

Prediction and identification of functional cisNATs in plants

Veerendra Gadekar

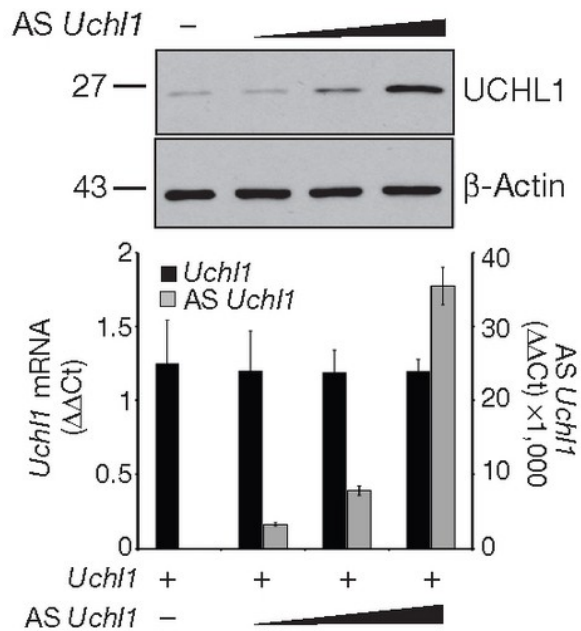
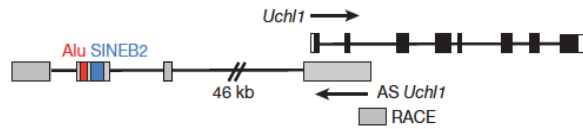
Department of Theoretical Chemistry

University of Vienna

Bled, Feb 17, 2017

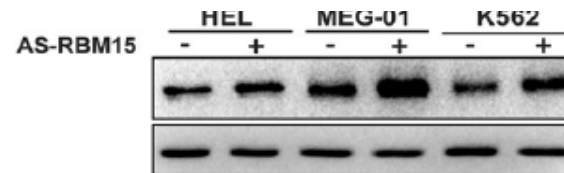
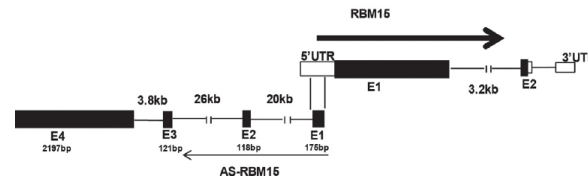
Functional antisense long noncoding RNAs

Mouse *Uchl1* / *cis*-NAT pair

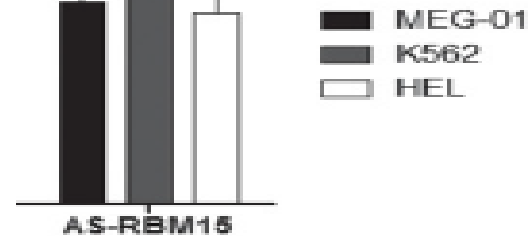


Carrieri et al. 2012

Human *RBM15* / *cis*-NAT pair

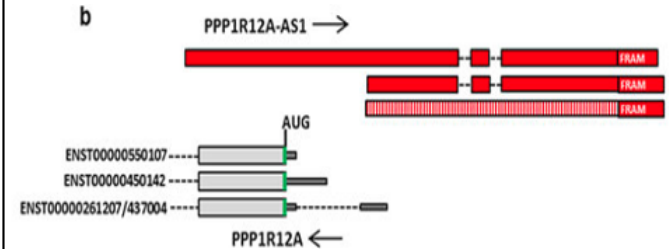


M15

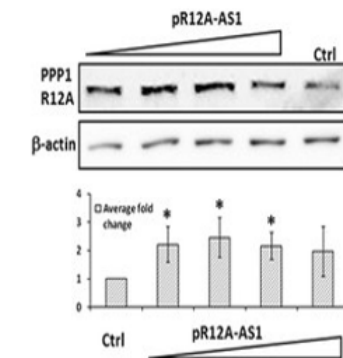


Tran et al. 2016

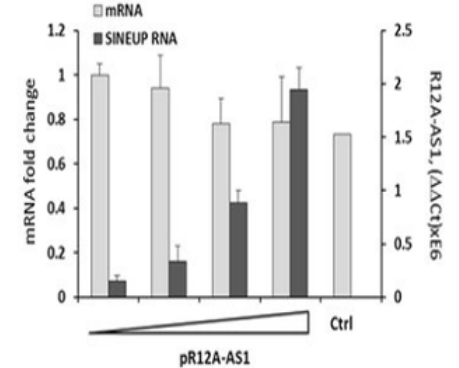
Human *PPP1R12A* / *cis*-NAT pair



d

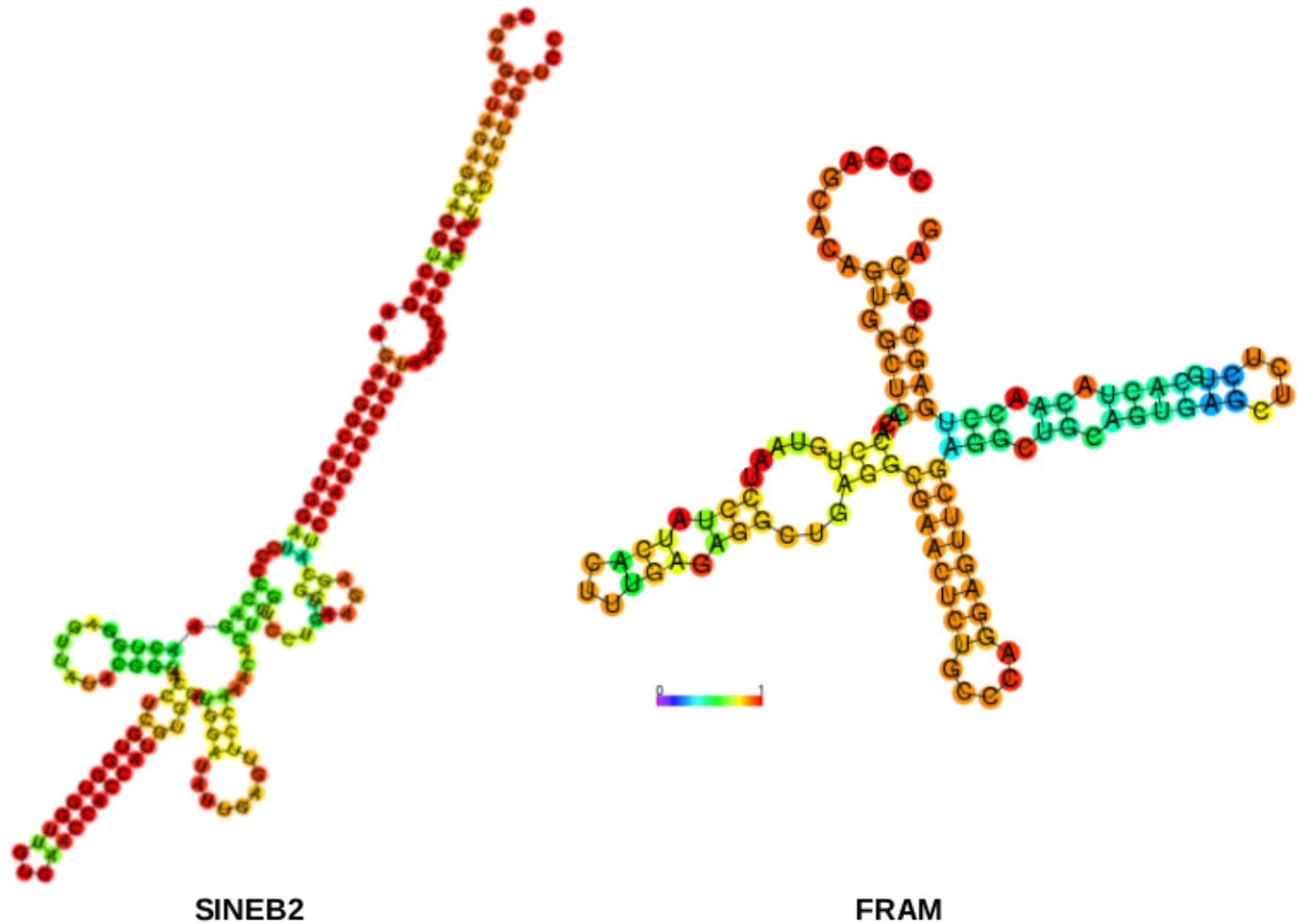


e



Schein et al. 2016

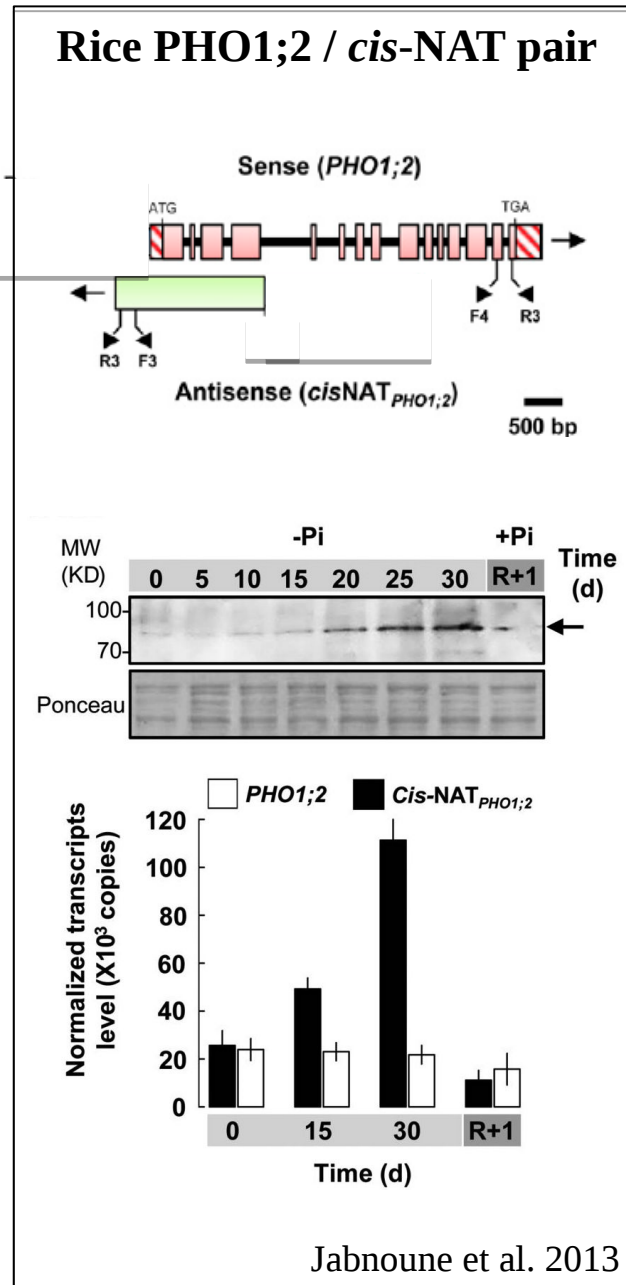
SINEB2 and FRAM are predicted to fold into a different structures



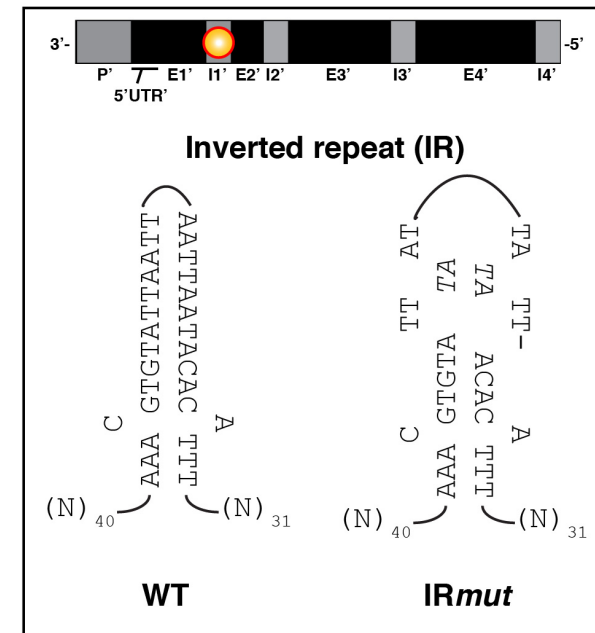
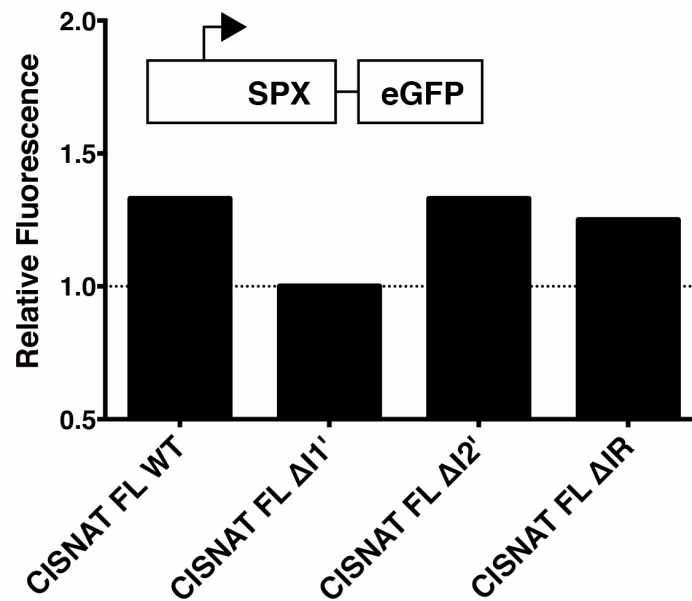
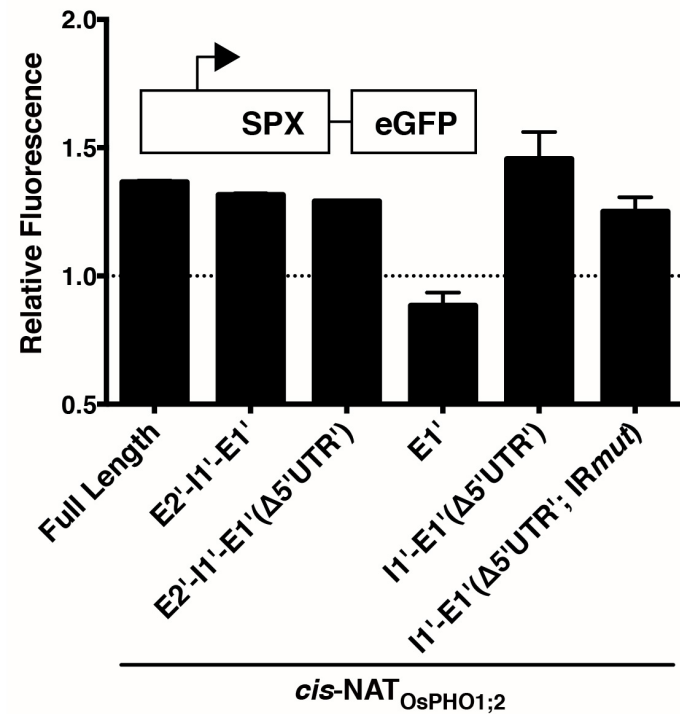
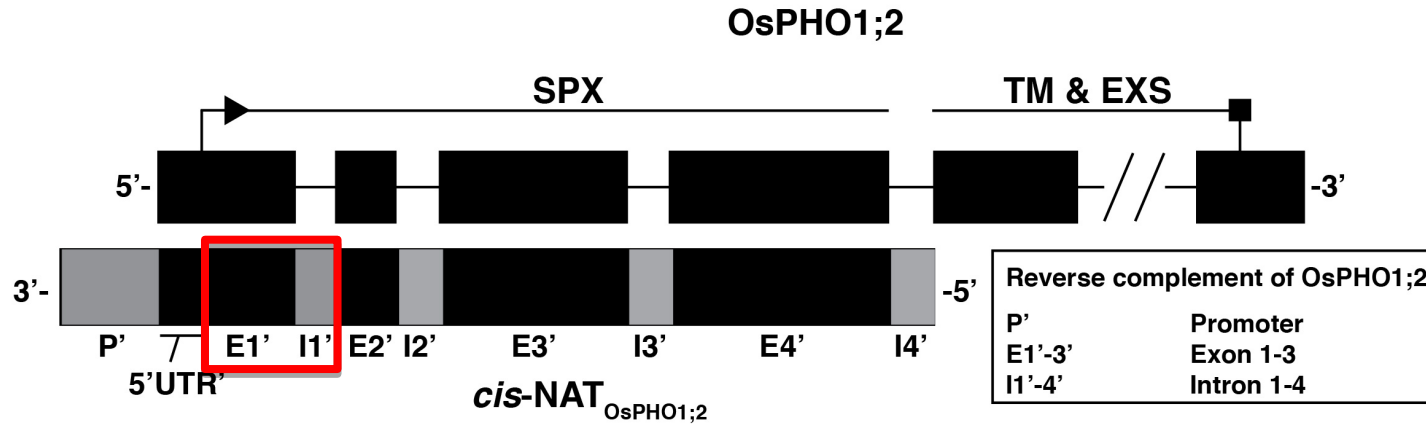
Supplementary Figure S5. Human FRAM and mouse SINEB2 elements obtain different secondary structures

Secondary structure of SINEs, predicted by RNAfold program. Prediction confidence for each base is indicated by the color code.

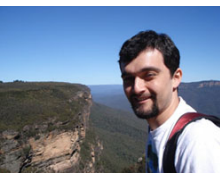
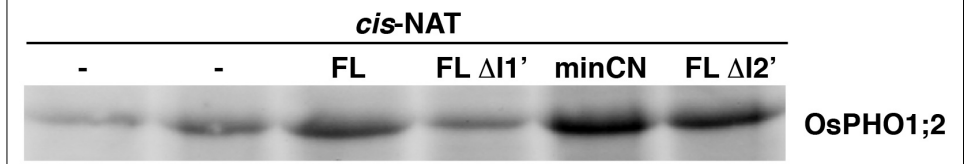
A Rice cis-Natural Antisense RNA



Dissection of *cis-NAT*_{OsPHO1;2} translation enhancement of OsPHO1;2



Arabidopsis protoplast

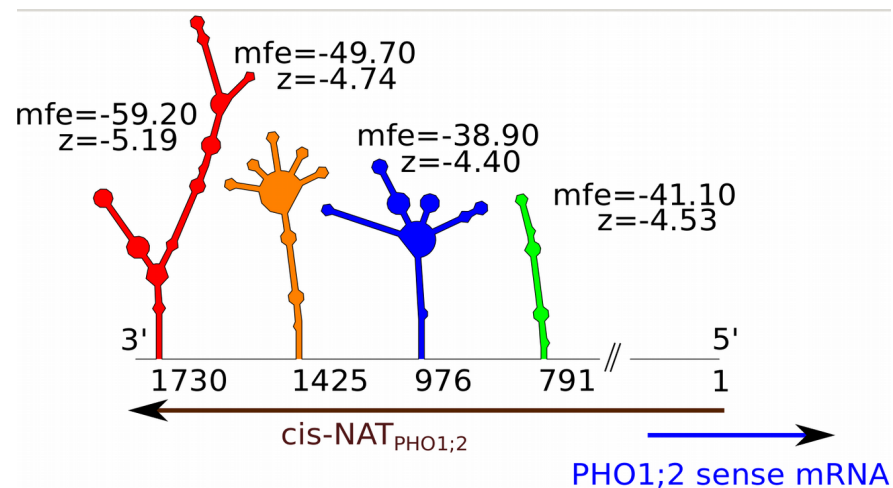


Aim

Mechanism underlying the translation enhancement of PHO1 by *cis*-NAT in rice

Objectives

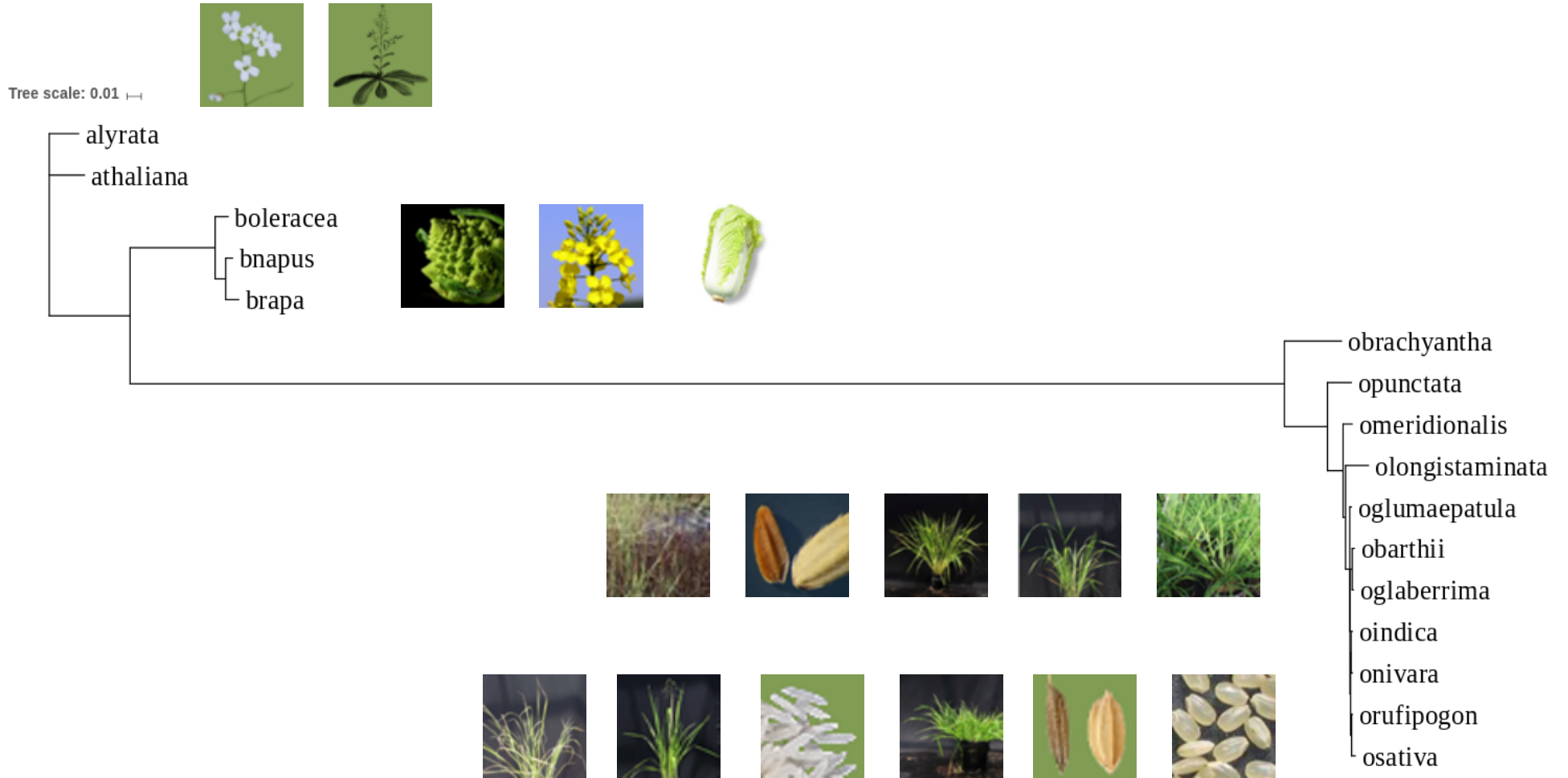
- Identification of similar putative *cis*-NATs in plant genomes
- Application of comparative genomics to identify evolutionary conserved structure elements
- Select candidate structures for experimental testing



Data set

Species	Genome Version	Genome size in Mbp	No. of chromosomes
<i>Oryza sativa</i>	IRGSP-1.0	~374	12
<i>Oryza rufipogon</i>	OR_W1943	~338	12
<i>Oryza punctata</i>	AVCL000000000	~394	12
<i>Oryza nivara</i>	AWHD000000000	~338	12
<i>Oryza meridionalis</i>	Oryza_meridionalis_v1.3	~336	12
<i>Oryza longistaminata</i>	O_longistaminata_v1.0	~326	60198 scaffold
<i>Oryza indica</i>	ASM465v1	~412	12
<i>Oryza glumaepatula</i>	ALNU020000000	~372	12
<i>Oryza glaberrima</i>	AGI1.1	~316	12
<i>Oryza brachyantha</i>	Oryza_brachyantha.v1.4b	~260	12
<i>Oryza barthii</i>	O.barthii_v1	~308	12
<i>Brassica rapa</i>	IVFCAASv1	~284	10
<i>Brassica oleracea</i>	v2.1	~489	9
<i>Brassica napus</i>	AST_PRJEB5043_v1	~738	20899 supercontigs
<i>Arabidopsis thaliana</i>	TAIR10	~135	5
<i>Arabidopsis lyrata</i>	v.1.0	~206	8

Phylogenetic tree



No. of annotated transcripts

Species	Protein coding	ncRNA	Other	Total
<i>Oryza sativai</i>	42132	55619	0	97751
<i>Oryza rufipogon</i>	47441	2020	758	50219
<i>Oryza punctata</i>	41060	28	5167	46255
<i>Oryza nivara</i>	48360	24	1648	50032
<i>Oryza meridionalis</i>	43455	2283	0	45738
<i>Oryza longistaminata</i>	31686	0	0	31686
<i>Oryza indica</i>	40745	2232	45461	88438
<i>Oryza glumaepatula</i>	46893	19	3264	50176
<i>Oryza glaberrima</i>	33164	1997	38717	73878
<i>Oryza brachyantha</i>	32037	18	2100	34155
<i>Oryza barthii</i>	41595	1681	665	43941
<i>Brassica rapa</i>	41025	1	1827	42853
<i>Brassica oleracea</i>	59220	5	0	59225
<i>Brassica napus</i>	101040	0	0	101040
<i>Arabidopsis thaliana</i>	48321	3917	1775	54013
<i>Arabidopsis lyrata</i>	32667	0	0	32667

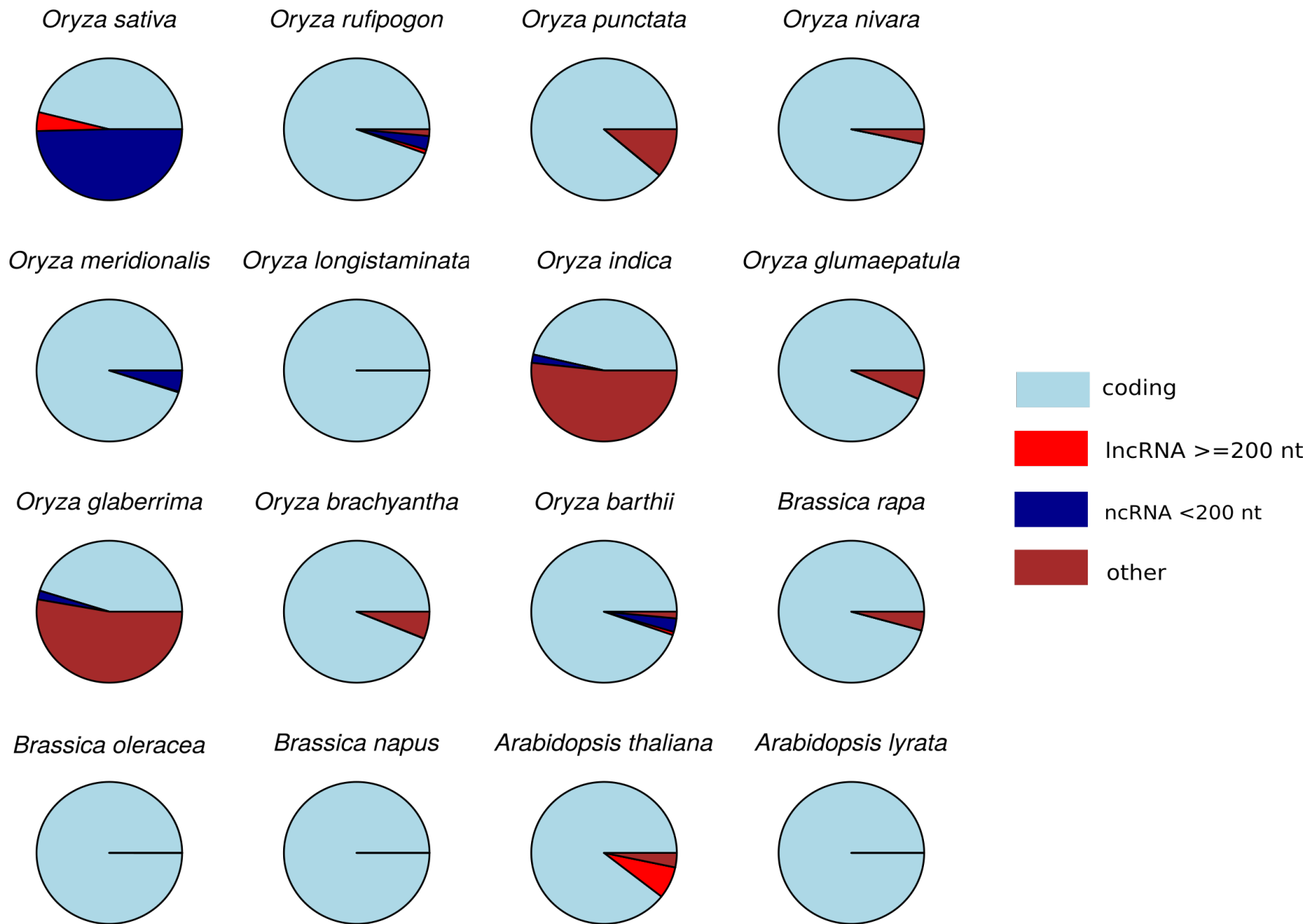
No. of annotated transcripts

Species	Protein coding	ncRNA	Other	Total
<i>Oryza sativa</i>	42132	55619	0	97751
<i>Oryza rufipogon</i>	47441	2020	758	50219
<i>Oryza punctata</i>	41060	28	5167	46255
<i>Oryza nivara</i>	48360	24	1648	50032
<i>Oryza meridionalis</i>	43455	2283	0	45738
<i>Oryza longistaminata</i>	31686	0	0	31686
<i>Oryza indica</i>	40745	2232	45461	88438
<i>Oryza glumaepatula</i>	46893	19	3264	50176
<i>Oryza glaberrima</i>	33164	1997	38717	73878
<i>Oryza brachyantha</i>	32037	18	2100	34155
<i>Oryza barthii</i>	41595	1681	665	43941
<i>Brassica rapa</i>	41025	1	1827	42853
<i>Brassica oleracea</i>	59220	5	0	59225
<i>Brassica napus</i>	101040	0	0	101040
<i>Arabidopsis thaliana</i>	48321	3917	1775	54013
<i>Arabidopsis lyrata</i>	32667	0	0	32667

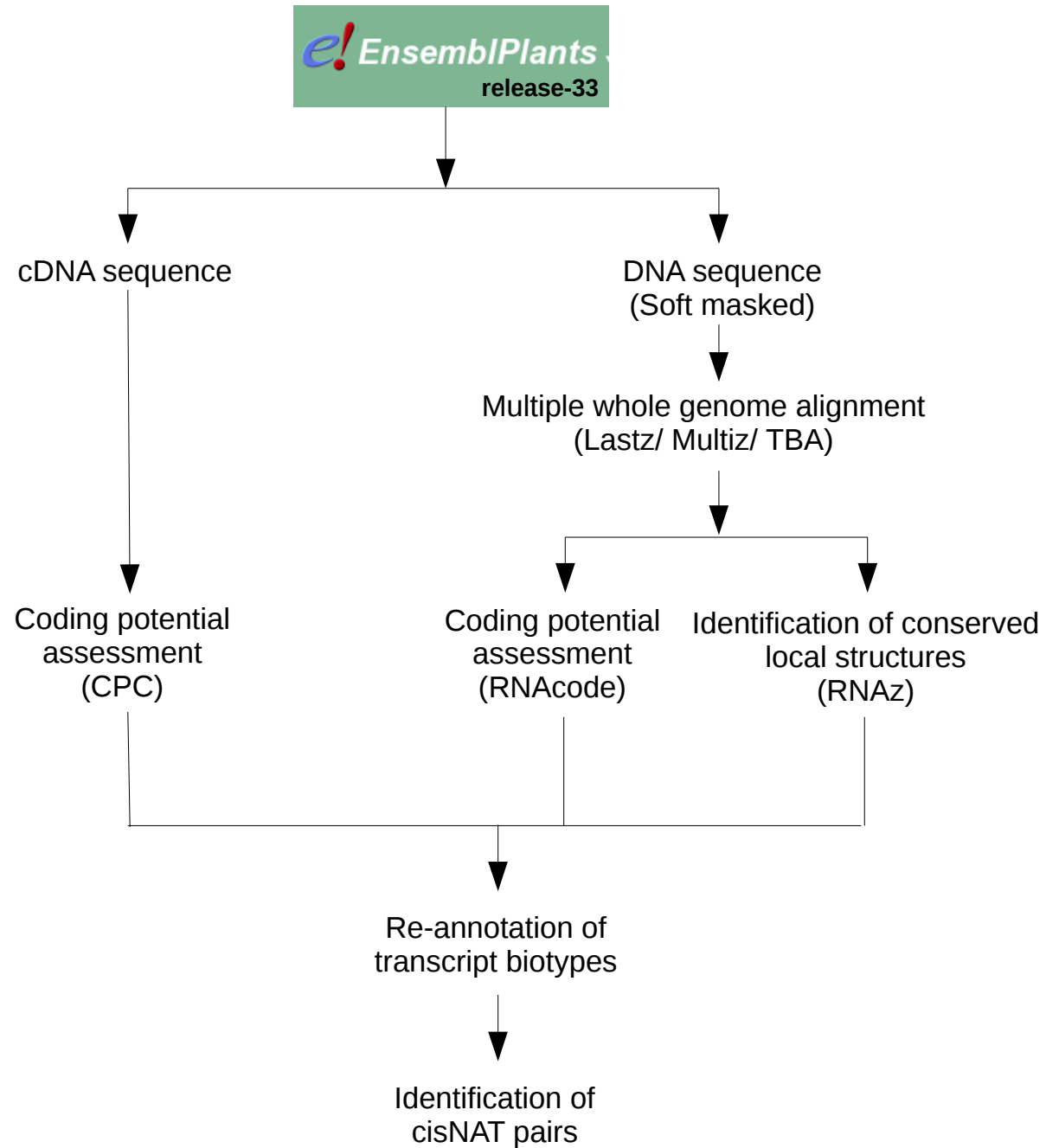
No. of annotated transcripts

Species	Protein coding	lncRNA nt \geq 200	ncRNA nt $<$ 200	Other
Oryza sativa	42132	3923	45310	0
Oryza rufipogon	47441	440	1580	758
Oryza punctata	41060	4	24	5167
Oryza nivara	48360	1	23	1648
Oryza meridionalis	43455	82	2201	0
Oryza longistaminata	31686	0	0	0
Oryza indica	40745	15	1669	45461
Oryza glumaepatula	46893	2	17	3264
Oryza glaberrima	33164	18	1459	38717
Oryza brachyantha	32037	1	17	2100
Oryza barthii	41595	326	1355	665
Brassica rapa	41024	0	1	1827
Brassica oleracea	59219	5	0	0
Brassica napus	101040	0	0	0
Arabidopsis thaliana	48321	3916	1	1775
Arabidopsis lyrata	32663	0	0	0

Percentage of different transcript biotypes



Work flow



Coding potential Calculator (CPC)

Quality of the ORF

- Implements *framefinder* software to identify the longest reading frame in the three forward frames
- Accesses the integrity of the ORF by checking if the ORF begins with a start codon and ends with an in-frame stop codon

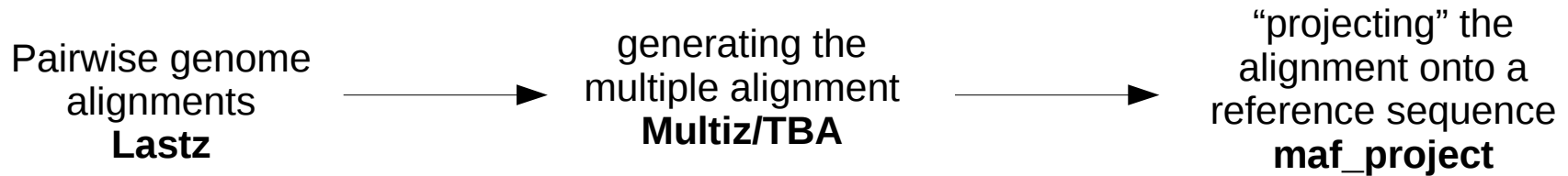
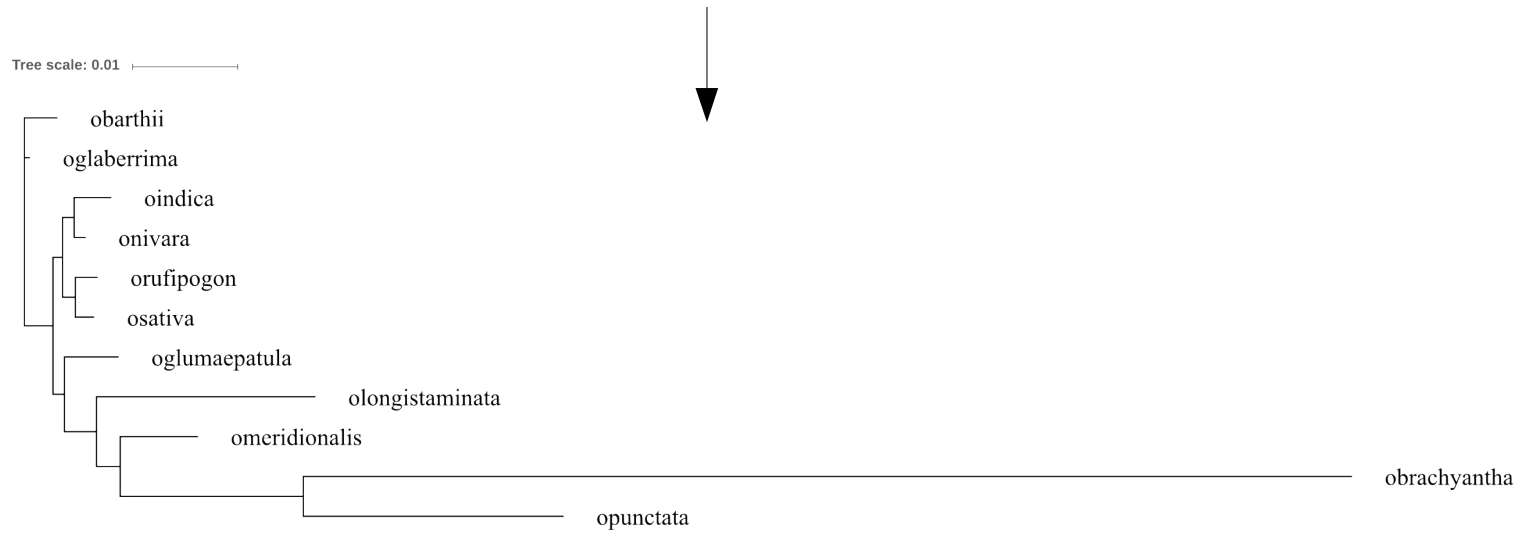
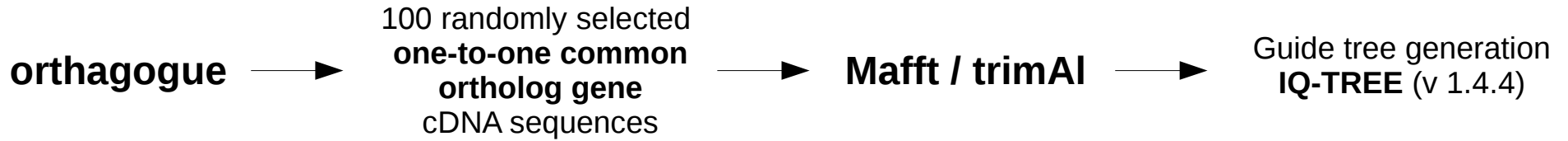
BLASTX against uniref

- Extracts the number of hits
- Quality of HSPs (Lower E-values)
- Distribution of HSPs with in the frame

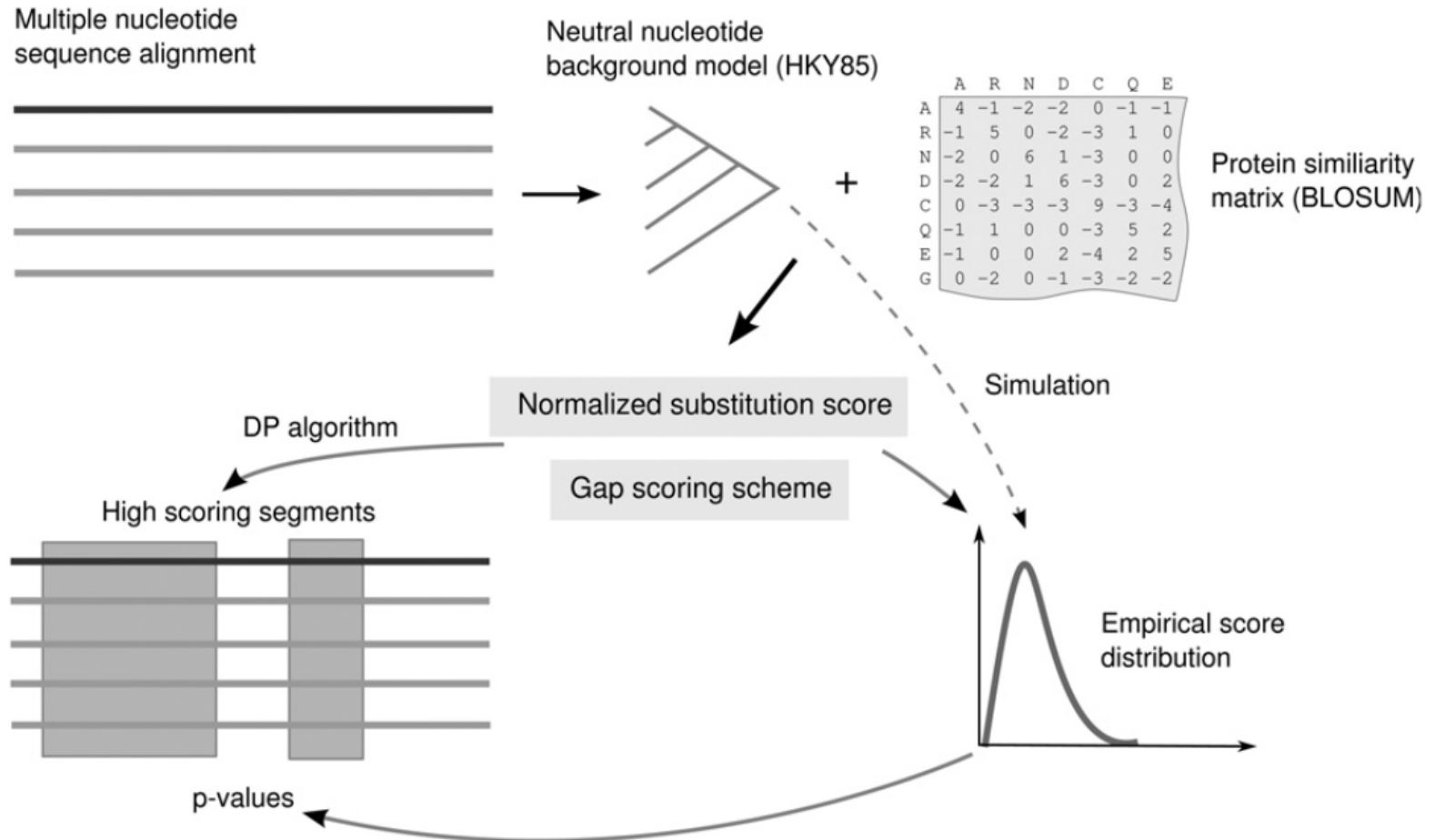
Support vector machine (SVM) machine learning classifier

- Transcripts with Score < 1.5 is classified as noncoding

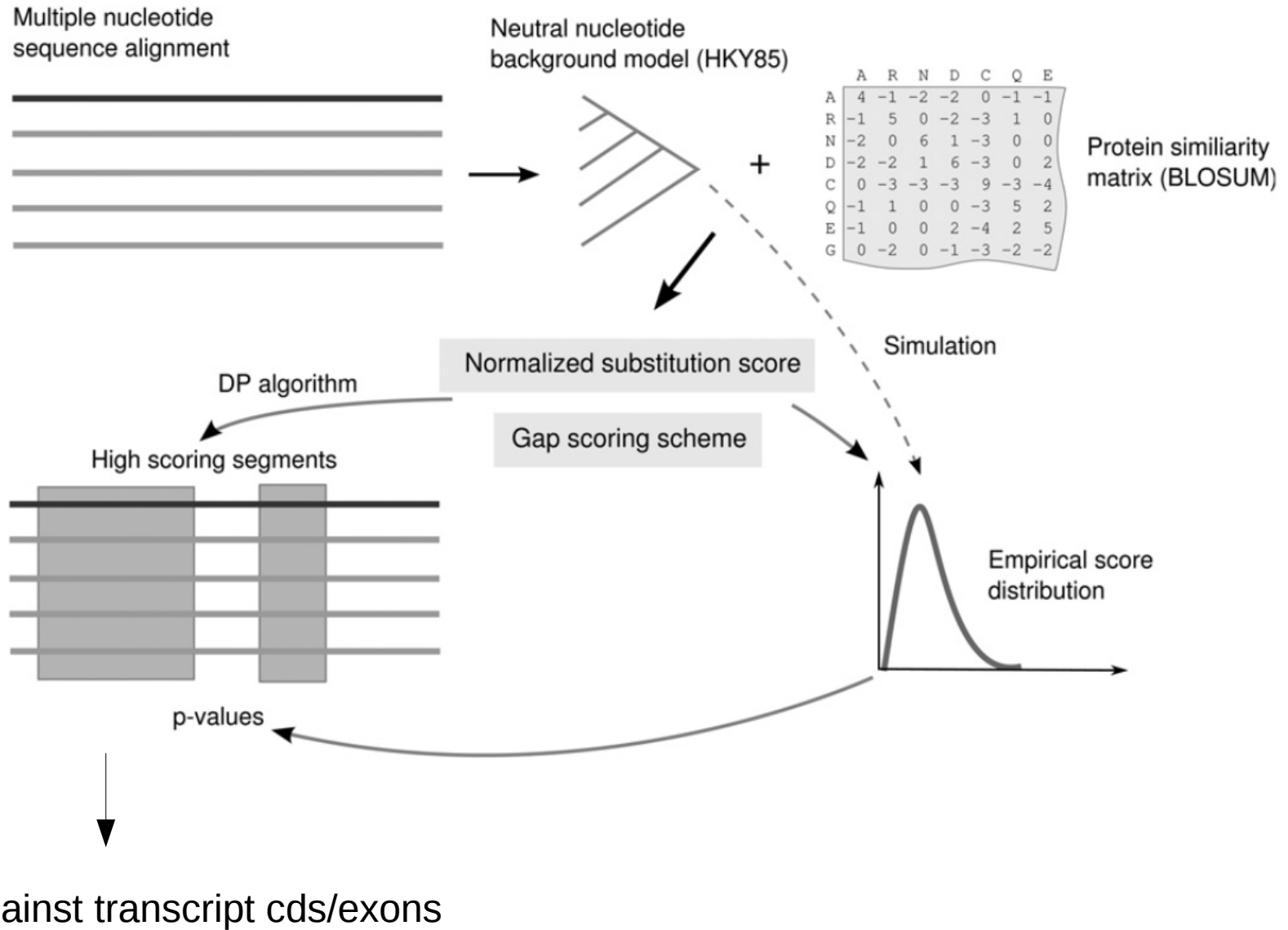
Multiple whole genome alignment



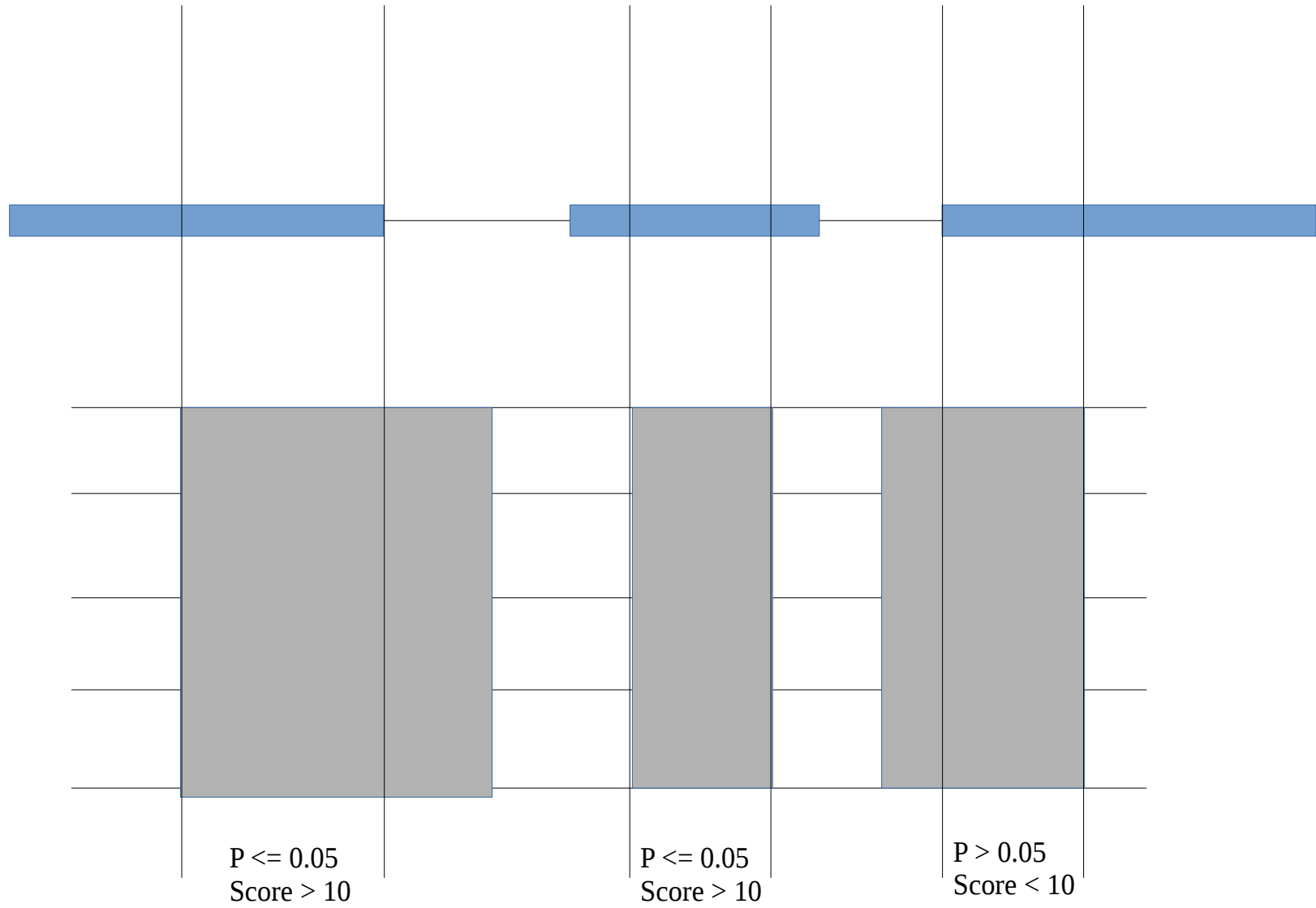
RNAcode



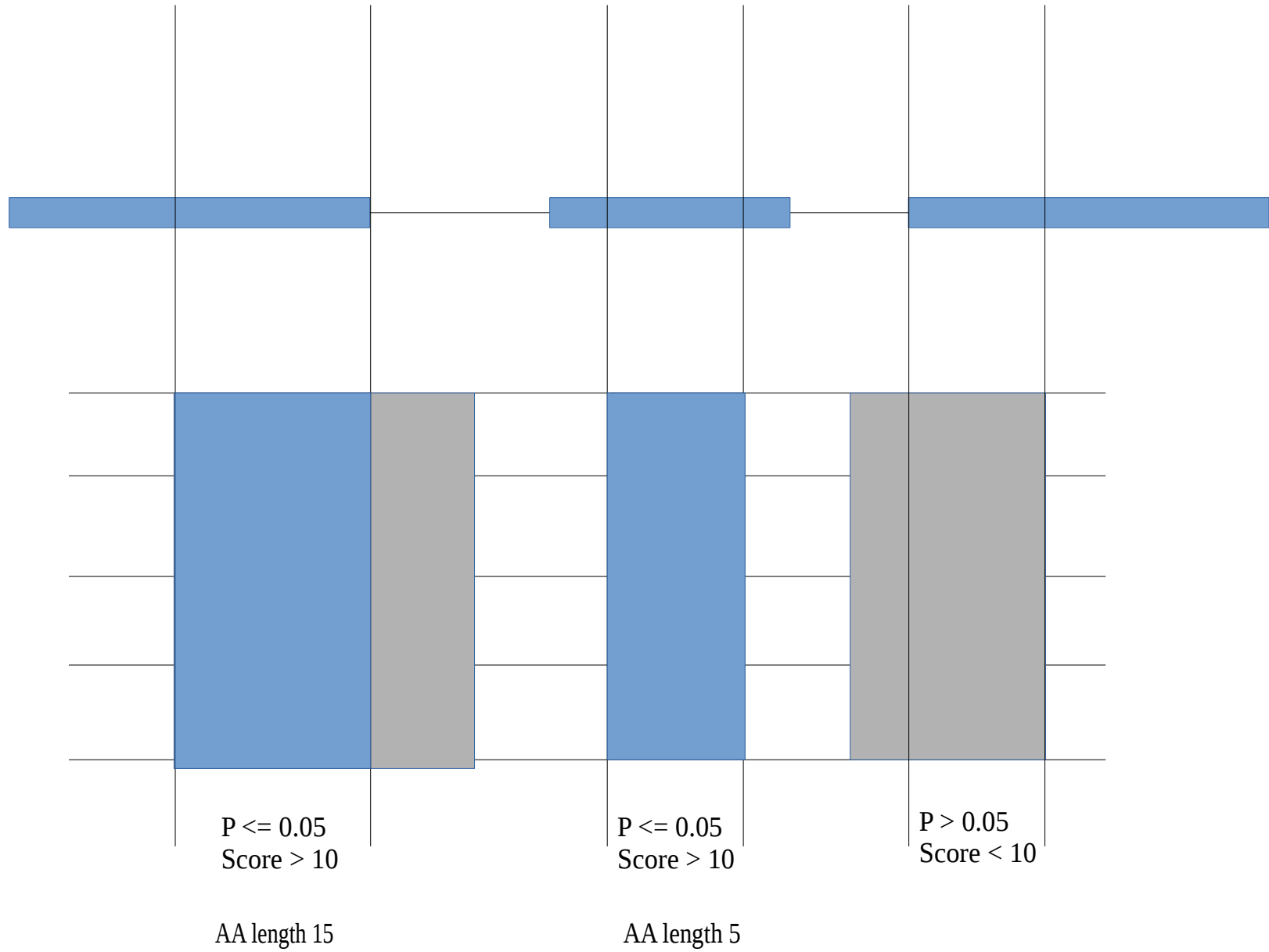
RNAcode



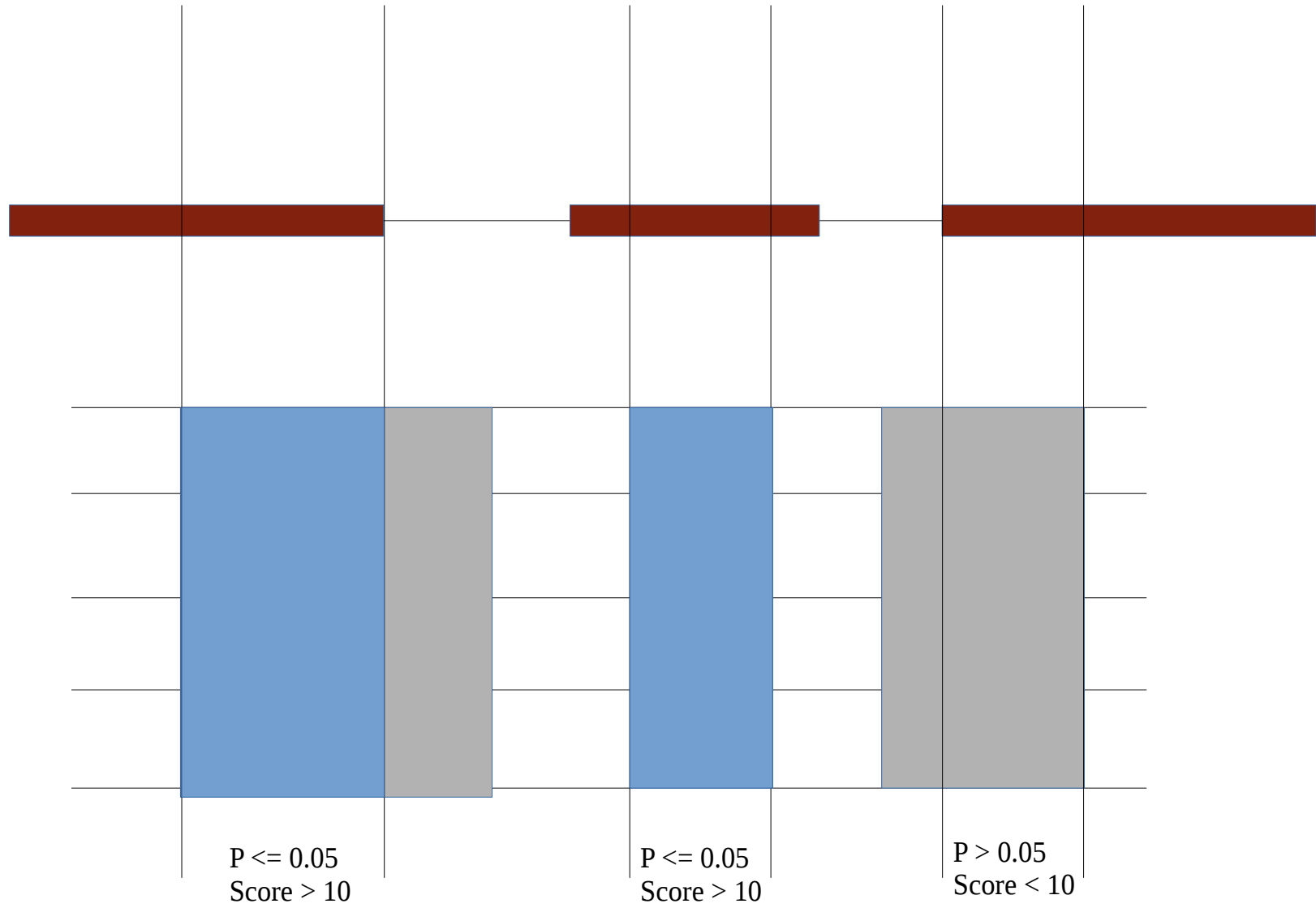
RNAcode



RNAcode

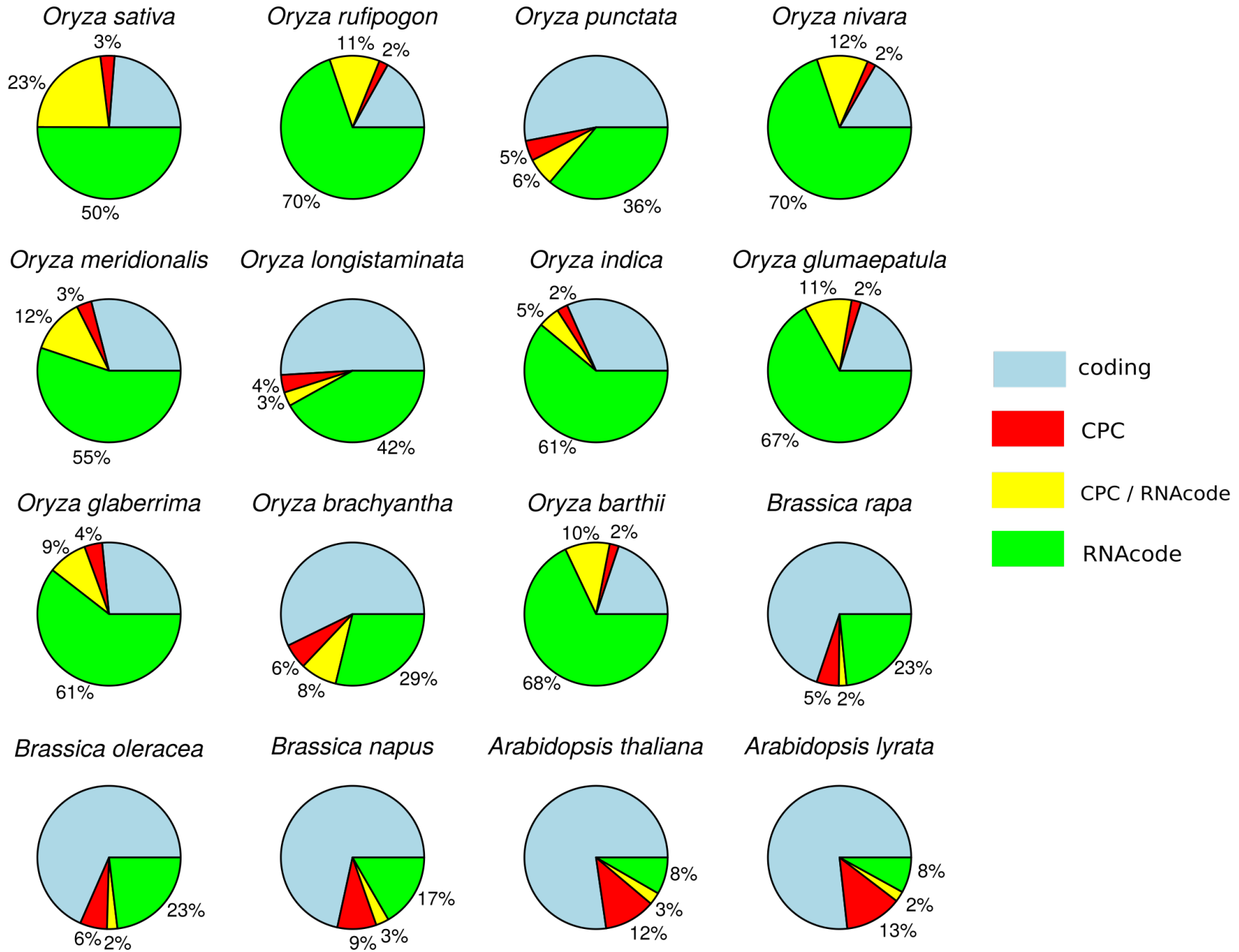


RNAcode

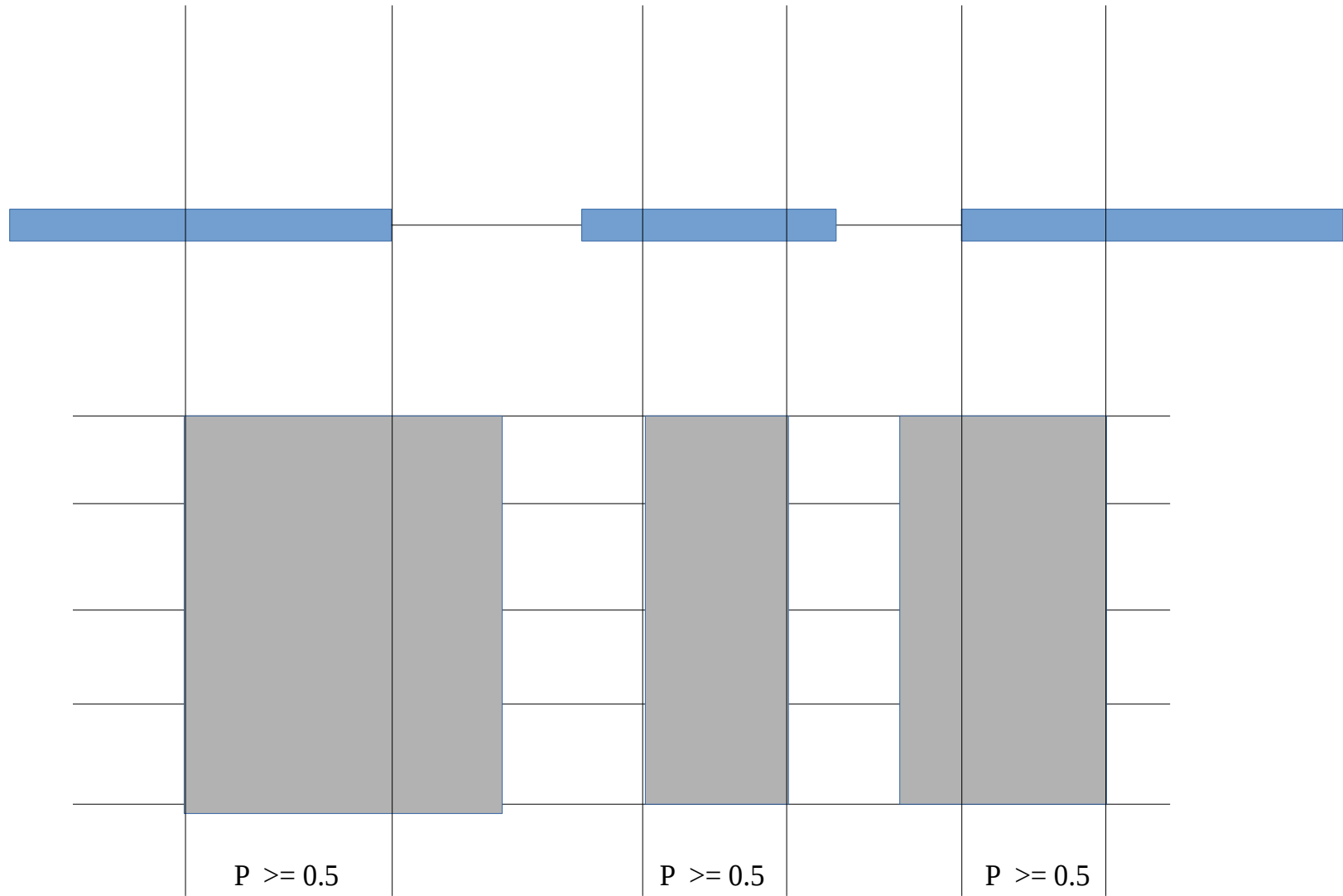


20 Amino Acids

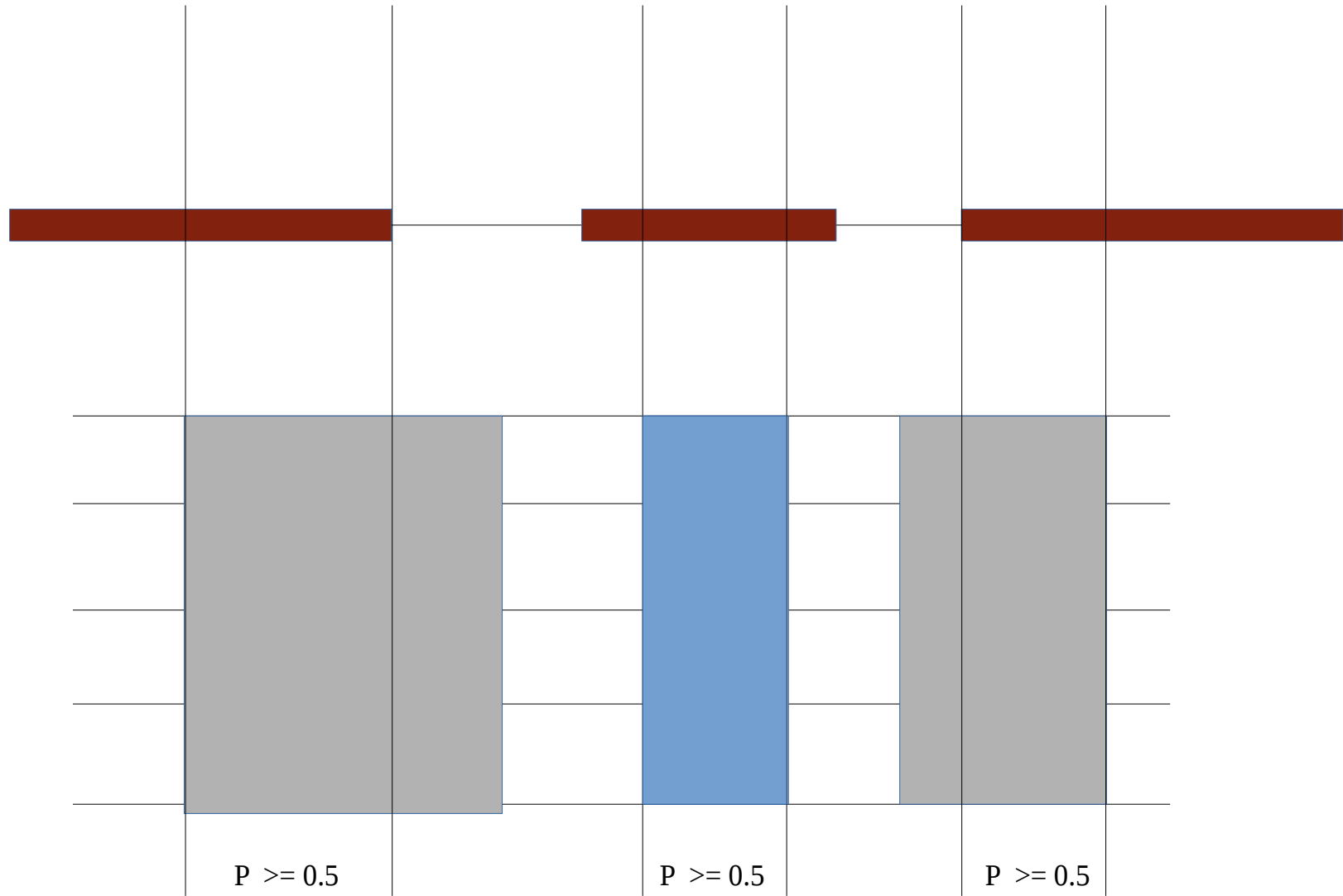
Coding potentials identified by RNAcode vs CPC



RNAz



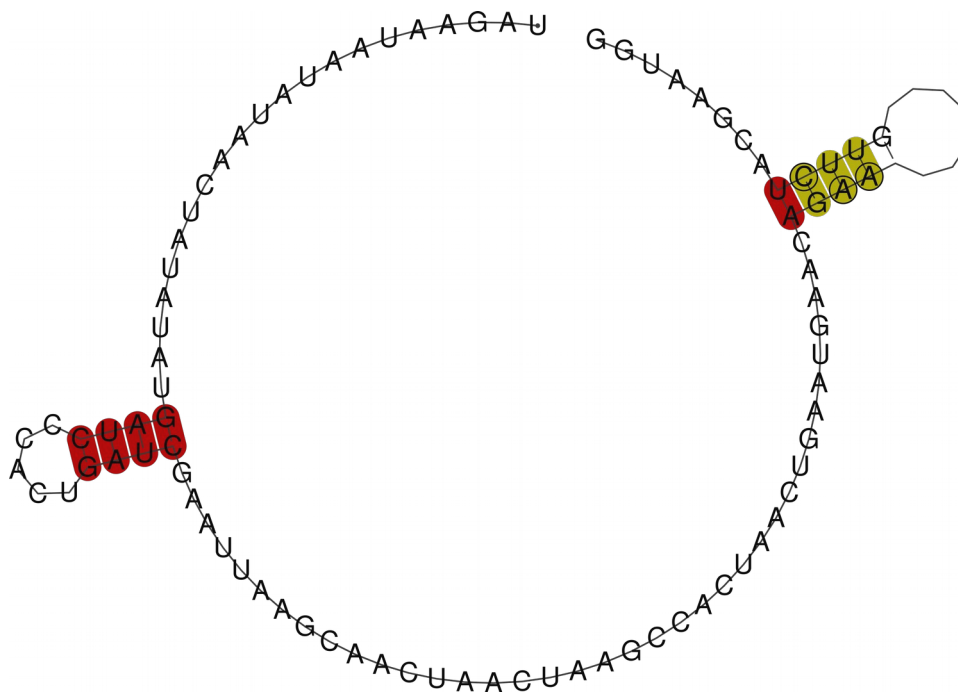
RNAz



100 % overlap

No. of identified lncRNA transcripts with conserved local structure

Species	Protein coding	RNAcode classified lncRNA	RNAz
<i>Oryza sativa</i>	42132	30779	14835
<i>Oryza rufipogon</i>	47441	33111	13690
<i>Oryza punctata</i>	41060	14828	3807
<i>Oryza nivara</i>	48360	33786	14531
<i>Oryza meridionalis</i>	43455	23971	9995
<i>Oryza longistaminata</i>	31686	13276	4257
<i>Oryza indica</i>	40745	24887	9102
<i>Oryza glumaepatula</i>	46893	31398	12618
<i>Oryza glaberrima</i>	33164	20073	6912
<i>Oryza brachyantha</i>	32037	9224	1258
<i>Oryza barthii</i>	41595	28263	10675
<i>Brassica rapa</i>	41024	9604	42
<i>Brassica oleracea</i>	59219	13690	46
<i>Brassica napus</i>	101040	16841	87
<i>Arabidopsis thaliana</i>	48321	4033	131
<i>Arabidopsis lyrata</i>	32663	2617	104



osativa.chr2
 omeridionalis.chr2
 opunctata.chr2
 oglumaepatula.chr2
 oglaberrima.chr2
 olongistaminata.scaKN540772

.....((((.....))).....
 TATAATAATATAACTATATATGATCCCACTGATCGAATTAAGCAACTAACTAAGCCACTA 60
 TAGAATAATATAACTATATATGATCCCACTGATCGAATTAAGCAACTAACTAAGCCACTA 60
 TAGAATAGTATAACTATAAATGATCCCACTGATCGAATTAAGCAACTAACTAAGCCACTA 60
 TATAATAATATAACTATATATGATCCCACTGATCGAATTAAGCAACTAACTAAGCCACTA 60
 TATAATAATATAACTATATATGATCCCACTGATCGAATTAAGCAACTAACTAAGCCACTA 60
 TAGAATAATATAACTATATATGATCCCACTGATCGAATTAAGCAACTAACTAAGCCACTA 60
10.....20.....30.....40.....50.....



osativa.chr2
 omeridionalis.chr2
 opunctata.chr2
 oglumaepatula.chr2
 oglaberrima.chr2
 olongistaminata.scaKN540772

.....((((.....))).....
 ACTGAATGAACAGGACACAAACAGTTCACGAATGG 96
 ACTGAATGAACAGAGCACAAACAGTTCACGAATGG 96
 ACTGAATGAACAGAACACAAACAGTTTAAAAATAG 96
 ACTGAATGAACAGGACACAAACAGTTCACGAATGG 96
 ACTGAATGAACAGGACACAAACAGTTCACGAATGG 96
 ACTGAATGAACAGAGCACAAACAGTTCACAAACTG 96
70.....80.....90.....



Total number of identified cisNAT pairs

Species	RNAcode based	Annotation based
<i>Oryza sativa</i>	970	445
<i>Oryza rufipogon</i>	293	1
<i>Oryza punctata</i>	329	0
<i>Oryza nivara</i>	384	0
<i>Oryza meridionalis</i>	622	18
<i>Oryza longistaminata</i>	0	0
<i>Oryza indica</i>	26	3
<i>Oryza glumaepatula</i>	420	0
<i>Oryza glaberrima</i>	7	7
<i>Oryza brachyantha</i>	346	1
<i>Oryza barthii</i>	162	1
<i>Brassica rapa</i>	4	0
<i>Brassica oleracea</i>	0	0
<i>Brassica napus</i>	0	0
<i>Arabidopsis thaliana</i>	839	552
<i>Arabidopsis lyrata</i>	6	0

Future work

- ✓ Identification of true protein coding and non-coding transcripts
- ✓ Genome wide screen for local conserved structures
- ✓ Identification of overlapping features
- Re-annotation of transcript biotypes considering the results from CPC, RNAcode and RNAz
- Enrichment analysis of local conserved structures in cisNATs
- Application of GraphClust / RNAscClust on the candidate cisNATs with conserved local structures for grouping them together for genome wide screen for similar regions
- Selection of candidate cisNATs for experimental testing
- Annotation of missing protein-coding genes

Acknowledgments



**Ivo Hofacker
Andrea Tanzer**

**Bernhard Thiel
Roman Ochsenreiter
Richard Neuboeck**

Uni Lausanne, Switzerland

**Yves Poirier
Jules Deforges
Rodrigo Siqueira Reis**

ETH Zurich, Switzerland

**Wilhelm Gruitsem
Katja Baerenfaller
Julia Svozil**

Funding



SWISS NATIONAL SCIENCE FOUNDATION

No. of annotated lncRNA transcripts classified as coding

Species	Annotated ncRNA	CPC Score ≥ 1.5	RNAcode Score>10; pvalue<0.05
Oryza sativa	55619	65	337
Oryza rufipogon	2020	12	34
Oryza punctata	28	0	0
Oryza nivara	24	0	0
Oryza meridionalis	2283	7	16
Oryza longistaminata	0	0	0
Oryza indica	2232	5	23
Oryza glumaepatula	19	0	0
Oryza glaberrima	1997	6	18
Oryza brachyantha	18	1	0
Oryza barthii	1681	3	15
Brassica rapa	1	0	0
Brassica oleracea	5	3	1
Brassica napus	0	0	0
Arabidopsis thaliana	3917	61	660
Arabidopsis lyrata	0	0	0