobin

#### 1. LG

((MYG\_ALLMI:0.346734,MYG\_EQUBU:0.137471):0.007576,MYG\_STRCA:0.201163,MYG\_RAT:0.154511);  
((MYG\_ALLMI:0.261681,MYG\_STRCA:0.102564):0.161860,MYG\_EQUBU:0.097480,MYG\_RAT:0.123639);  
((MYG\_ALLMI:0.336742,MYG\_RAT:0.141416):0.021991,MYG\_EQUBU:0.131997,MYG\_STRCA:0.201699);

2. WAG

((MYG\_ALLMI:0.311230,MYG\_EQUBU:0.142668):0.004823,MYG\_STRCA:0.180253,MYG\_RAT:0.157309);  
((MYG\_ALLMI:0.243448,MYG\_STRCA:0.095948):0.152621,MYG\_EQUBU:0.098243,MYG\_RAT:0.113698);  
((MYG\_ALLMI:0.301248,MYG\_RAT:0.142007):0.021703,MYG\_EQUBU:0.135894,MYG\_STRCA:0.180156);

3. Poisson

((MYG\_ALLMI:0.336832,MYG\_EQUBU:0.124078):0.007050,MYG\_STRCA:0.198378,MYG\_RAT:0.144539);  
((MYG\_ALLMI:0.250671,MYG\_STRCA:0.097718):0.148530,MYG\_EQUBU:0.091079,MYG\_RAT:0.120739);  
((MYG\_ALLMI:0.320826,MYG\_RAT:0.120287):0.033088,MYG\_EQUBU:0.115644,MYG\_STRCA:0.194298);

## Parvalbumin

1. LG

((PRVA\_RAT:0.179552,PRVA\_AMPME:0.326468):0.176410,PRVA\_LATCH:0.292100,PRVA\_TRISE:0.267793);  
((PRVA\_RAT:0.231982,PRVA\_LATCH:0.348190):0.032714,PRVA\_AMPME:0.343972,PRVA\_TRISE:0.317235);  
((PRVA\_RAT:0.244997,PRVA\_TRISE:0.339128):0.000004,PRVA\_LATCH:0.360526,PRVA\_AMPME:0.364043);

2. WAG

((PRVA\_RAT:0.179067,PRVA\_AMPME:0.292863):0.153355,PRVA\_LATCH:0.272352,PRVA\_TRISE:0.250651);  
((PRVA\_RAT:0.222339,PRVA\_LATCH:0.323013):0.027584,PRVA\_AMPME:0.312089,PRVA\_TRISE:0.299048);  
((PRVA\_RAT:0.232508,PRVA\_TRISE:0.315793):0.000004,PRVA\_LATCH:0.332321,PRVA\_AMPME:0.326638);

3. Poisson

((PRVA\_RAT:0.199973,PRVA\_AMPME:0.266704):0.141846,PRVA\_LATCH:0.268246,PRVA\_TRISE:0.230462);  
((PRVA\_RAT:0.231054,PRVA\_LATCH:0.311383):0.038143,PRVA\_AMPME:0.287142,PRVA\_TRISE:0.274608);  
((PRVA\_RAT:0.237759,PRVA\_TRISE:0.284504):0.018877,PRVA\_LATCH:0.313572,PRVA\_AMPME:0.293818);

## Likelihoods (lnL)

### Two taxa

#### HPPK

LG: -735.662167

WAG: -741.261095

Poisson: -791.086604

### Lysozyme

LG: -601.541213

WAG: -594.482931

Poisson: -619.852608

### Myoglobin

LG: -665.949741

WAG: -668.857905

Poisson: -723.725963

### Parvalbumin

LG: -482.837648

WAG: -478.894401

Poisson: -532.923162



### **Three taxa**

#### **HPPK**

LG: -972.537795

WAG: -977.544549

Poisson: -1042.765084

#### **Lysozyme**

LG: -787.024360

WAG: -776.966605

Poisson: -826.511054

#### **Myoglobin**

LG: -794.556943

WAG: -796.066717

Poisson: -872.272315

#### **Parvalbumin**

LG: -629.383650

WAG: -625.072729

Poisson: -697.673948

## Four taxa

### HPPK

Tree #	LG	WAG	Poisson
1	-1197.547	-1201.159	-1285.722
2	-1197.597	-1201.081	-1286.574
3	-1180.046*	-1182.973*	-1262.192*

### Lysozyme

Tree #	LG	WAG	Poisson
1	-955.261	-942.007	-1006.594
2	-930.941*	-919.064*	-986.027*
3	-955.260	-942.007	-1006.594

### Myoglobin

Tree #	LG	WAG	Poisson
1	-937.633	-941.064	-1045.844
2	-911.131*	-913.278*	-1009.918*
3	-936.682	-939.841	-1041.132

### Parvalbumin

Tree #	LG	WAG	Poisson
1	-775.110*	-768.997*	-866.085*
2	-783.536	-776.476	-874.209
3	-784.903	-777.338	-875.069

# Plots

## Two taxa

a)  $Z_s$  (see paper for  $PE_d$  and  $PE_s$ )

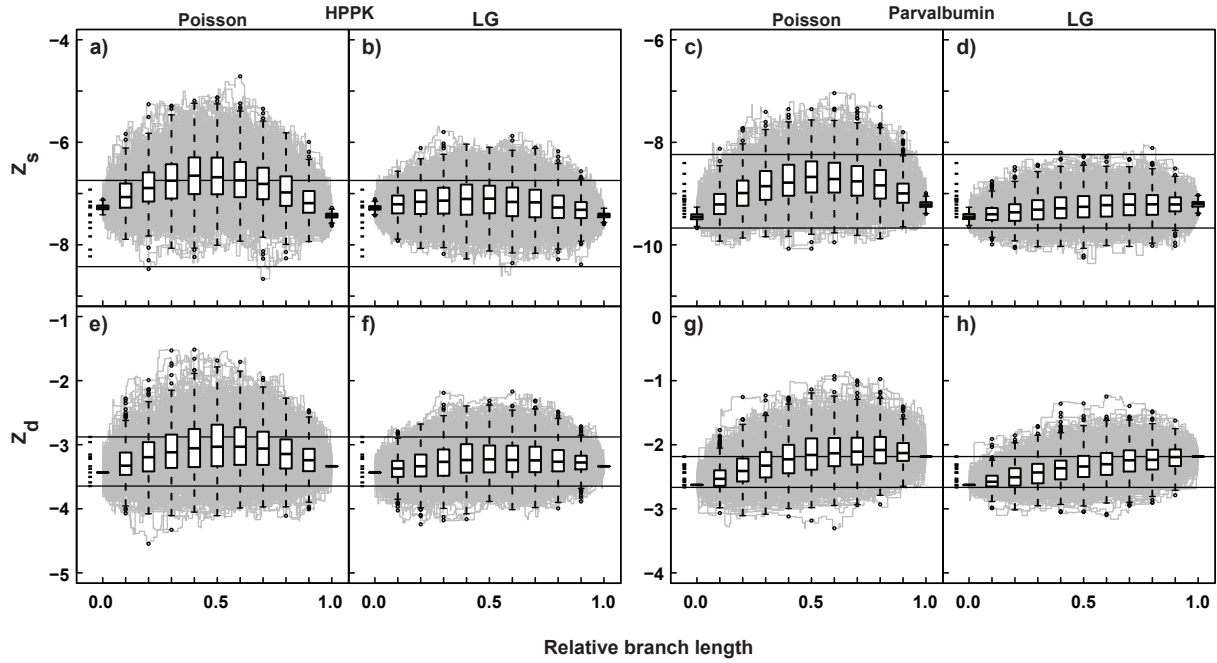


Figure 1. Constrained mappings ( $Z_s$  based on 10,000 shuffles).

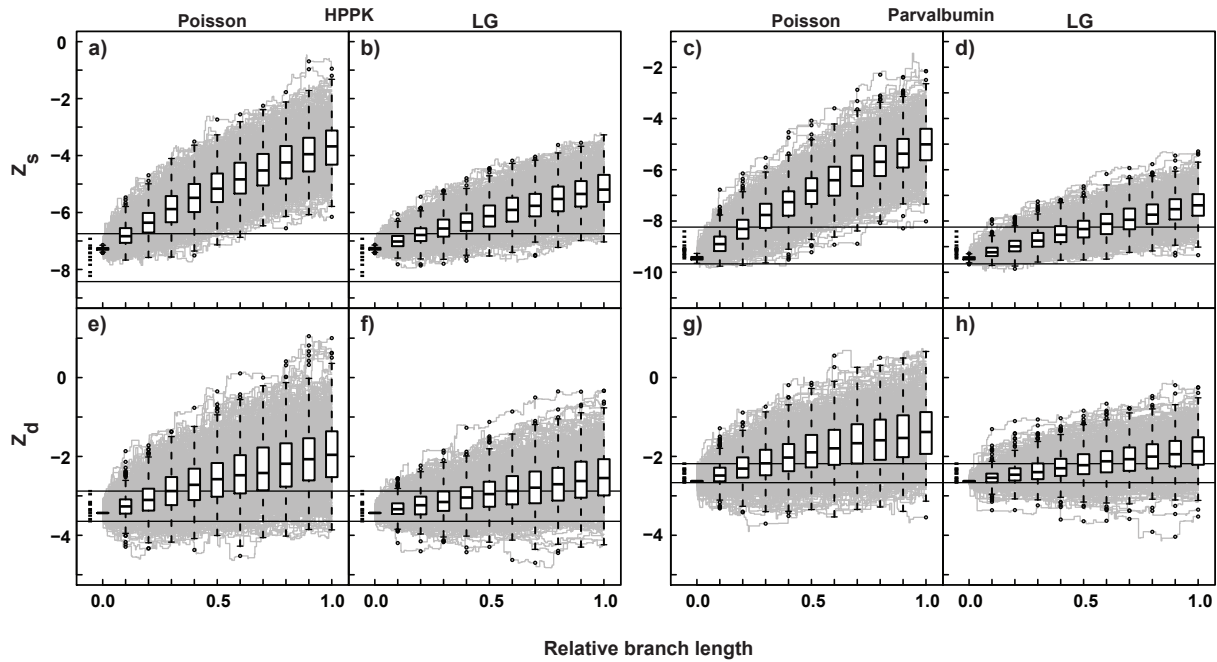


Figure 2. Unconstrained mappings ( $Z_s$  based on 10,000 shuffles).

In the remaining plots the  $PE_s$  and  $Z_s$  scores are based on 1,000 shuffles. Results presented in the paper are based on 10,000 shuffles which greatly reduced the variance. The median  $PW_s$  and  $Z_s$  scores were very similar for 1,000 and 10,000 sequence shuffles.

### **Three taxa**

#### **a) $PE_s$**

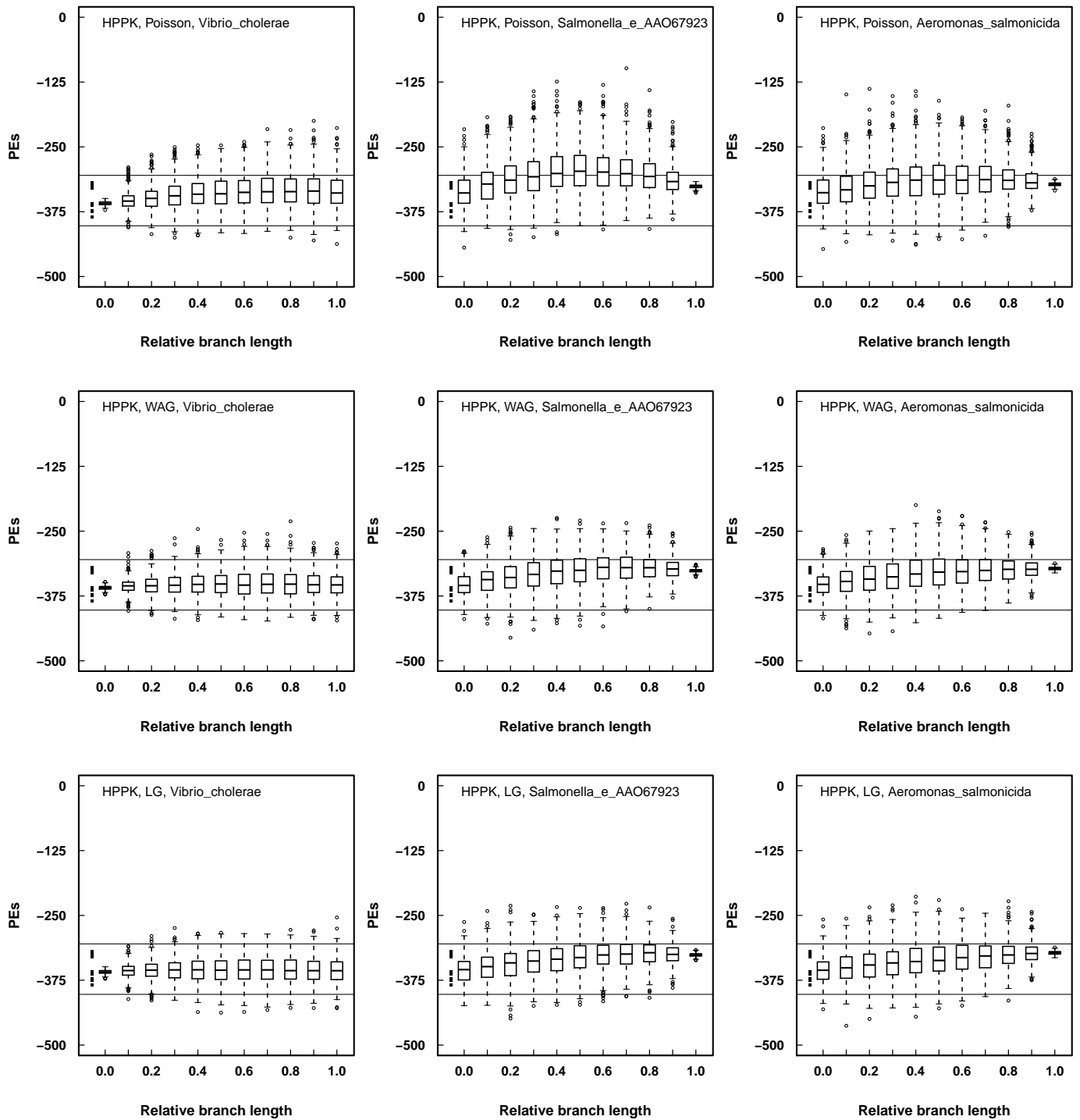


Figure 3. HPPK,  $PE_s$

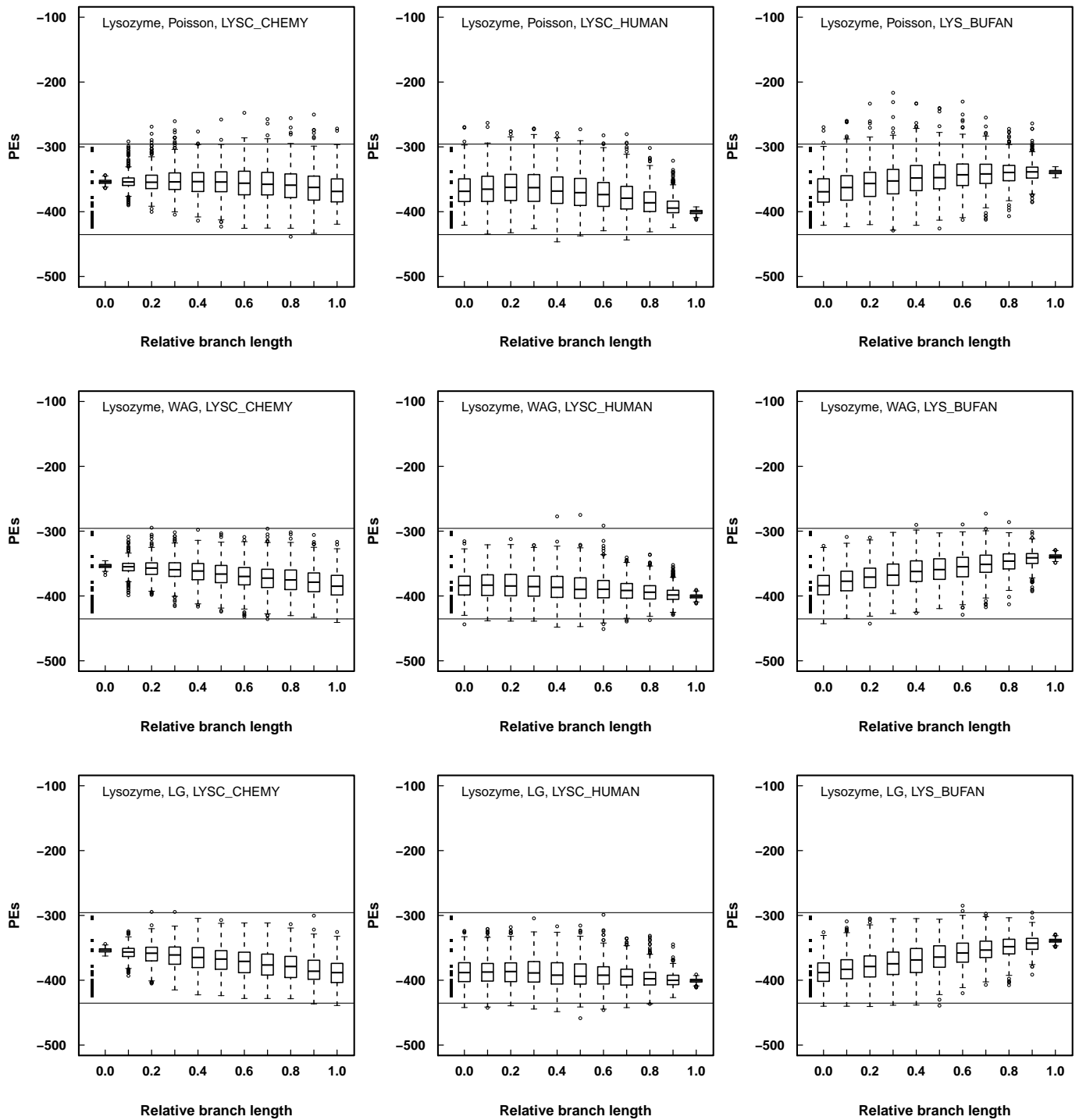


Figure 4. Lysozyme,  $PE_s$

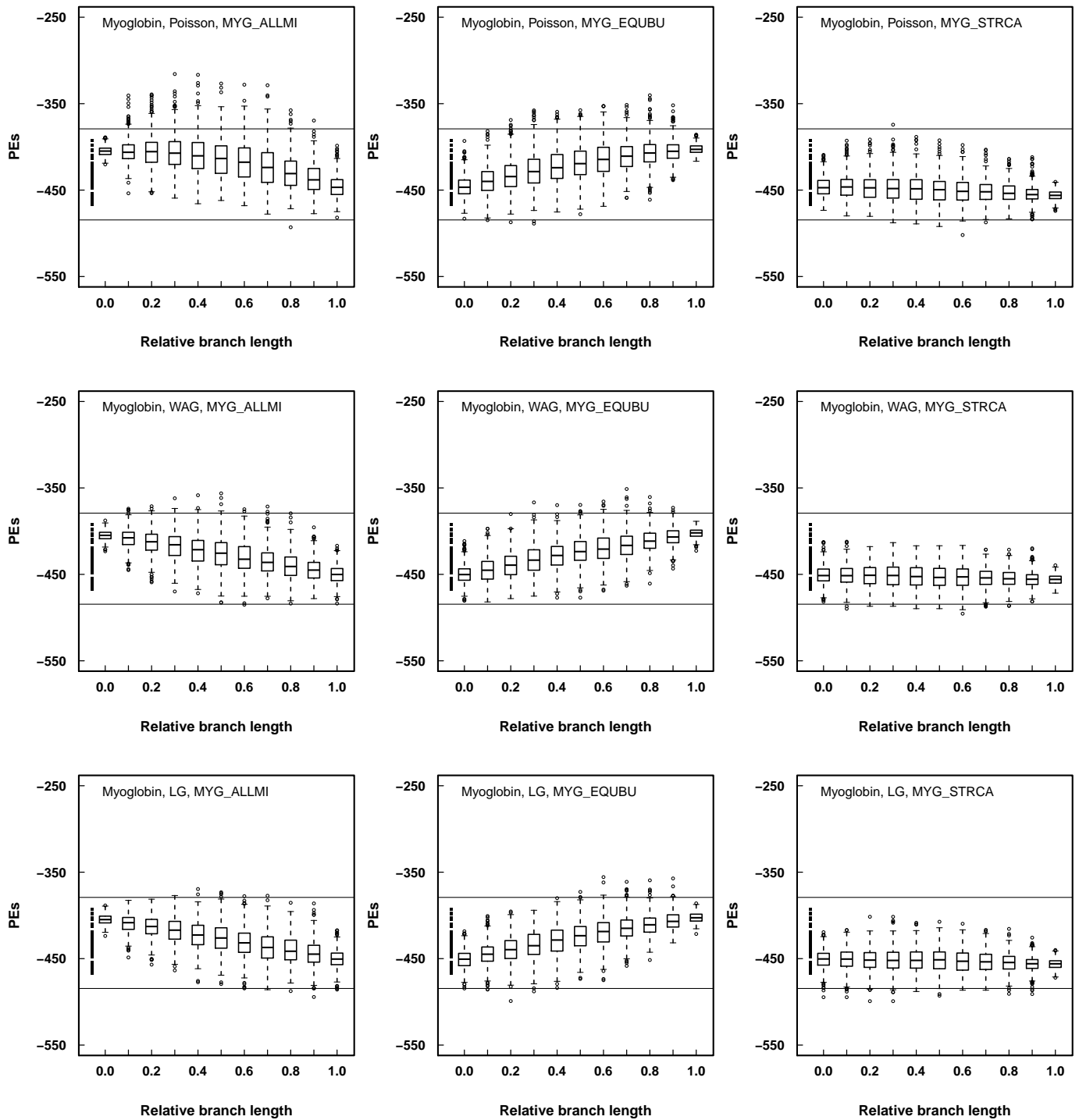


Figure 5. Myoglobin,  $PE_s$

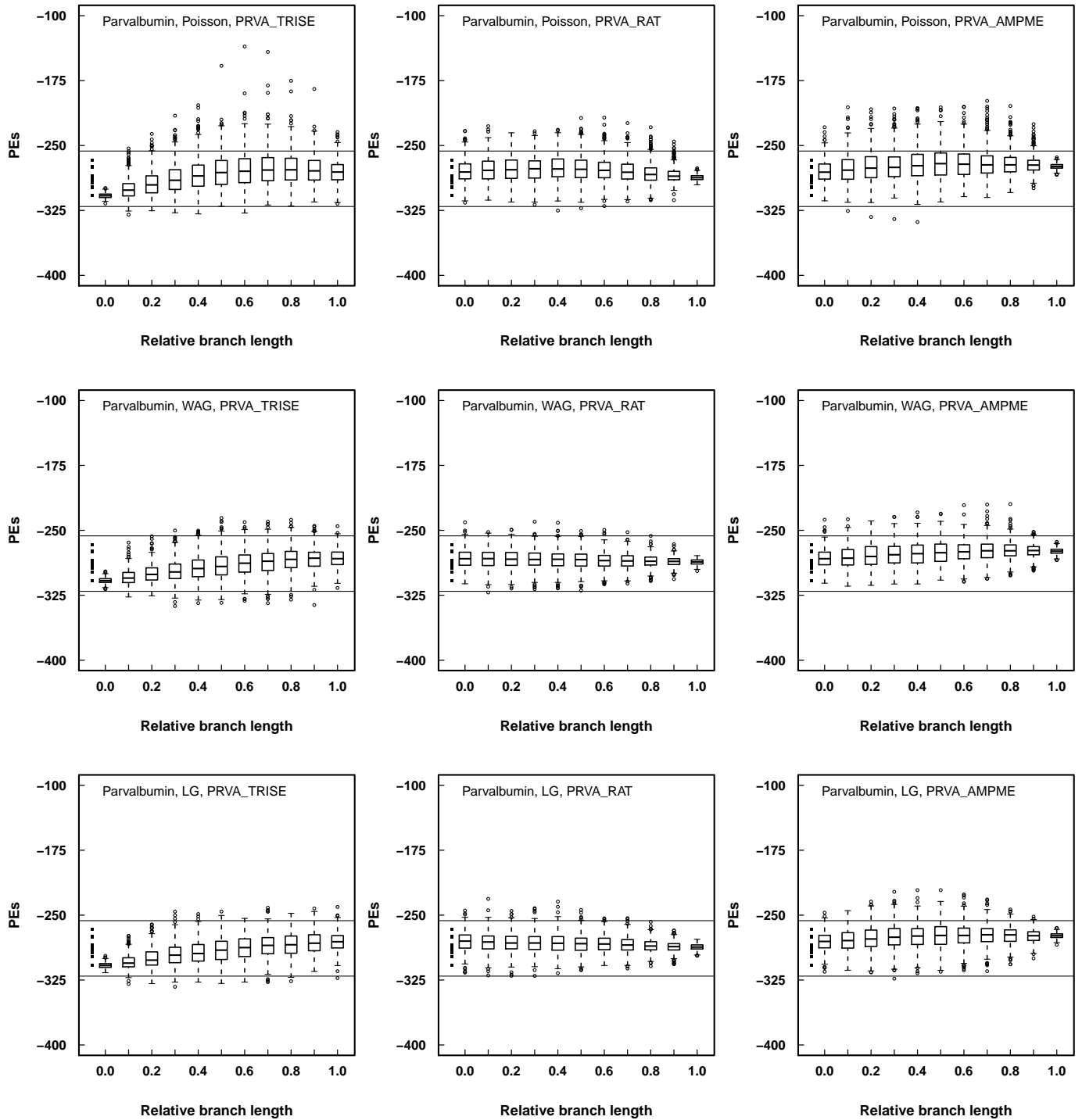


Figure 6. Parvalbumin,  $PE_s$



b)  $Z_s$

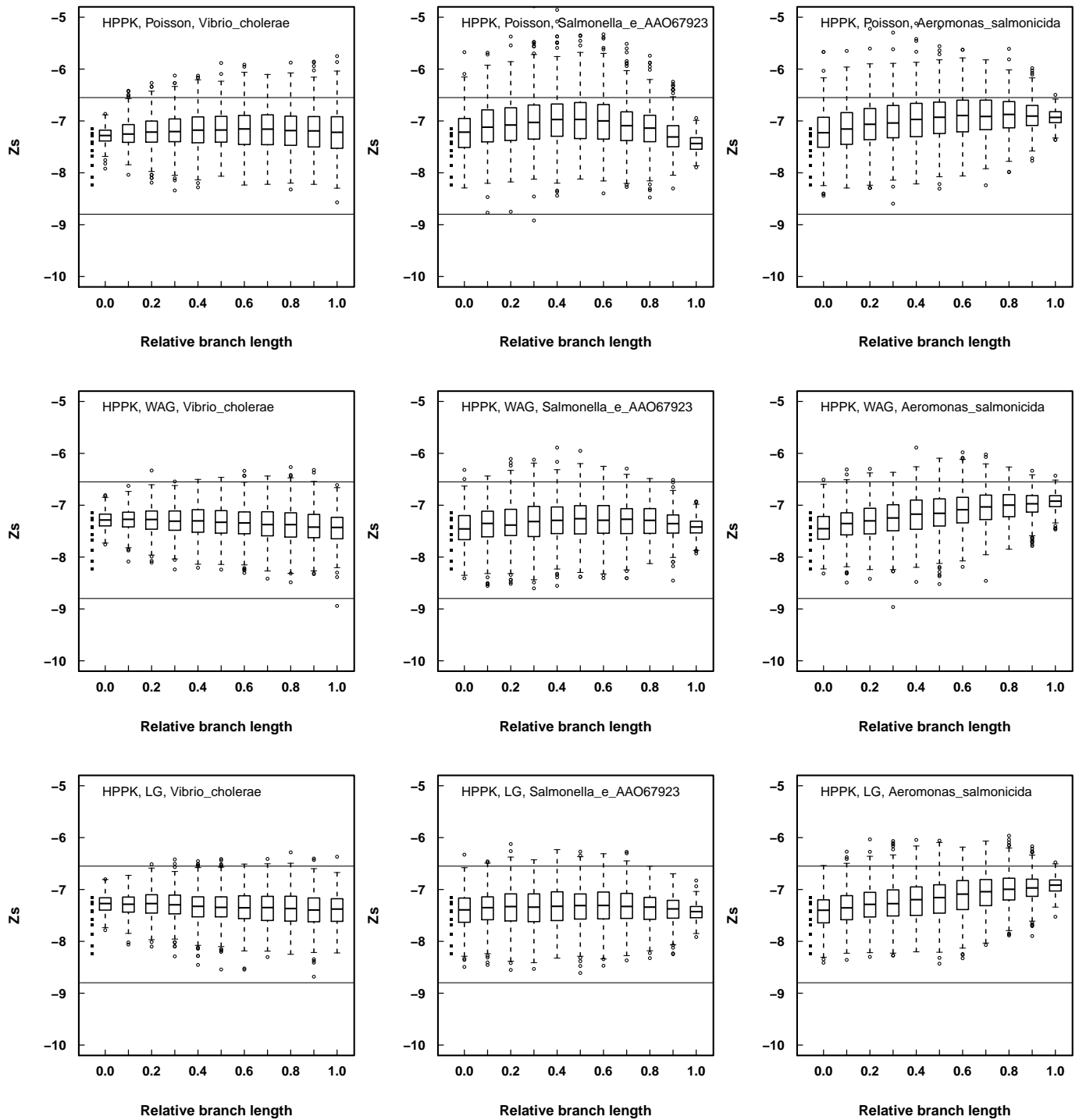


Figure 7. HPPK,  $Z_s$

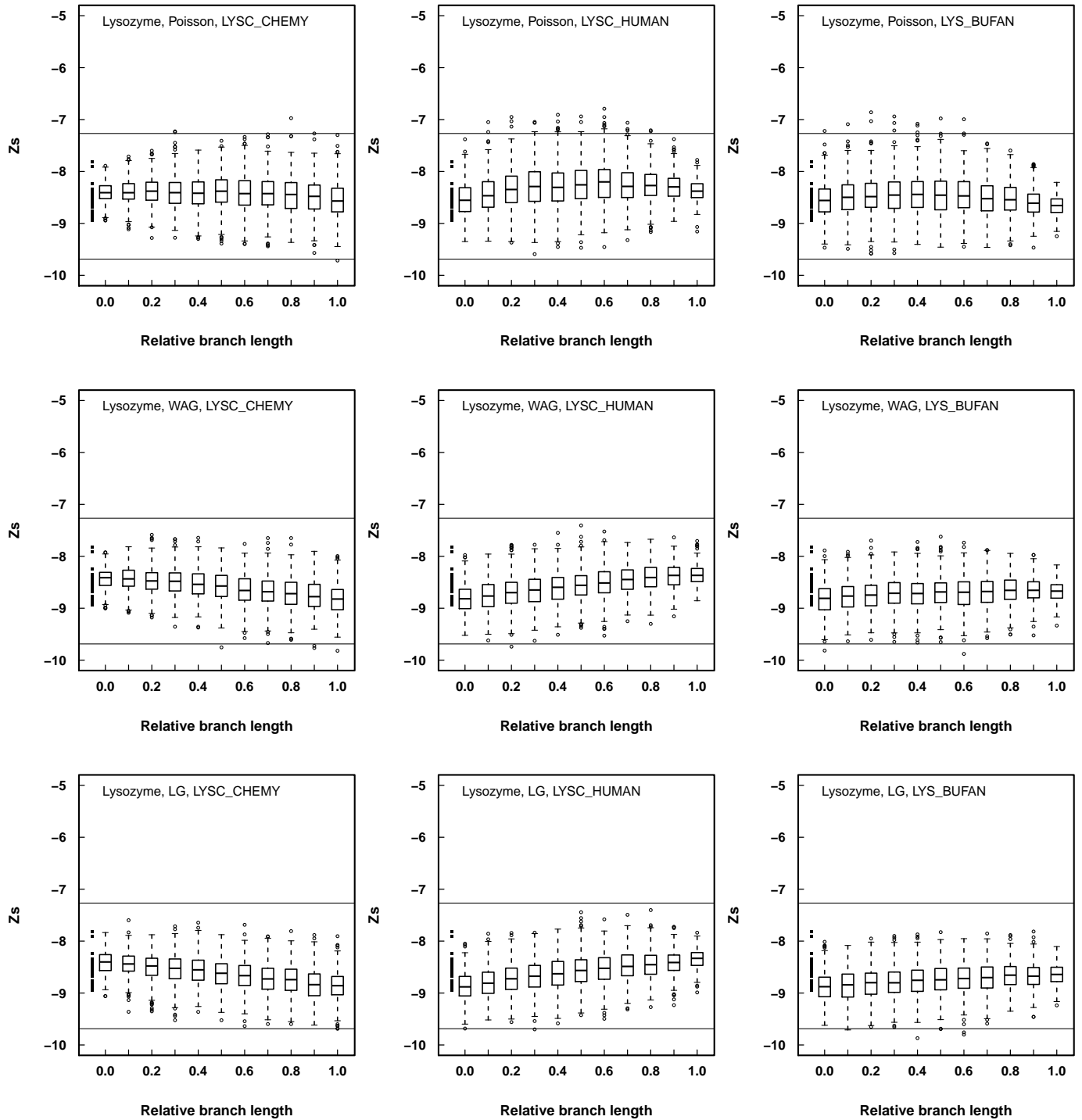


Figure 8. Lysozyme,  $Z_s$

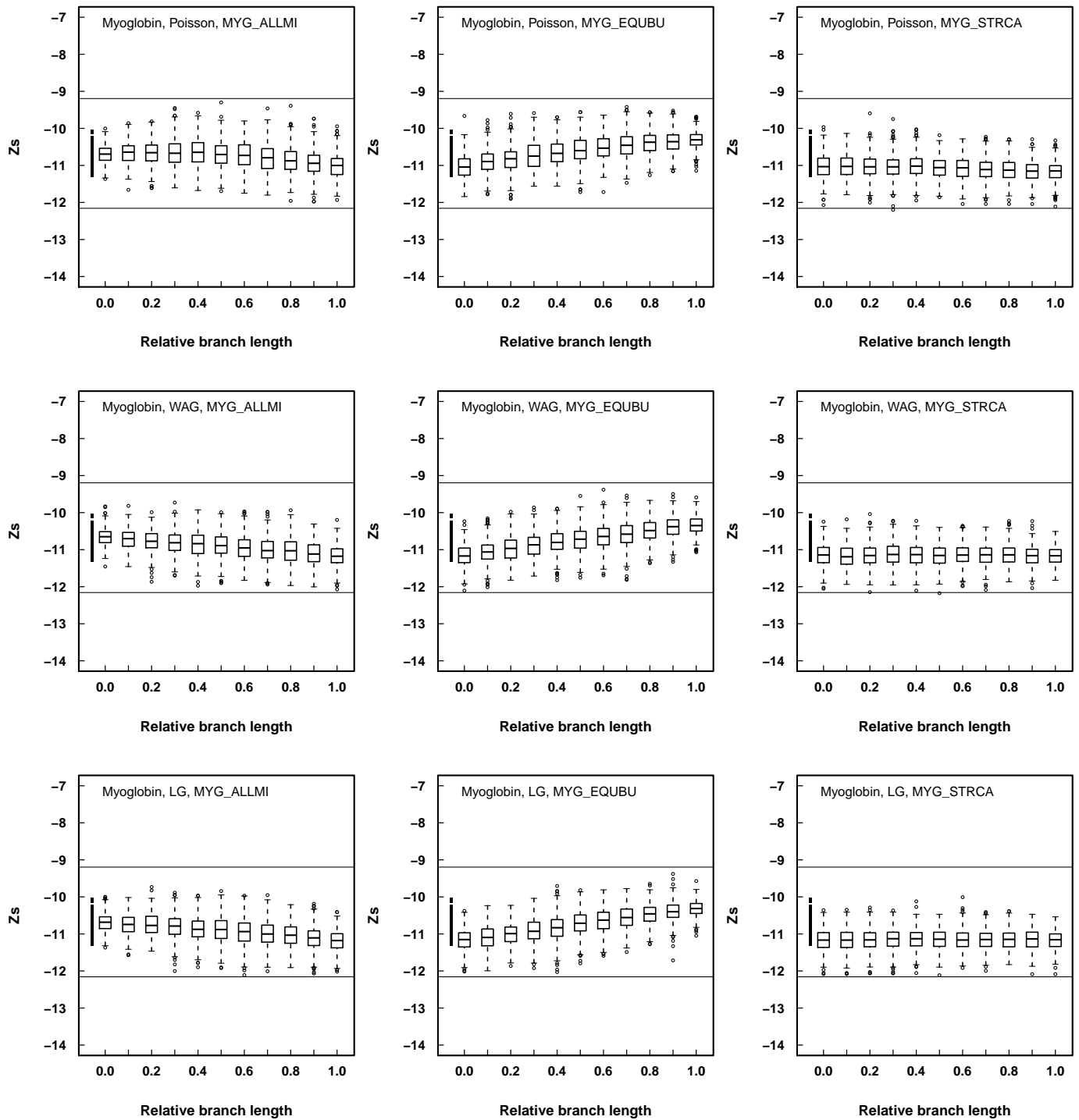


Figure 9. Myoglobin,  $Z_s$

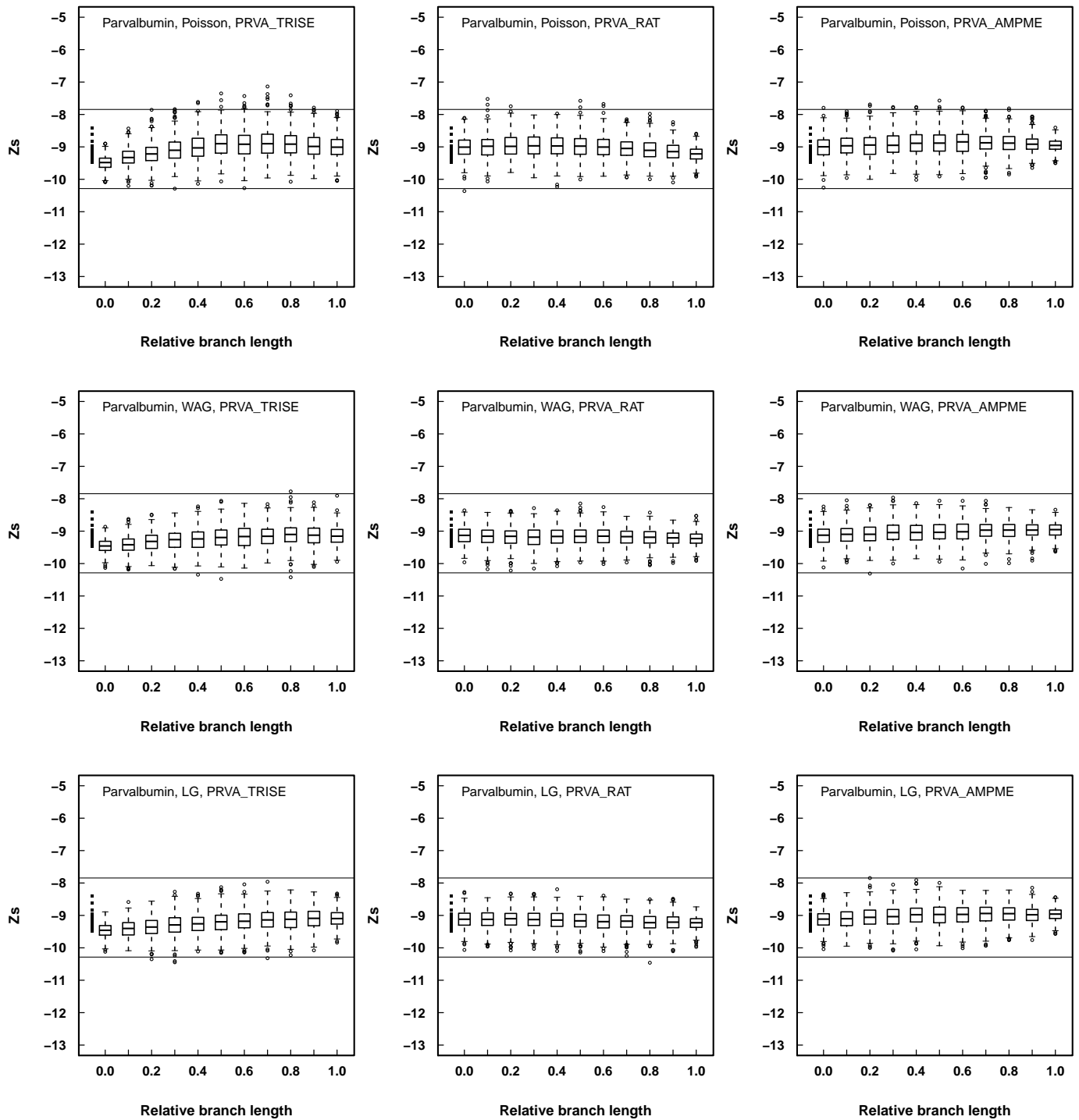


Figure 10. Parvalbumin,  $Z_s$

Four taxa

a)  $PE_s$

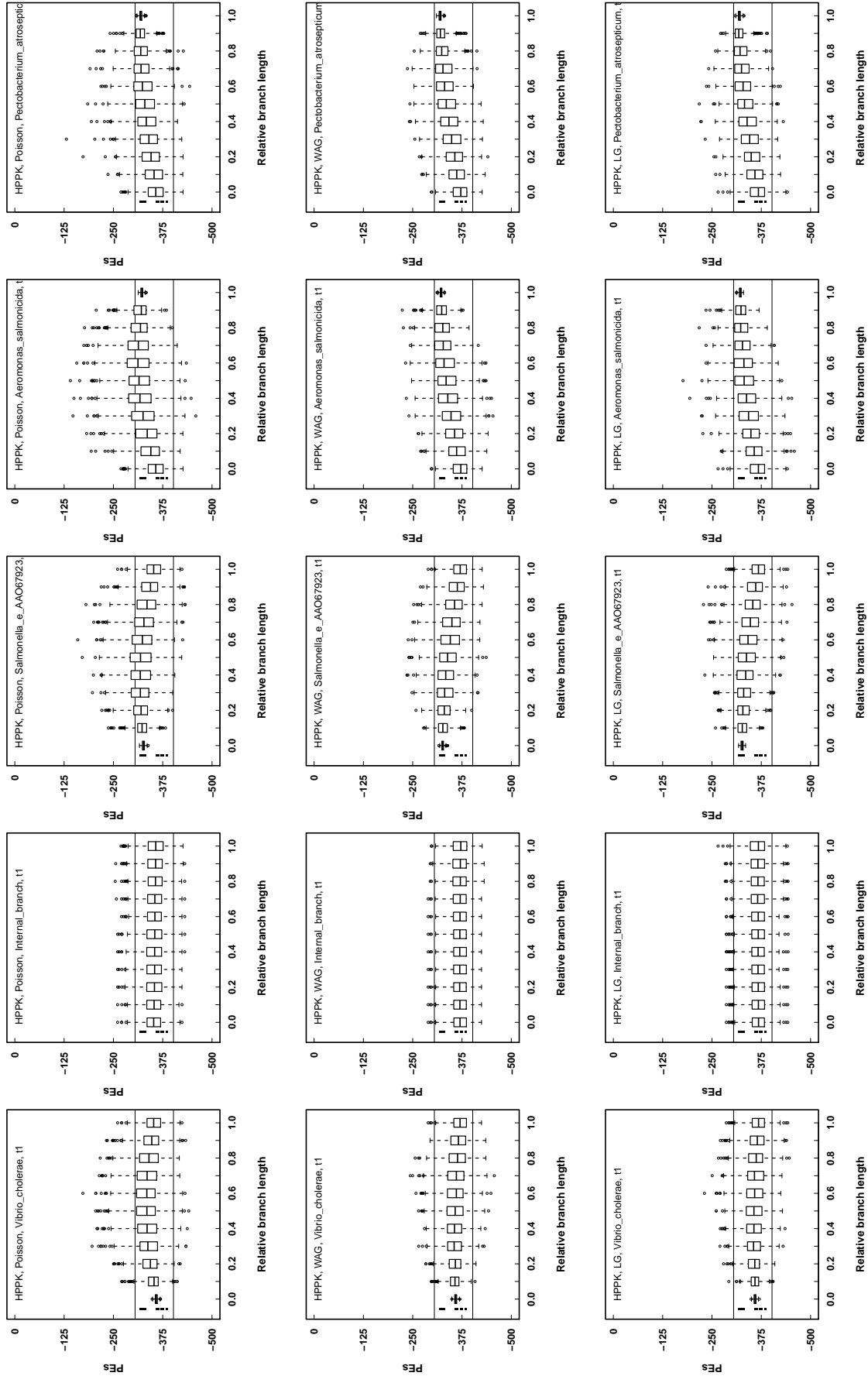


Figure 11. HPPK,  $PE_s$ , Tree 1

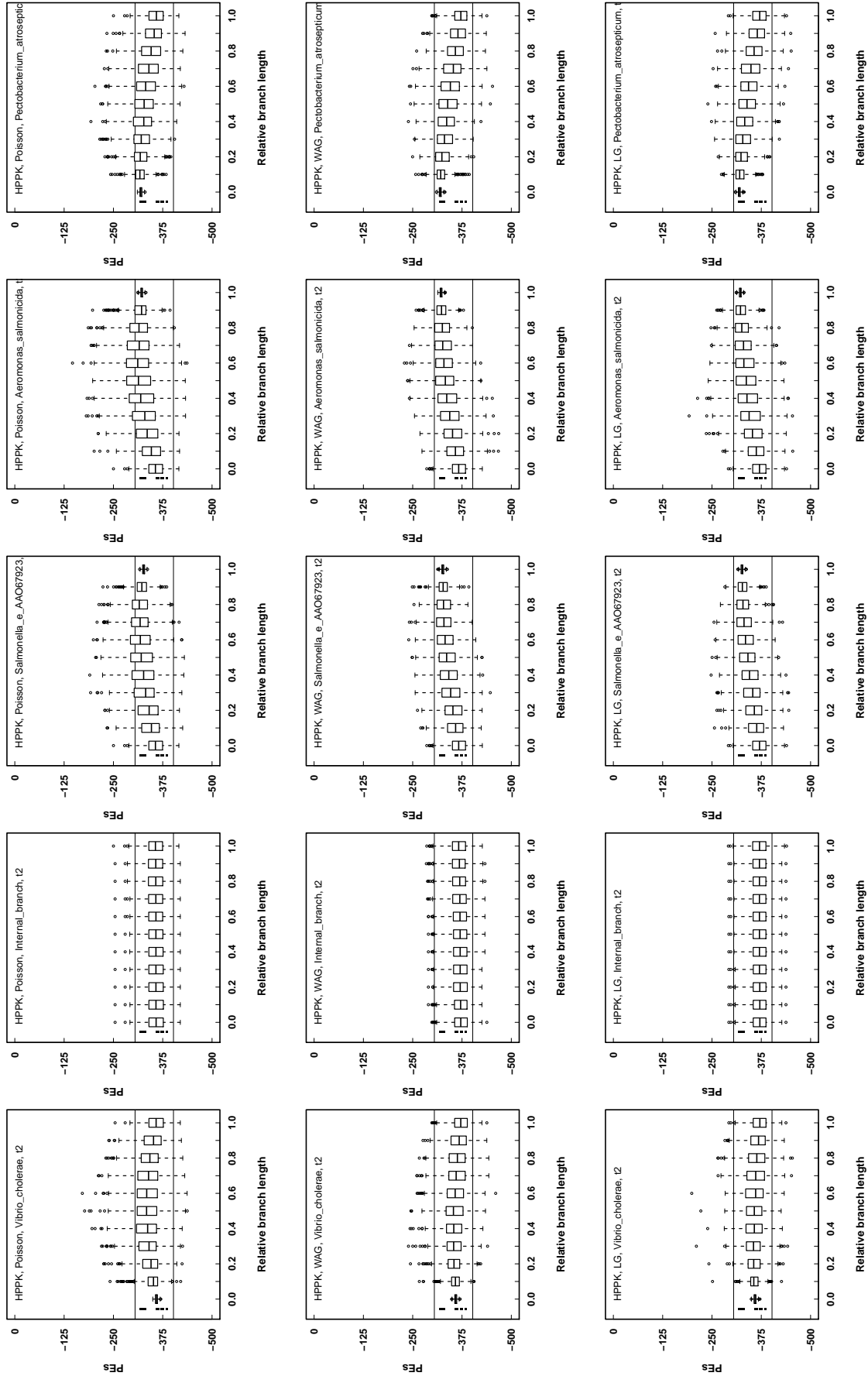


Figure 12. HPPK,  $PE_s$ , Tree 2



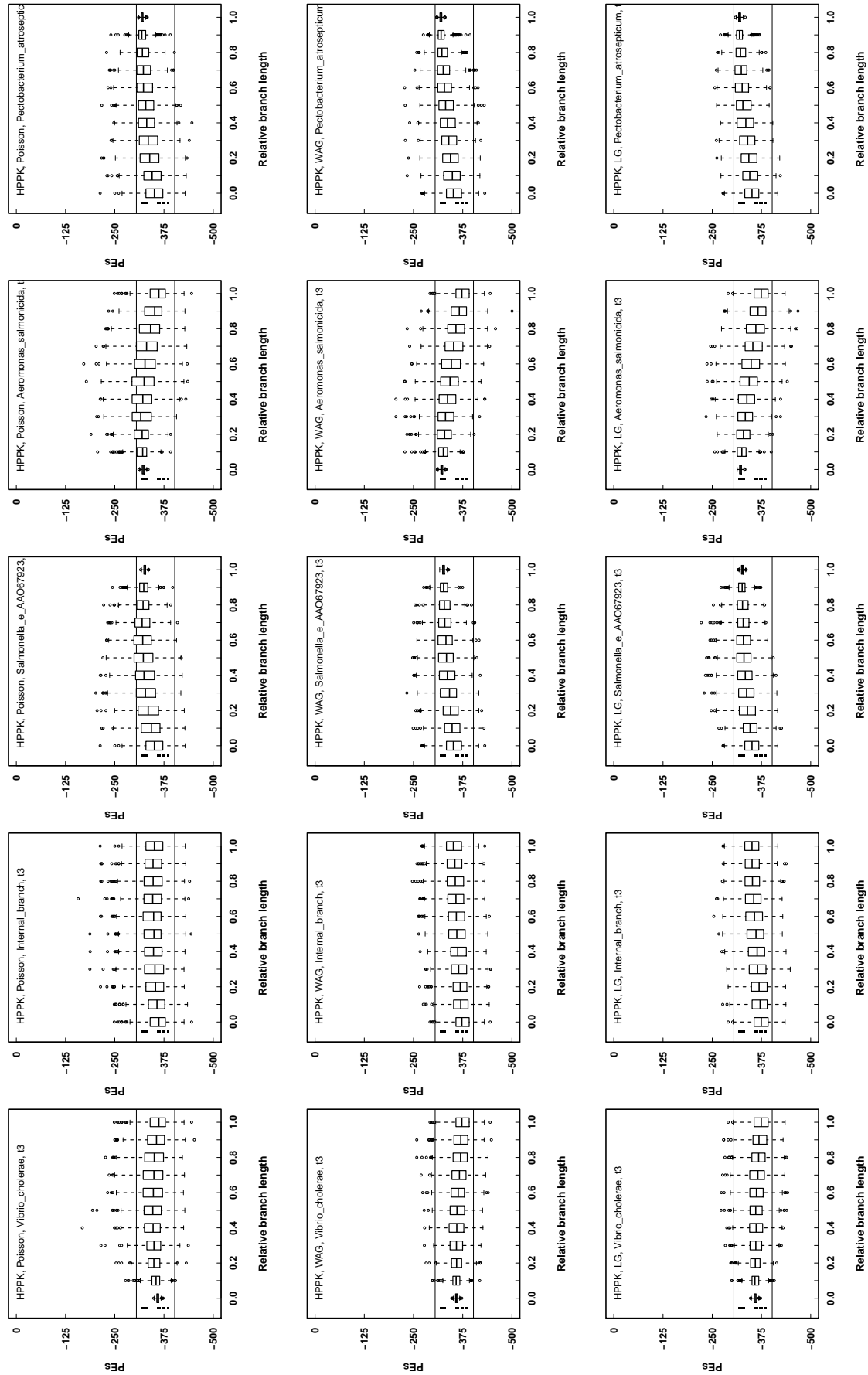


Figure 13. HPPK,  $PE_s$ , Tree 3

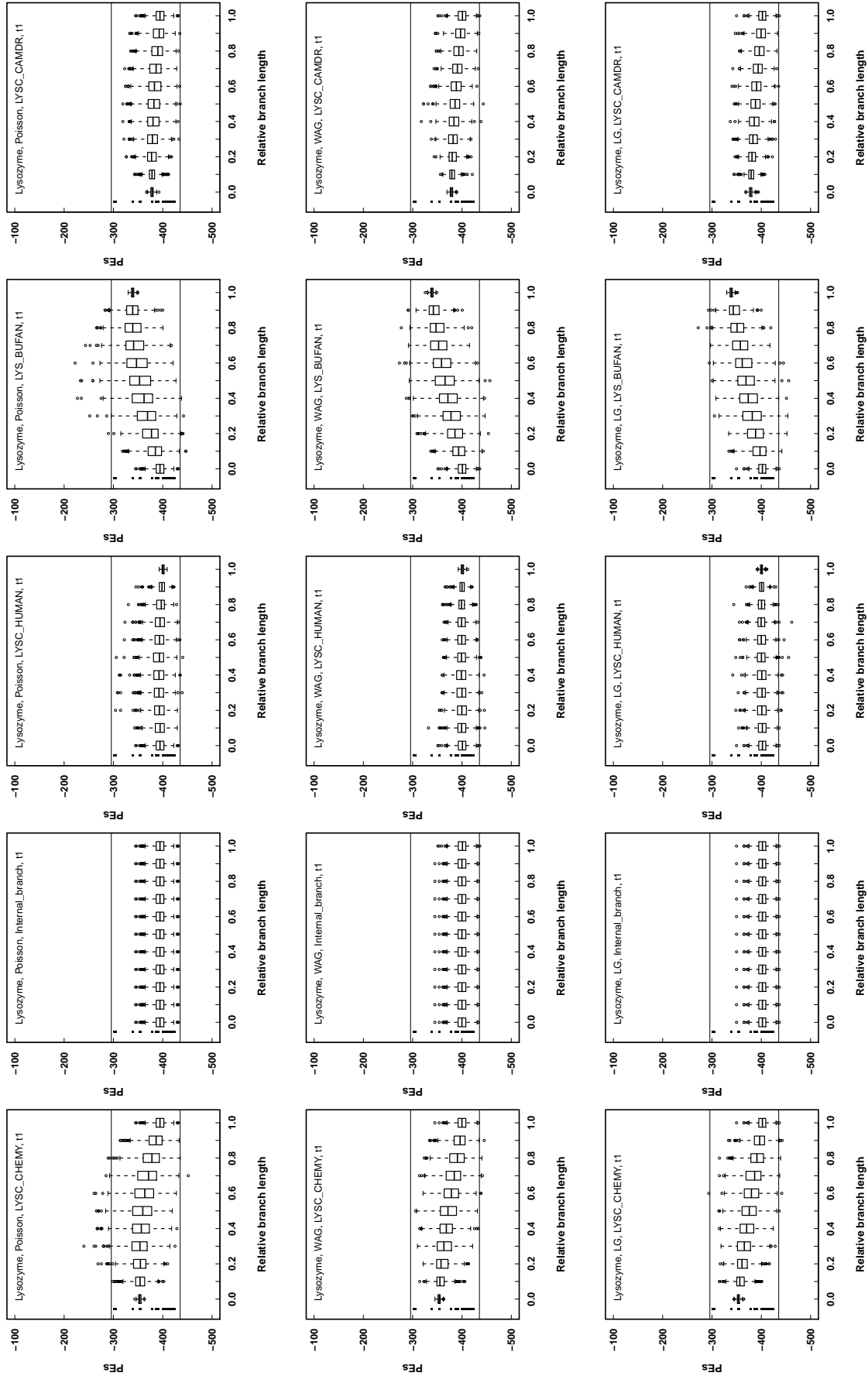


Figure 14. Lysozyme,  $PE_s$ , Tree 1

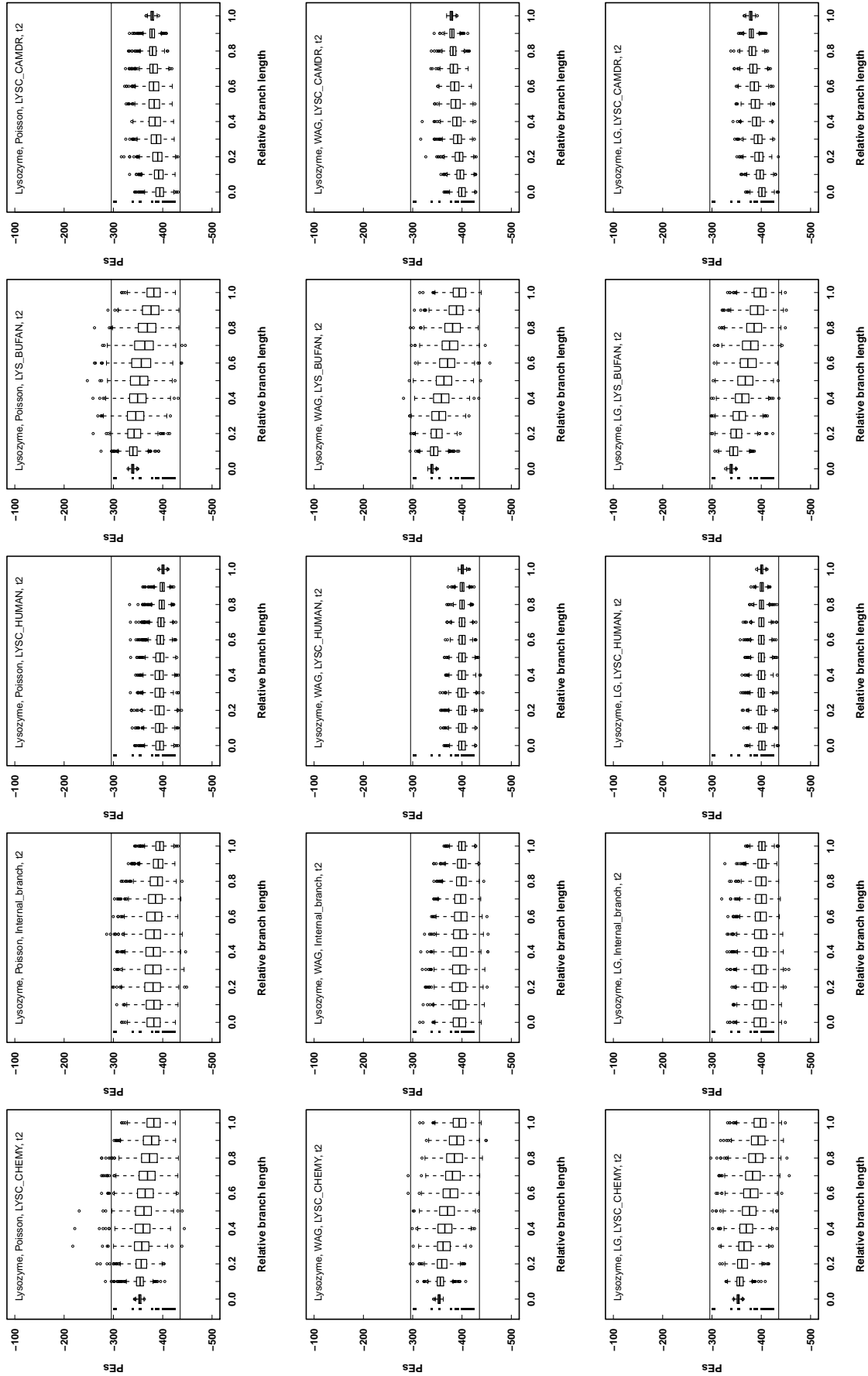


Figure 15. Lysozyme,  $PE_s$ , Tree 2

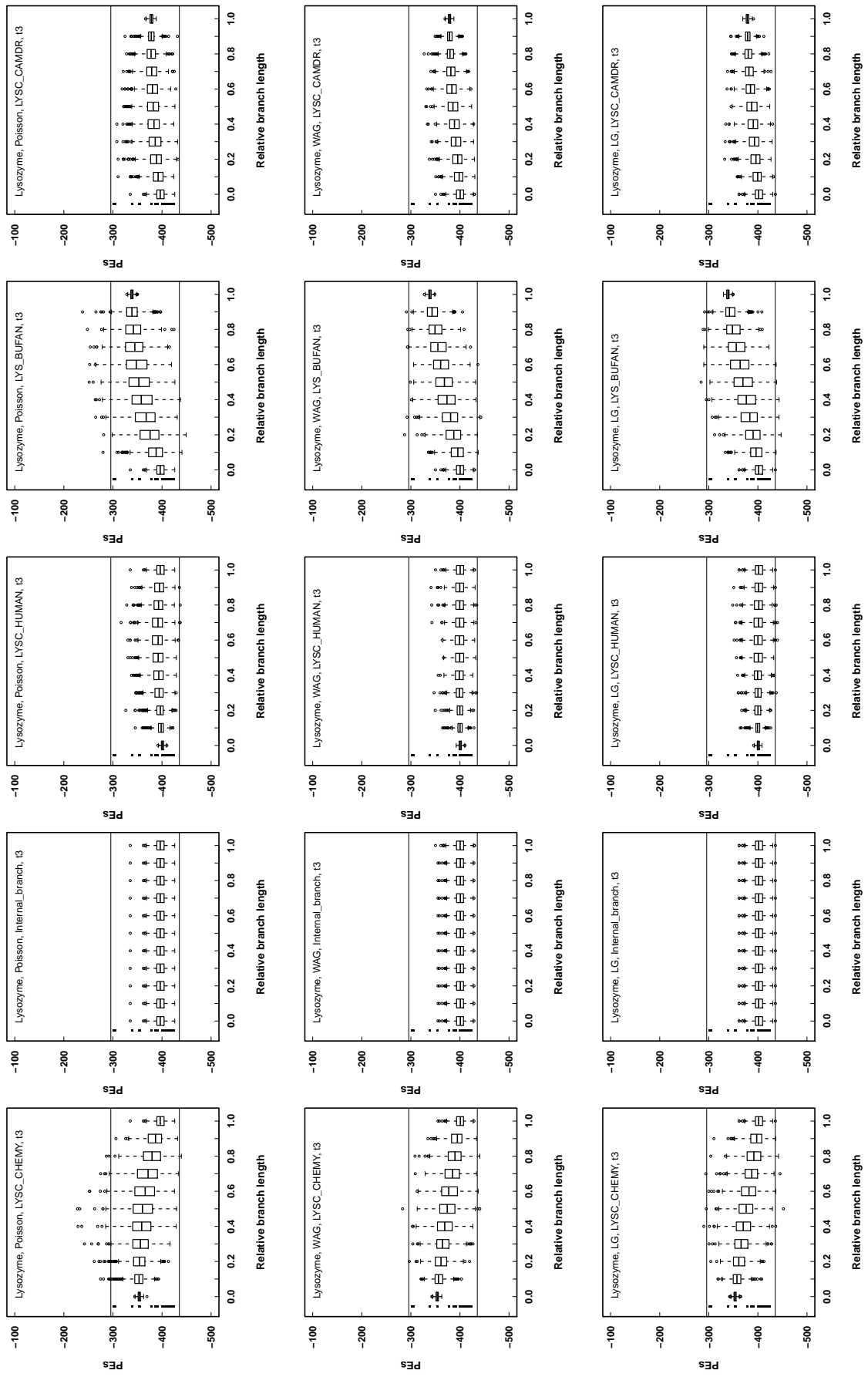


Figure 16. Lysozyme,  $PE_s$ , Tree 3

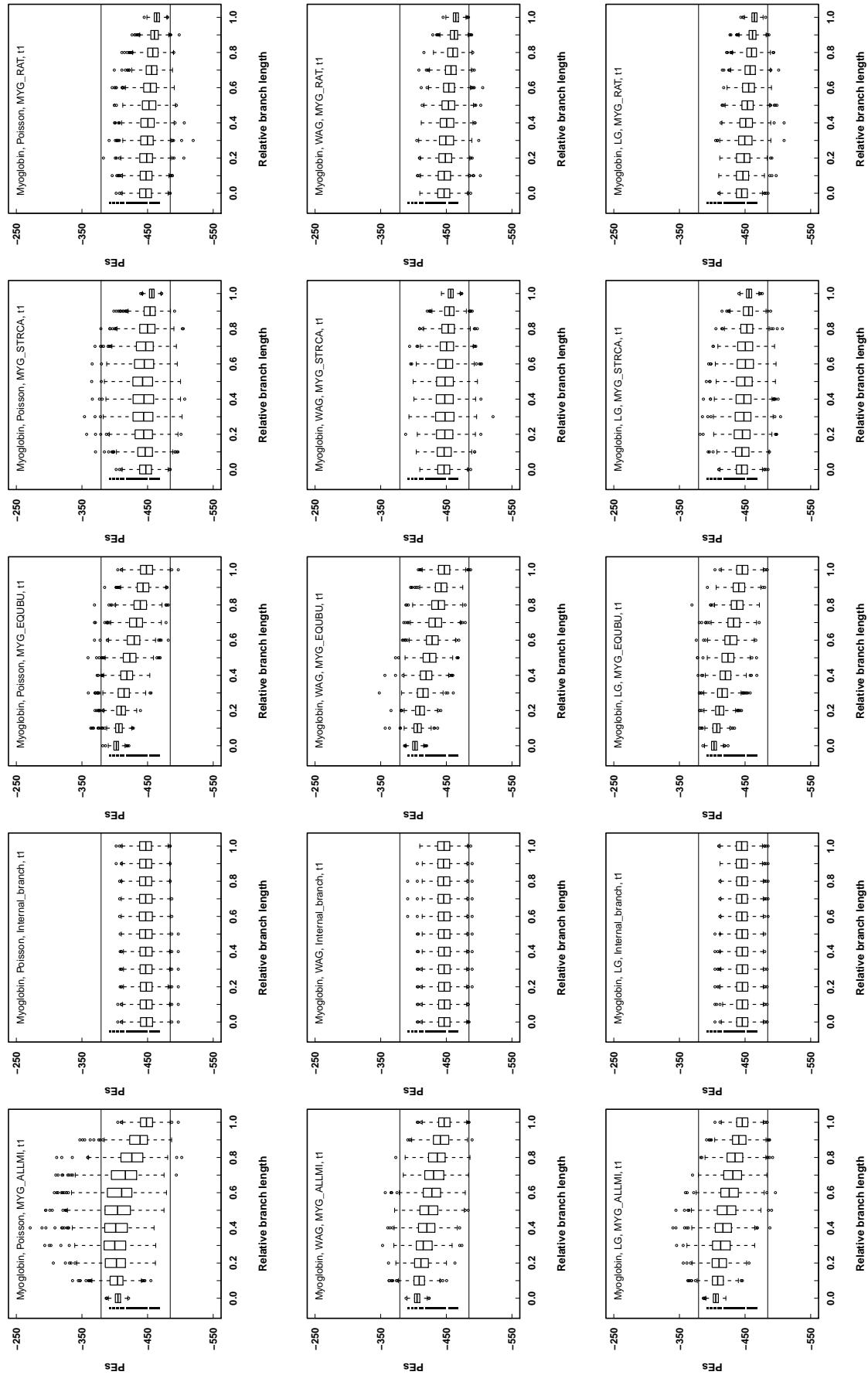


Figure 17. Myoglobin,  $PE_s$ , Tree 1

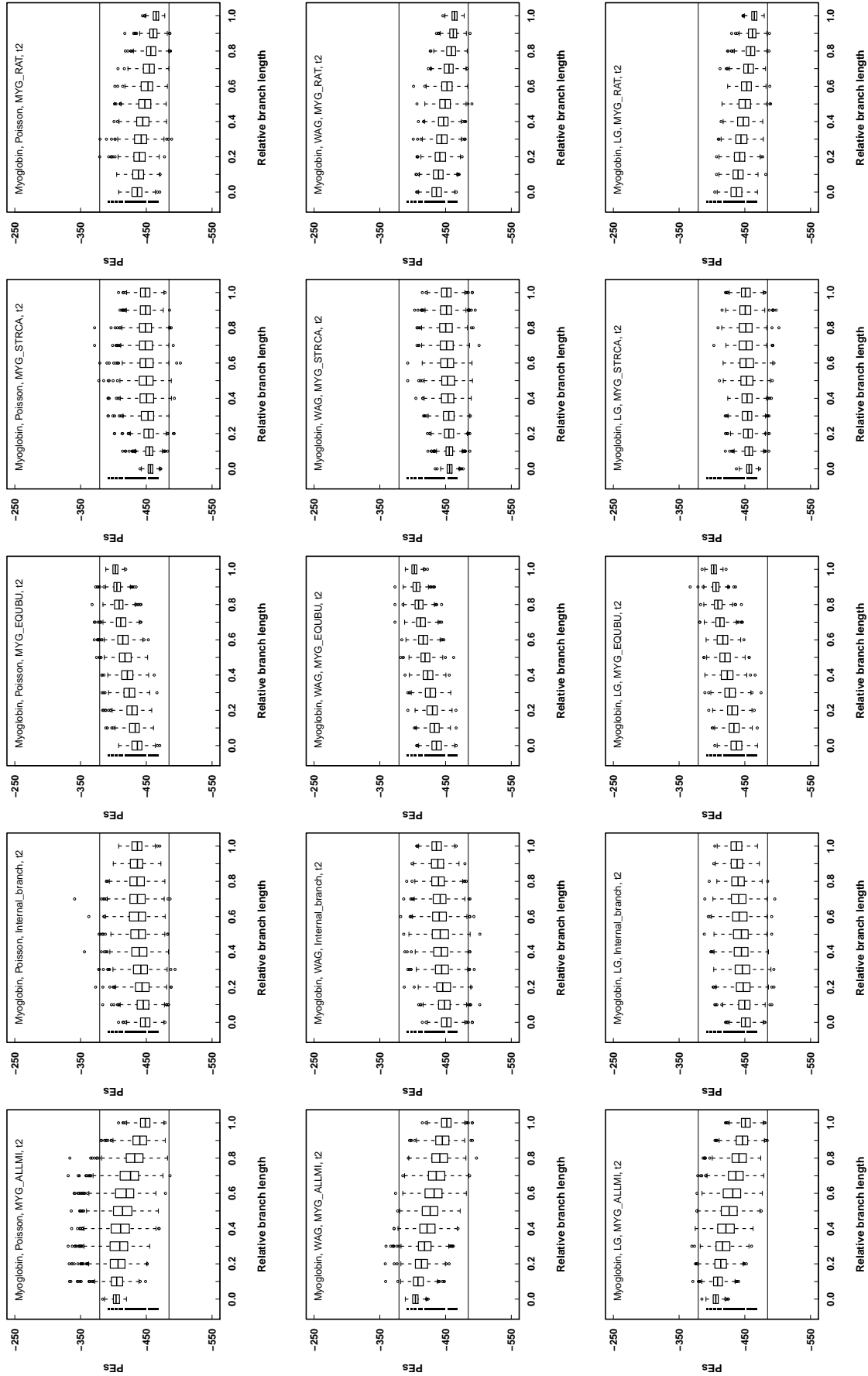


Figure 18. Myoglobin,  $PE_s$ , Tree 2

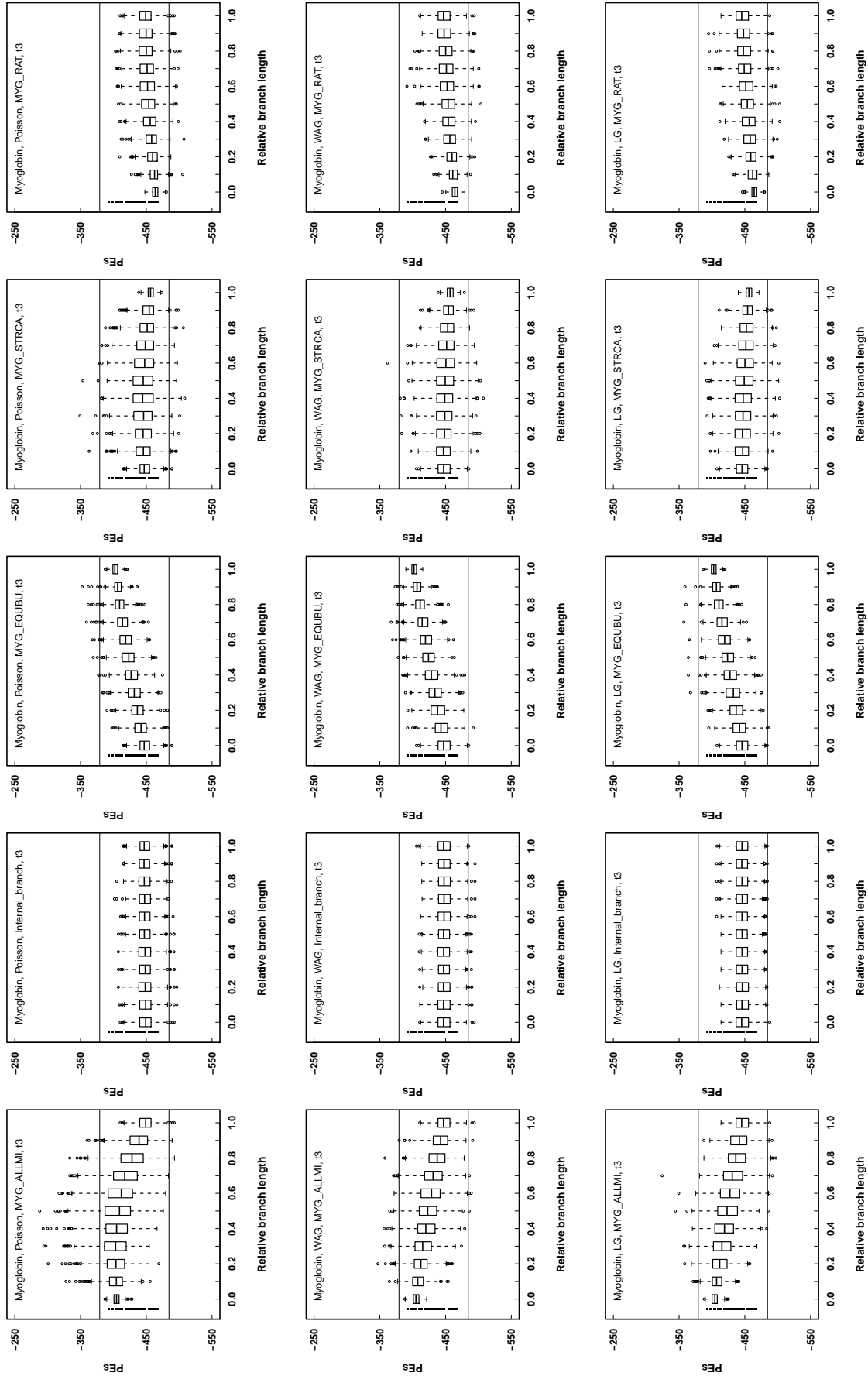


Figure 19. Myoglobin,  $PE_s$ , Tree 3

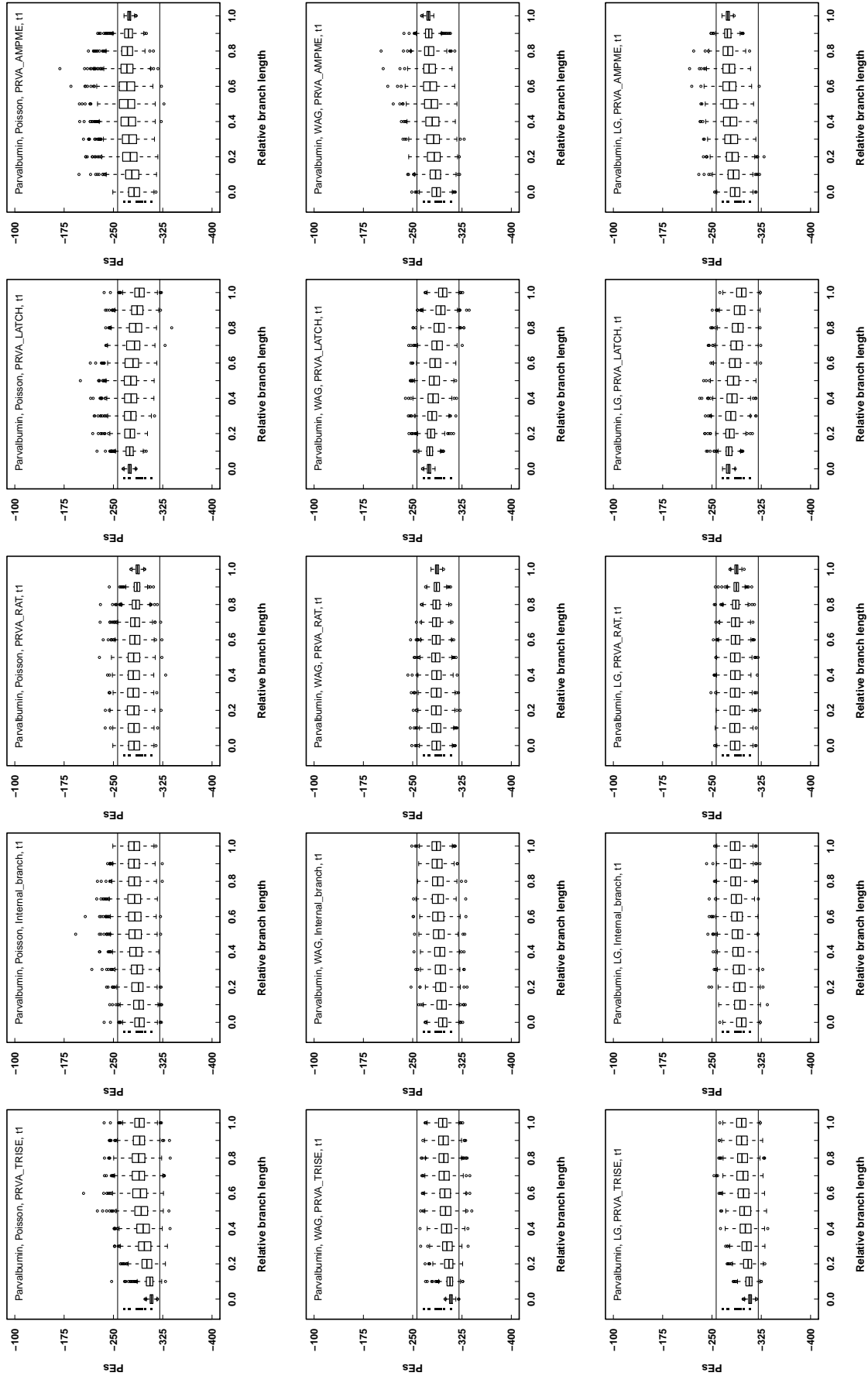


Figure 20. Parvalbumin,  $PE_s$ , Tree 1



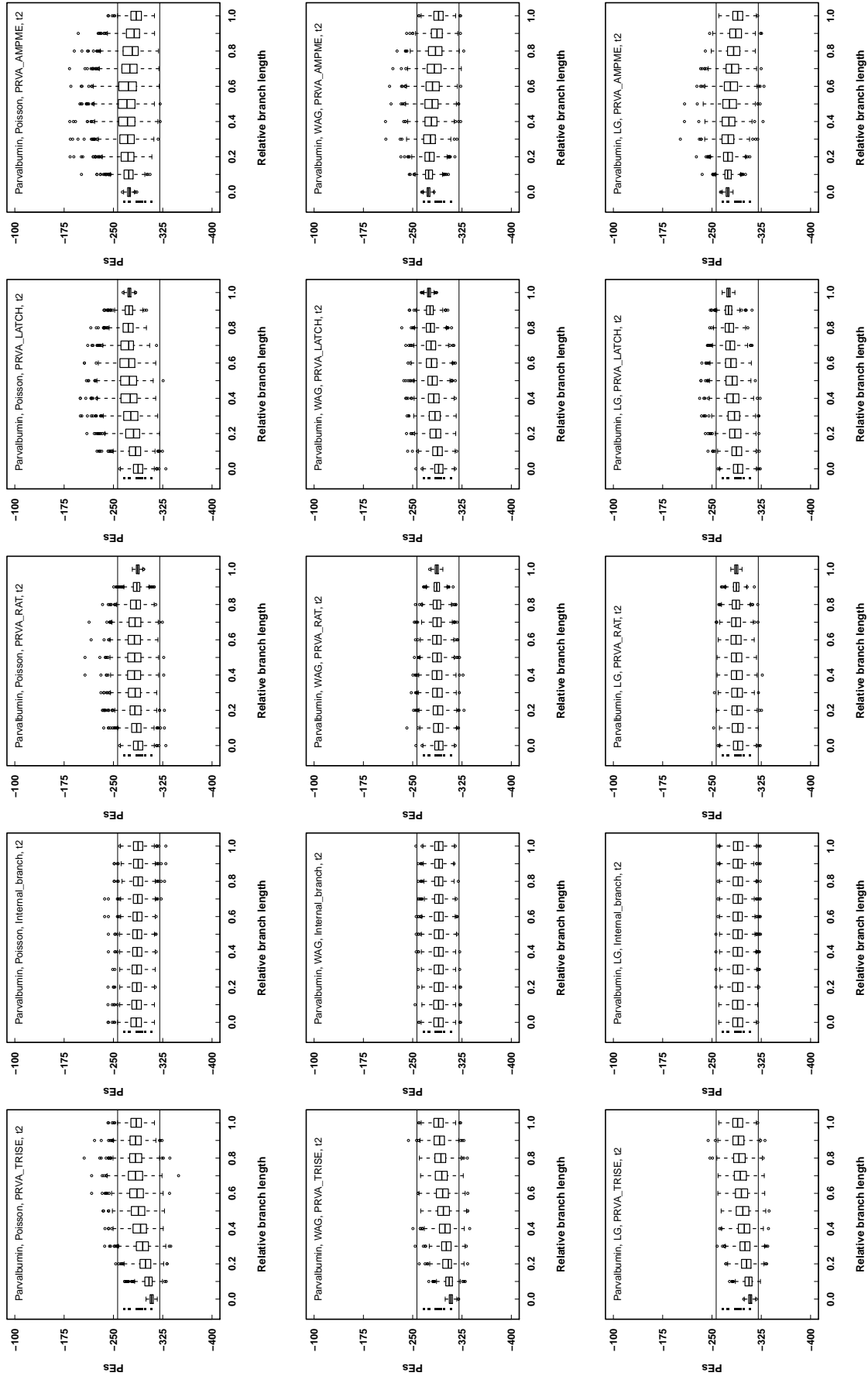


Figure 21. Parvalbumin,  $PE_s$ , Tree 2

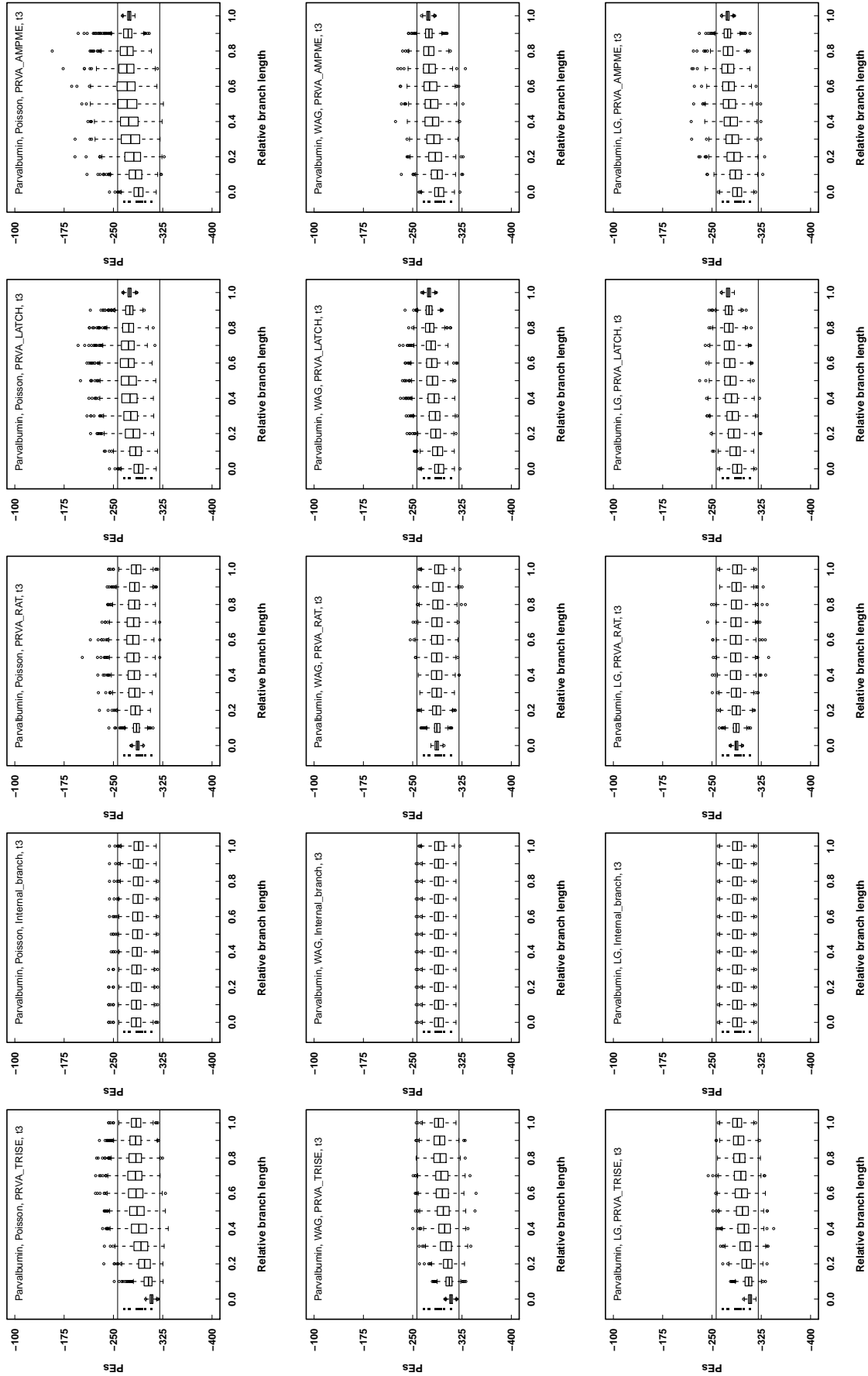


Figure 22. Parvalbumin,  $PE_s$ , Tree 3

b)  $Z_s$

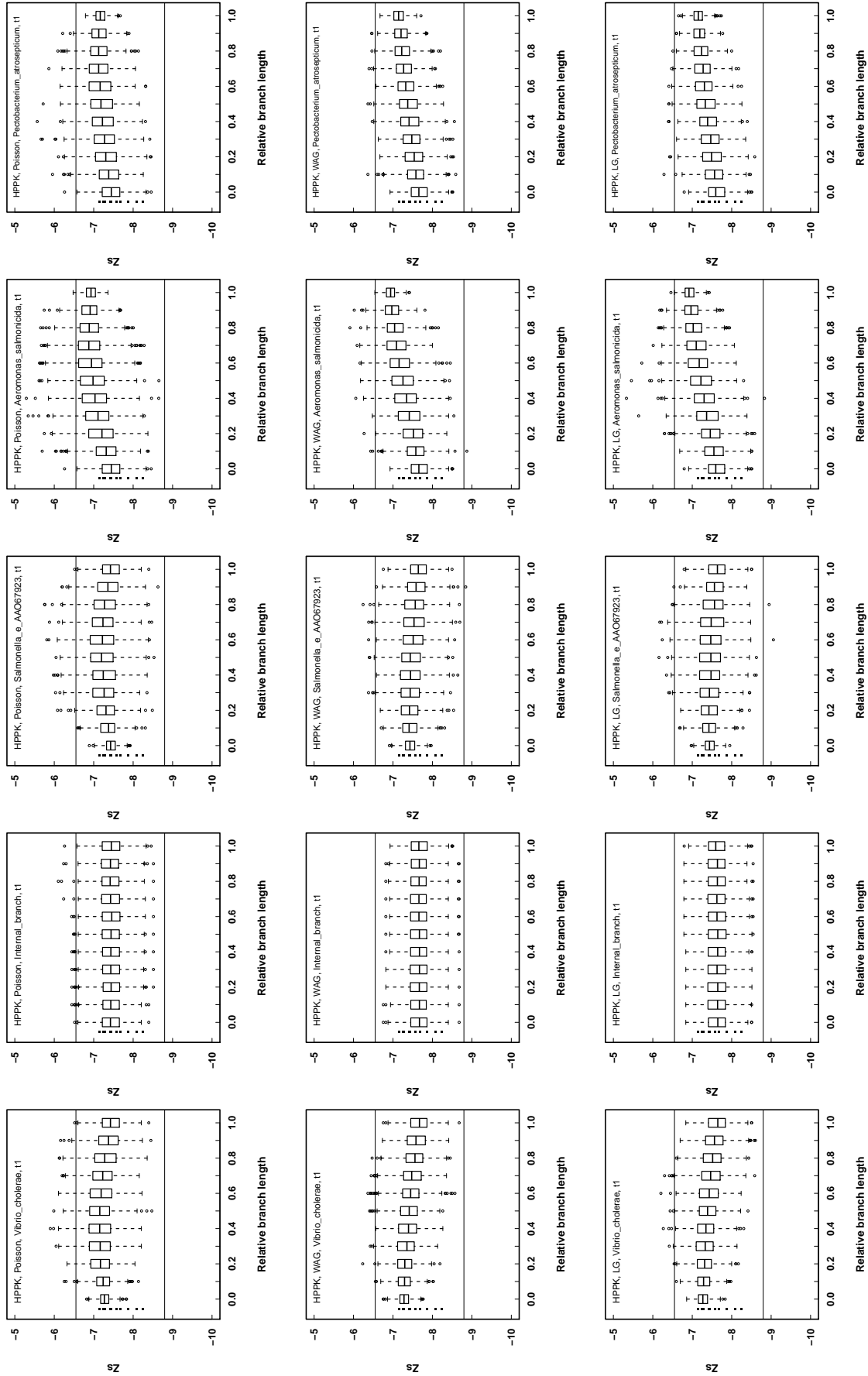


Figure 23. HPPK,  $Z_s$ , Tree 1

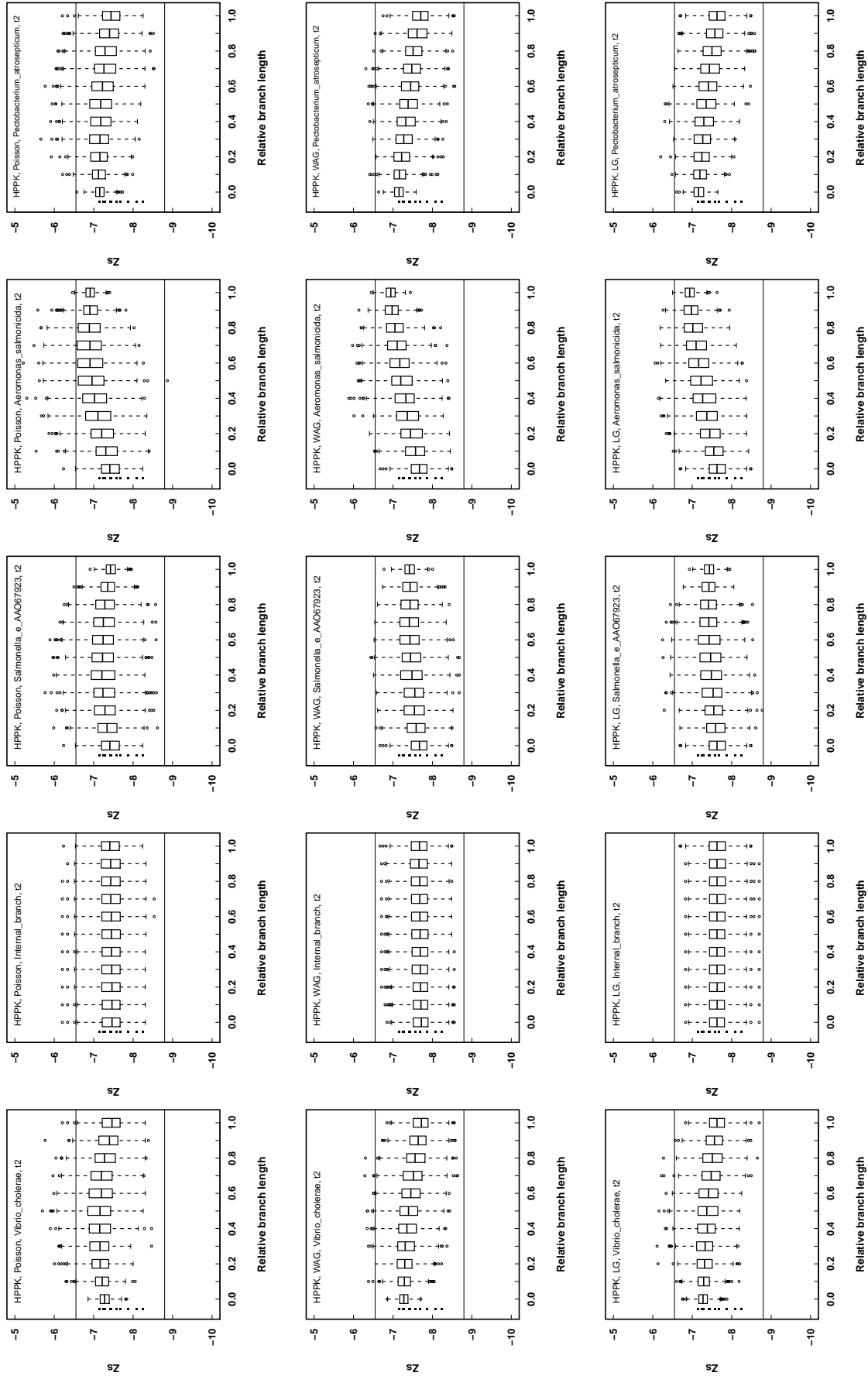


Figure 24. HPPK,  $Z_s$ , Tree 2

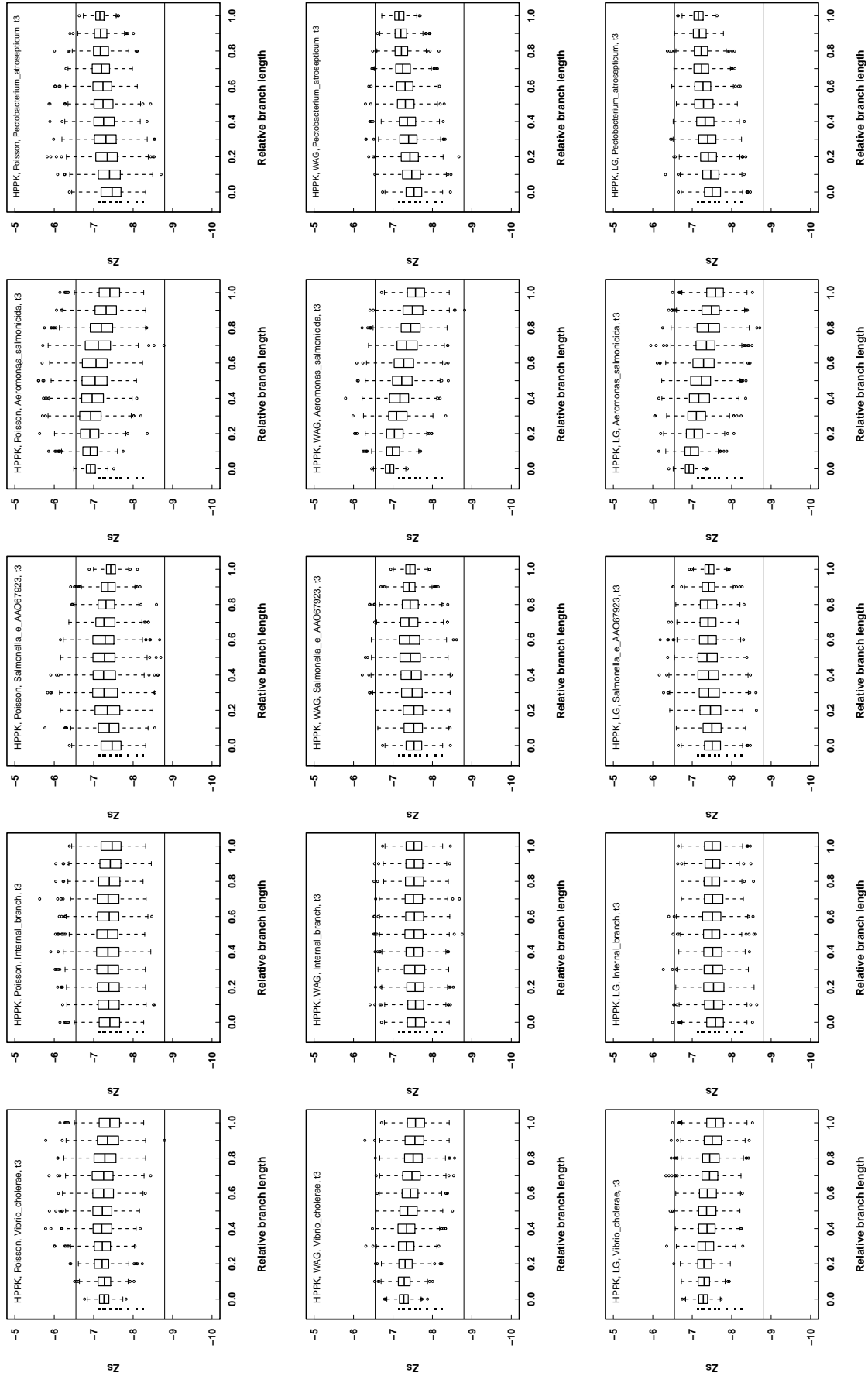


Figure 25. HPPK,  $Z_s$ , Tree 3

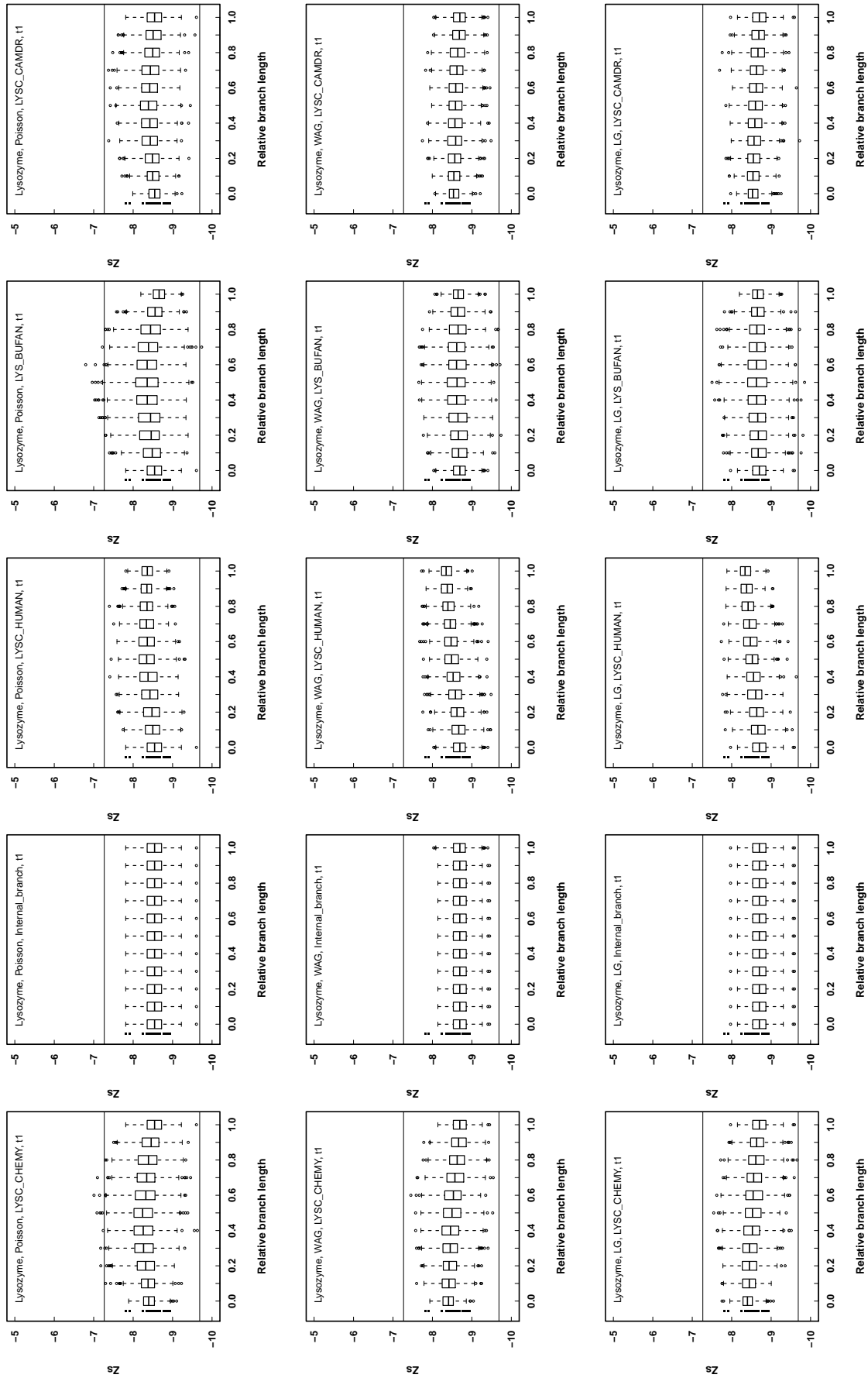


Figure 26. Lysozyme,  $Z_s$ , Tree 1

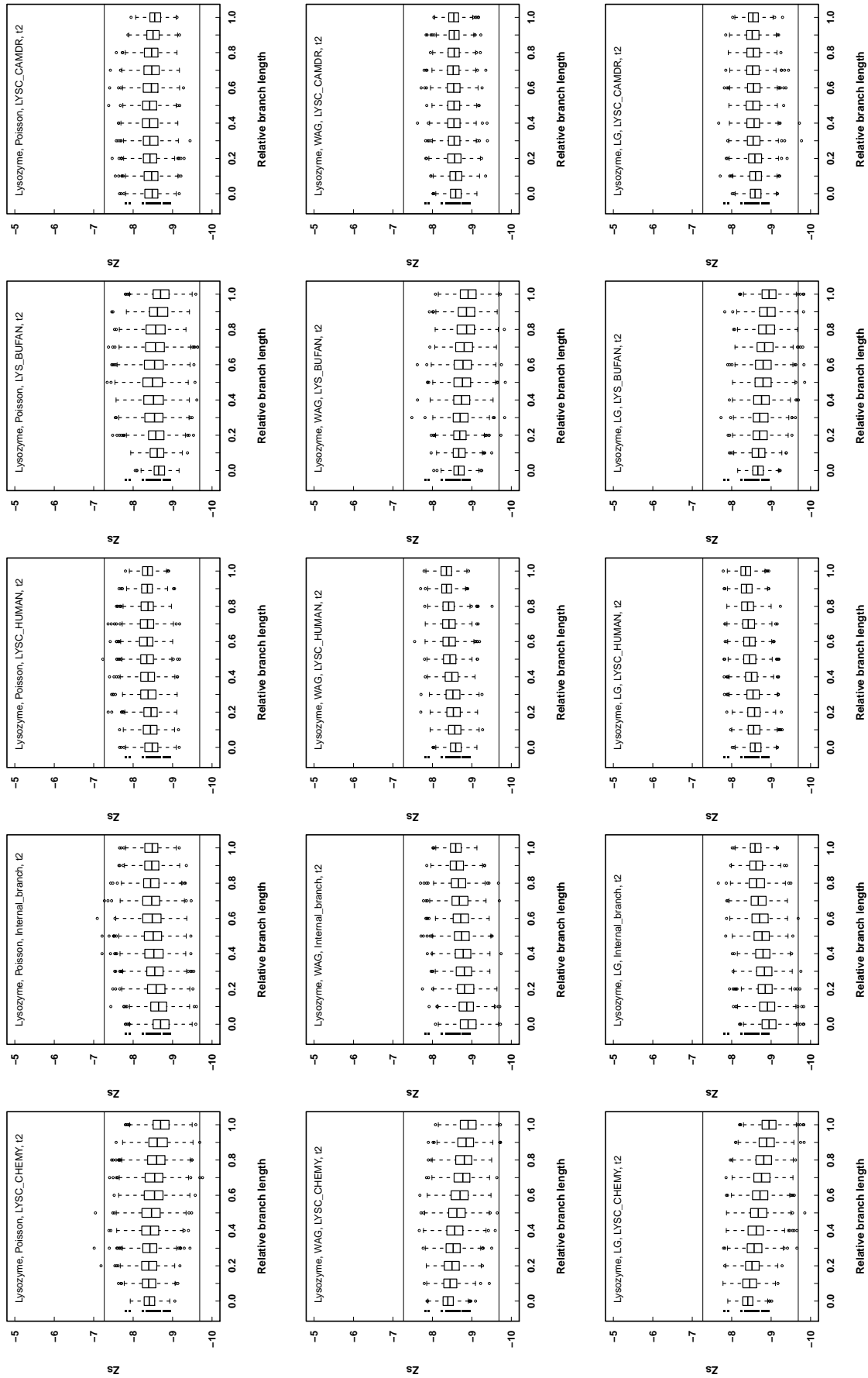


Figure 27. Lysozyme,  $Z_s$ , Tree 2



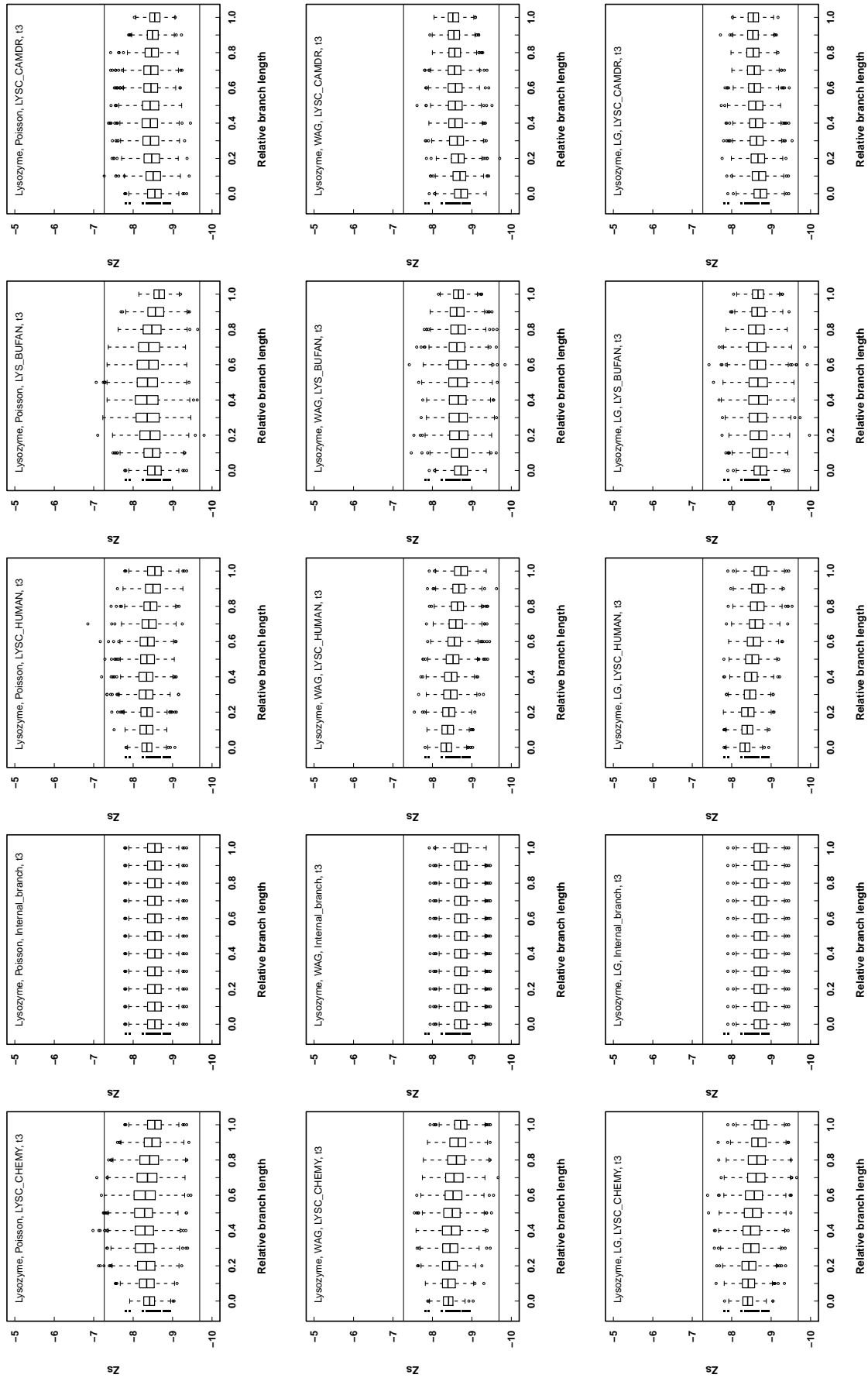


Figure 28. Lysozyme,  $Z_s$ , Tree 3

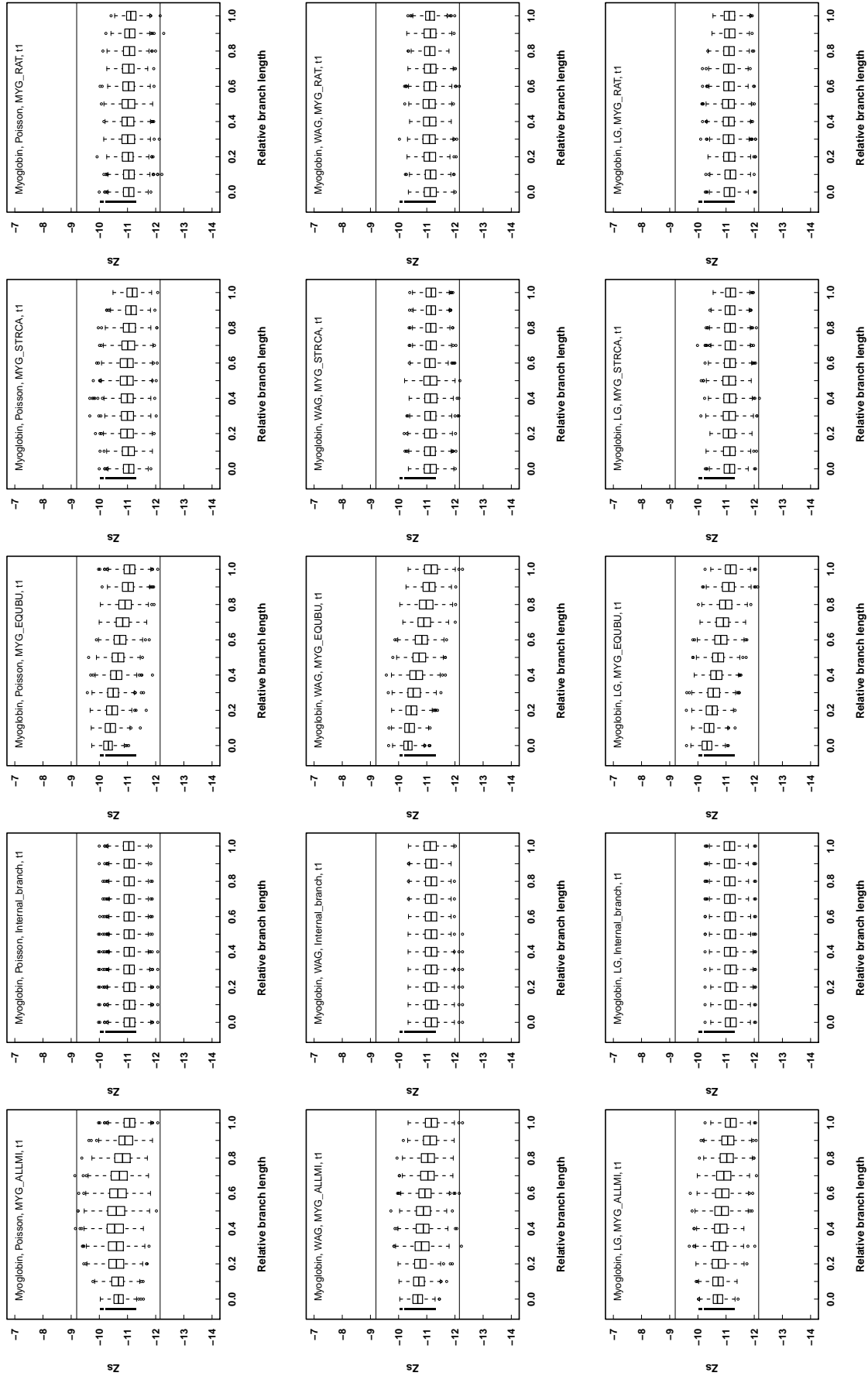


Figure 29. Myoglobin,  $Z_s$ , Tree 1

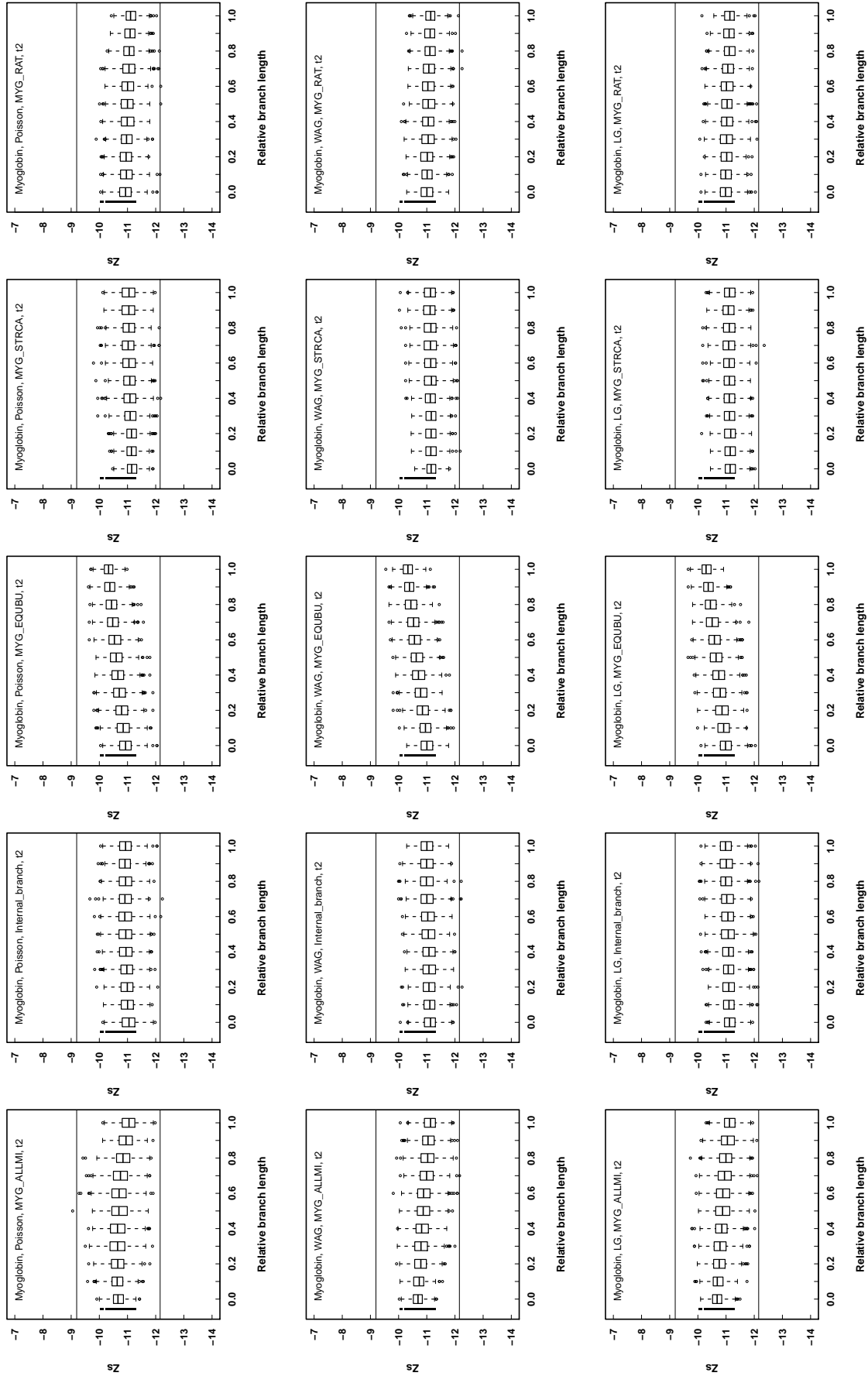


Figure 30. Myoglobin,  $Z_s$ , Tree 2

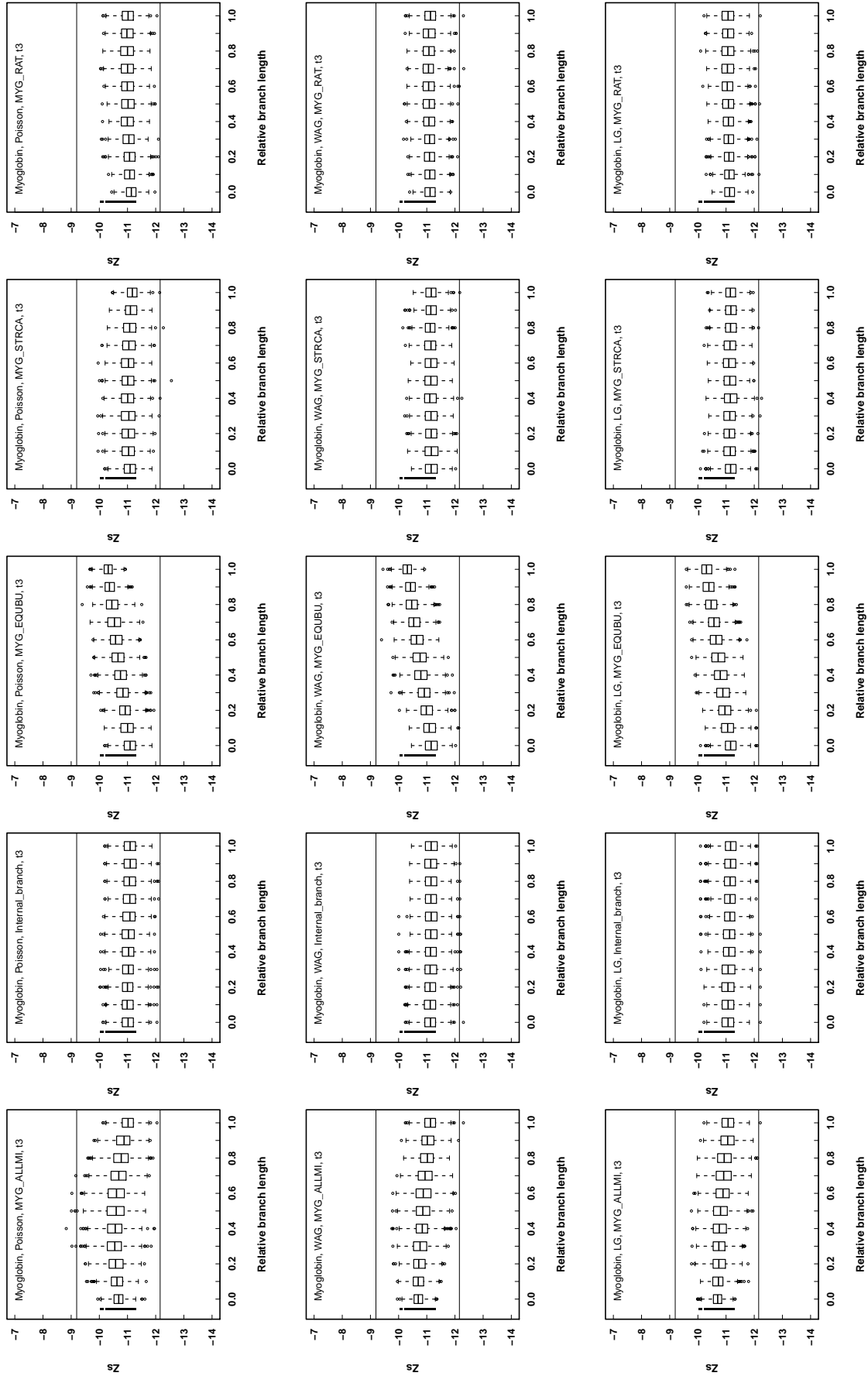


Figure 31. Myoglobin,  $Z_s$ , Tree 3

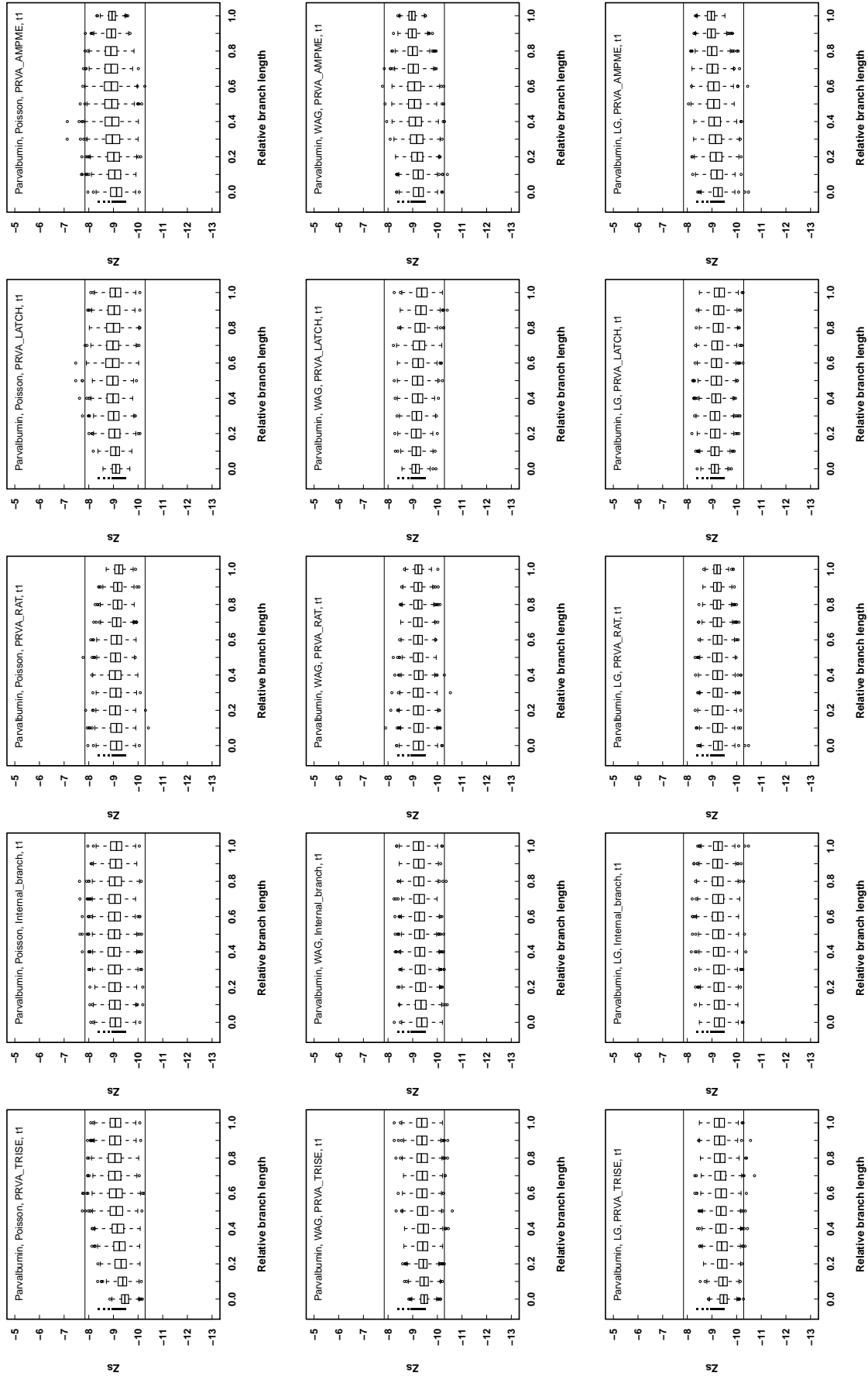


Figure 32. Parvalbumin,  $Z_s$ , Tree 1

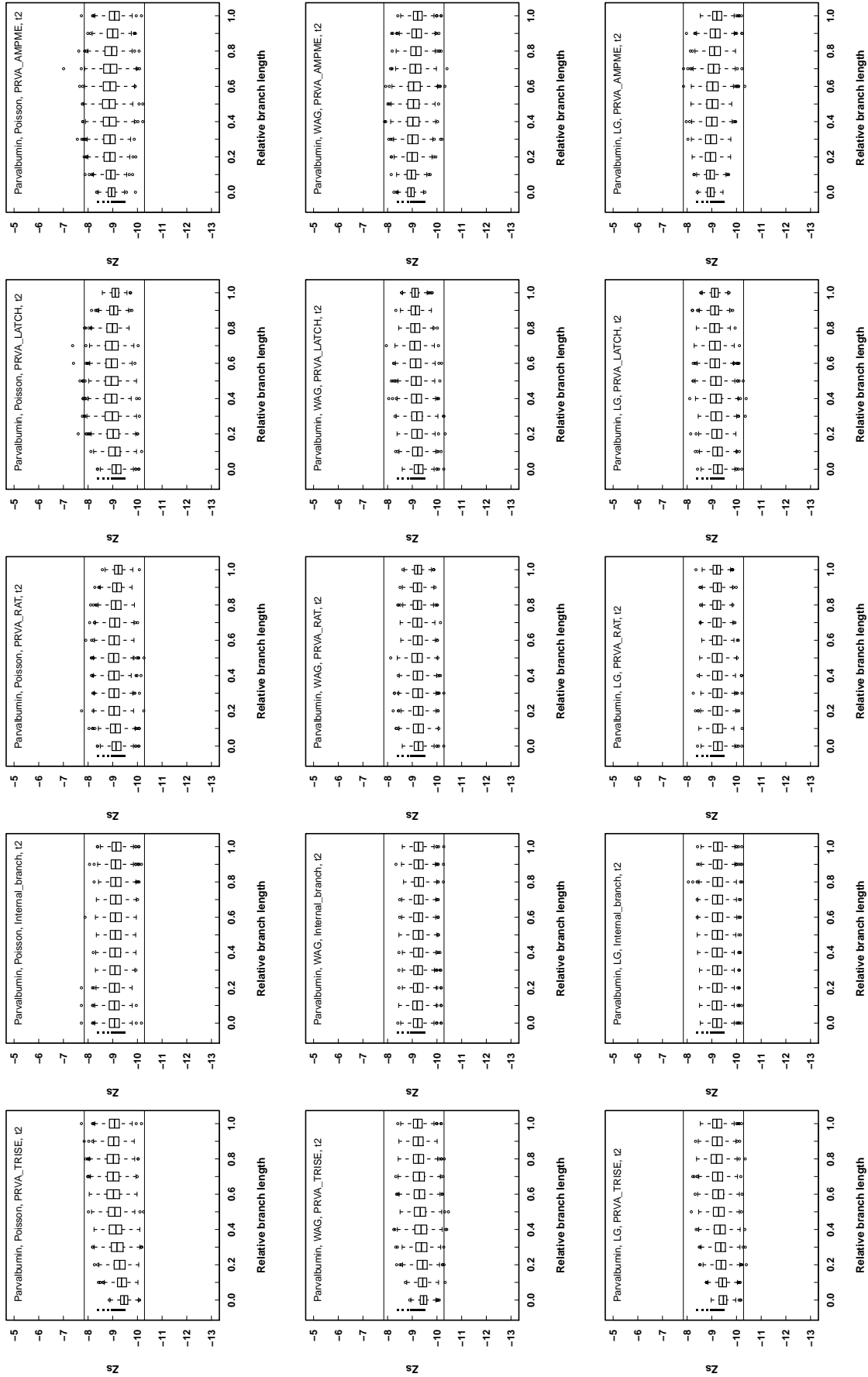


Figure 33. Parvalbumin,  $Z_s$ , Tree 2

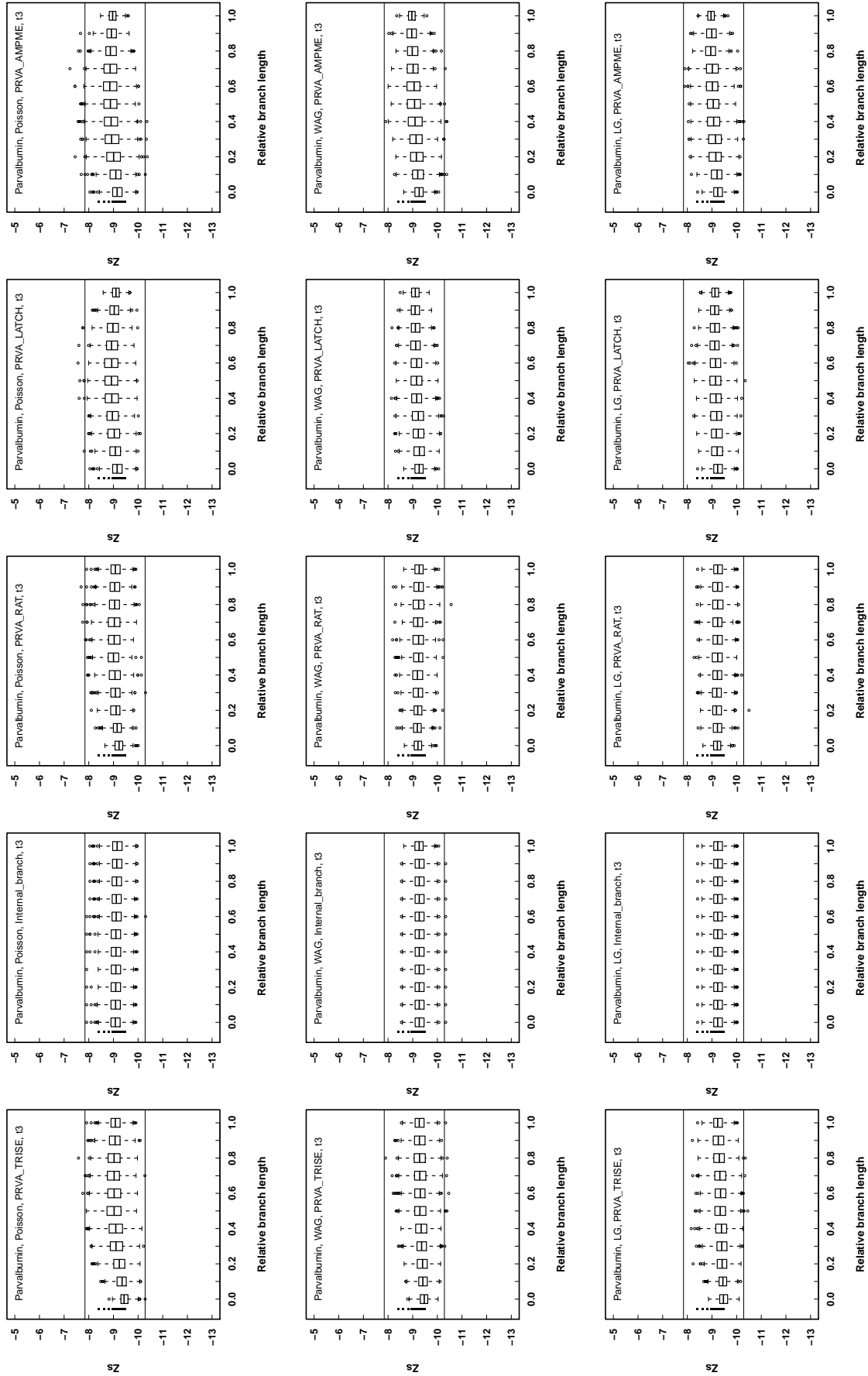
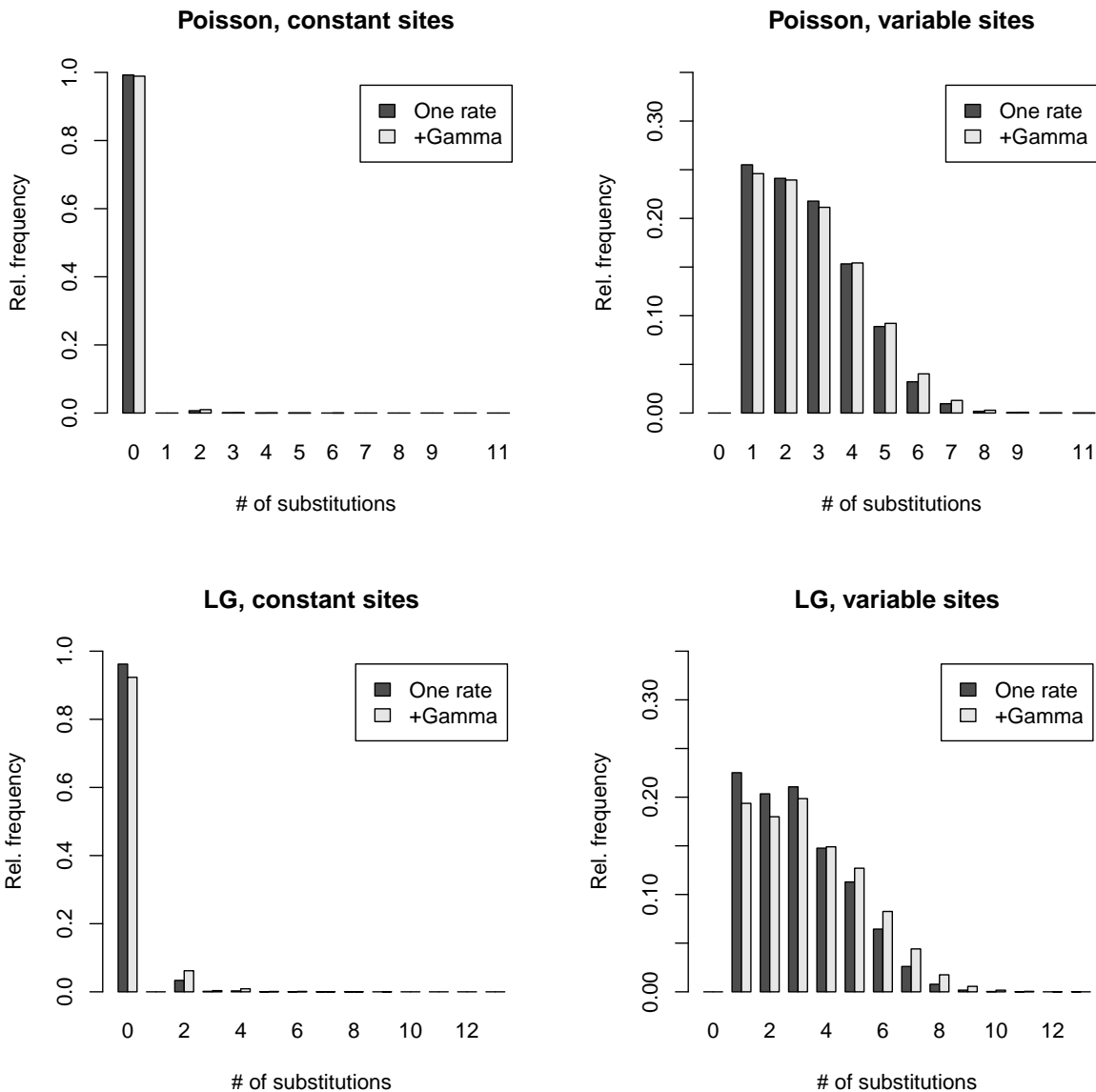


Figure 34. Parvalbumin,  $Z_s$ , Tree 3

## Rate variation (discrete $\Gamma$ )



**Figure 35.** HPPK, 4 taxa, tree 3. Comparison of the number of substitutions by site type (conserved across all four taxa or variable) with and without rate heterogeneity. Paths were sampled with site-specific rates, estimated with PAML (using 30 categories, Yang, 2007). Rate variation led to increased branch lengths and more multiple substitutions. It was demonstrated in the paper that, for our data sets, decrease in structural compatibility was mainly associated with amino acids that were not observed in the terminals. This may be an explanation why we observed slightly more sequences with worse sequence:structure fit when rate heterogeneity was used. Results were similar for all trees and proteins.



## References

Yang, Z. (2007). Paml 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol*, **24**(8), 1586–91.