

生物信息学：导论与方法

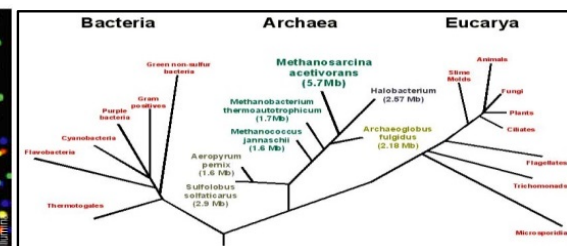
Bioinformatics: Introduction and Methods

Ge Gao 高歌 & Liping Wei 魏丽萍

Center for Bioinformatics, Peking University

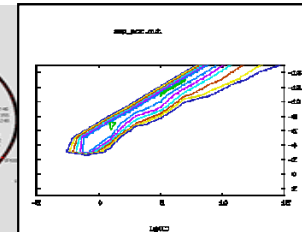


<https://www.coursera.org/course/pkubioinfo>



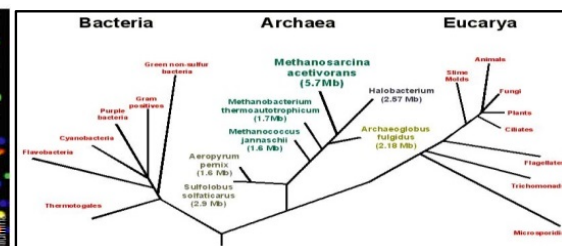
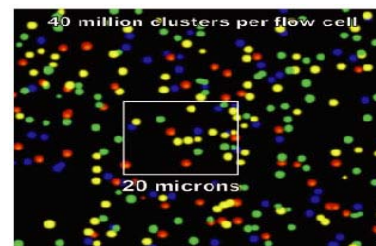
美国芝加哥大学生态与演化生物学系 龙漫远
北京大学生物信息学中心 魏丽萍

Manyuan Long, Department of Ecology and Evolutionary Biology, University of Chicago
Liping Wei, Center for Bioinformatics, Peking University





TAACCCTAACCCTAACCCTAACCCTAACCCTA
CCTAACCCTAACCCTAACCCTAACCCTAACC
CCCTAACCCTAACCCTAACCCTAACCCTAAC
AACCCTAACCCTAACCCTAACCCTAACCCTA
ACCCTAACCCTAACCCTAACCCTAACCCTAAC
CTACCCTAACCCTAACCCTAACCCTAACCCTA
ACCCTAACCCTAACCCTAACCCTAACCCTAAC

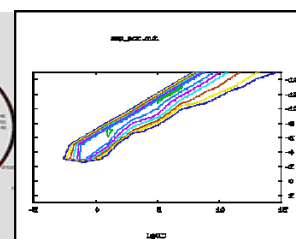
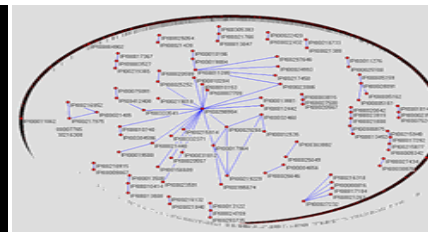
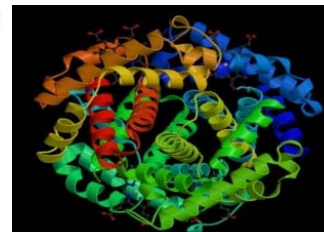
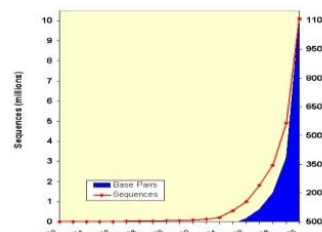
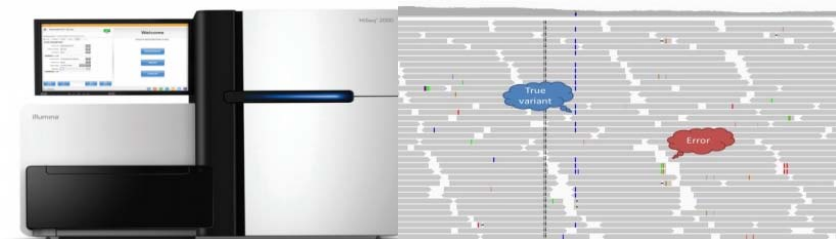


Unit 3: A Human-Specific *de novo* Gene Associated with Addiction

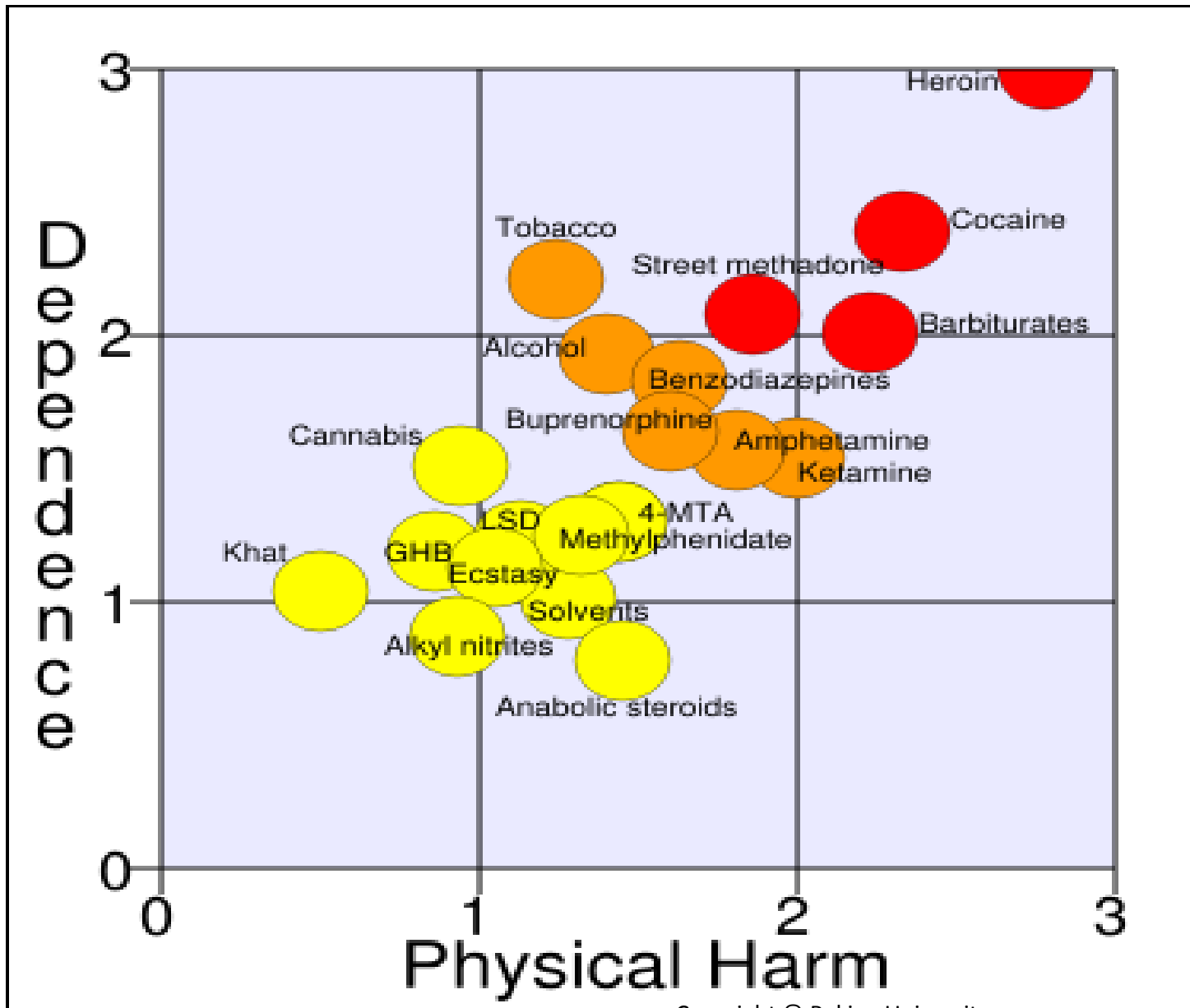
北京大学生物信息学中心 魏丽萍

Liping Wei, Ph.D.

Center for Bioinformatics, Peking University

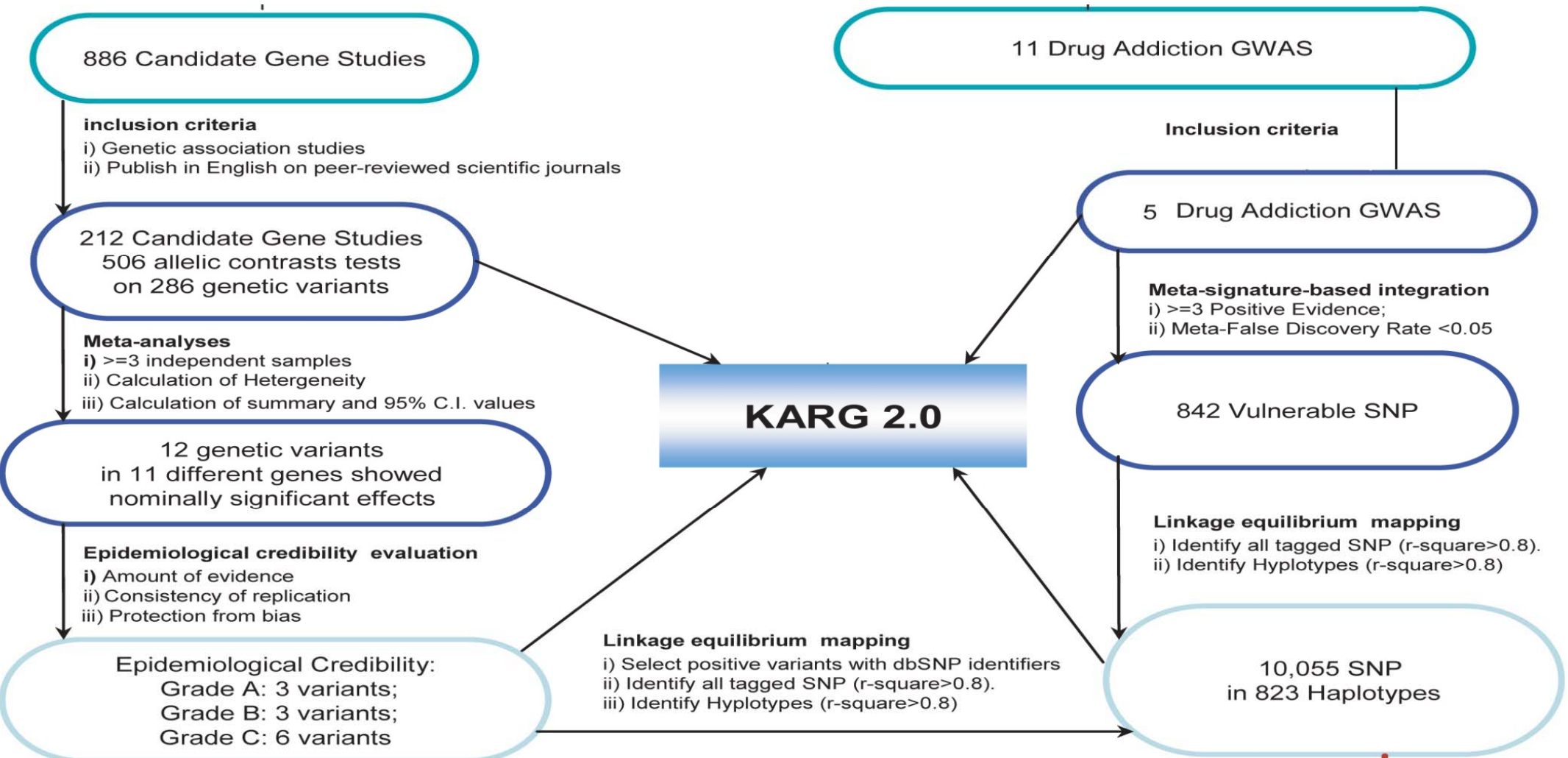


Addiction



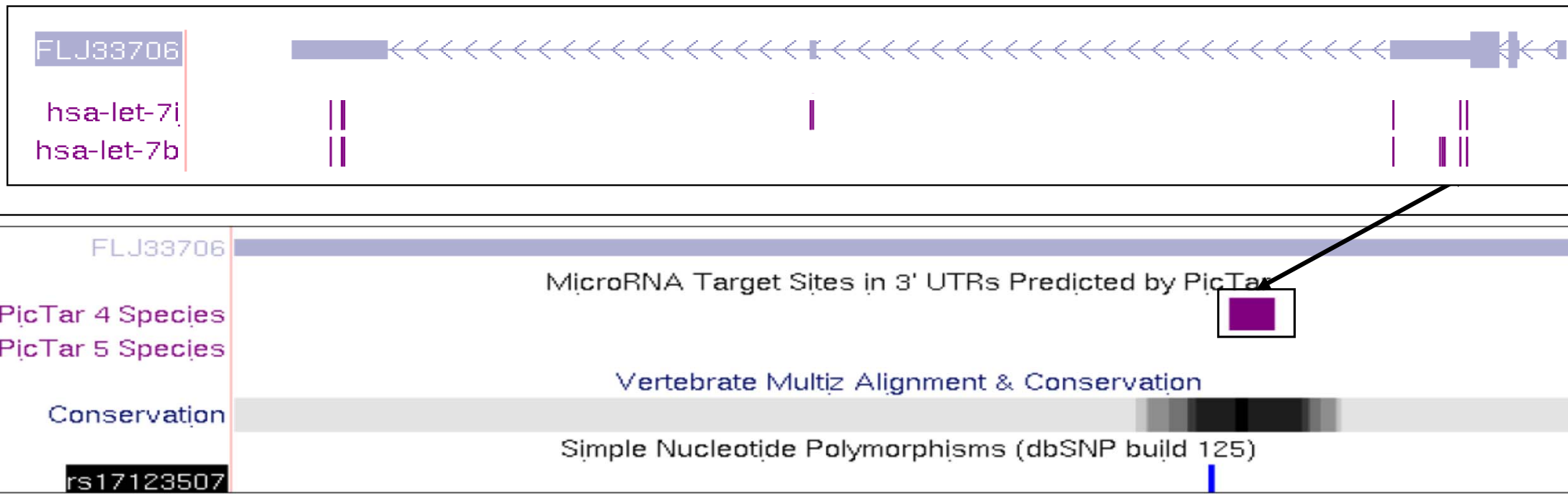
Nutt, *NEJM*, '07

Collaboration with NIH/NIDA to analyze addiction GWAS data

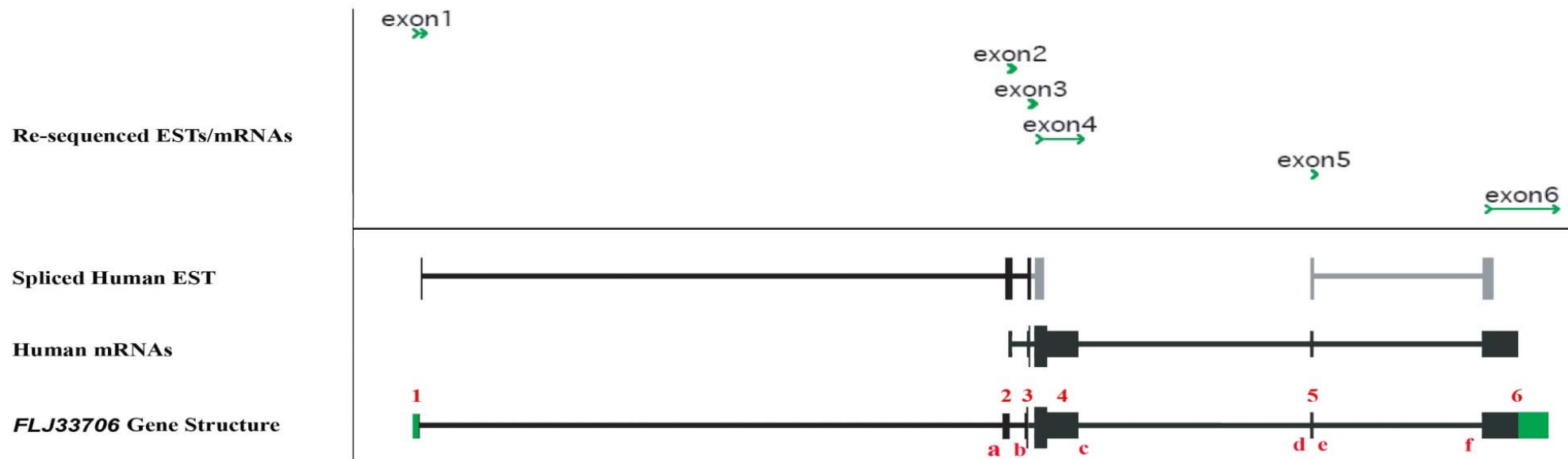


An SNP on the 3'UTR of *FLJ33706* is statistically significant in two GWAS of nicotine addiction and implicated in two linkage analyses.

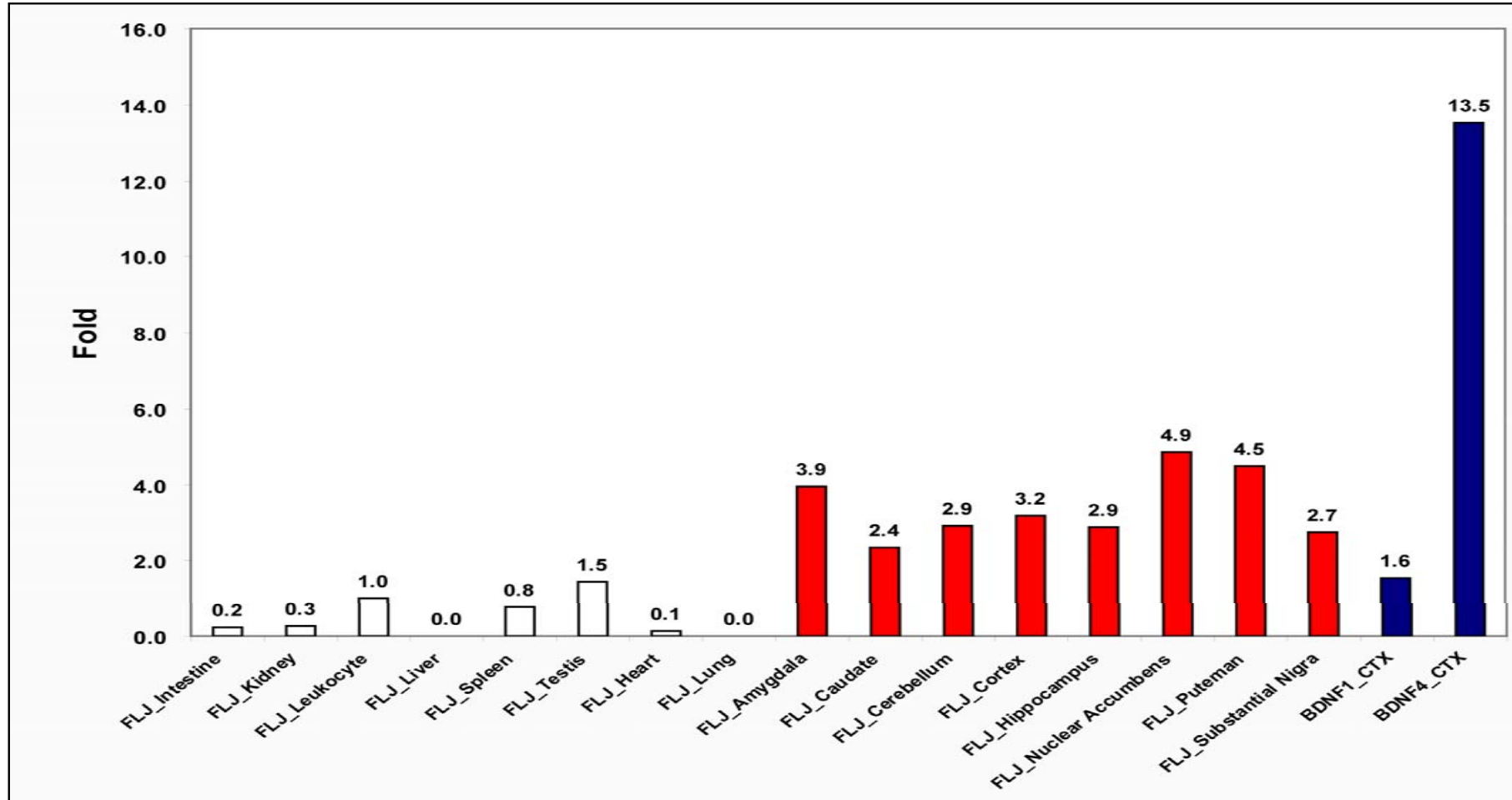
It is located in the middle of 12 binding sites of miRNA *let-7*.



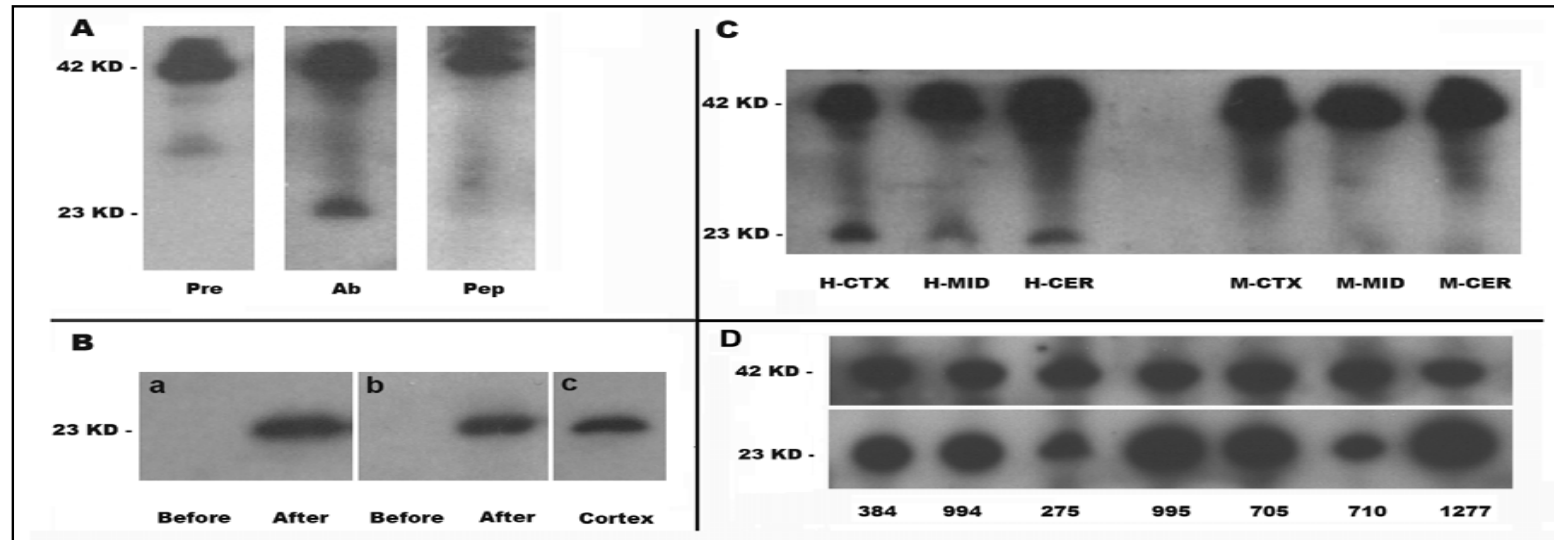
***FLJ33706* is a human-specific *de novo* protein-coding gene**



TaqMan-based Real-Time PCR showed that *FLJ33706* mRNA is enriched in human brain regions

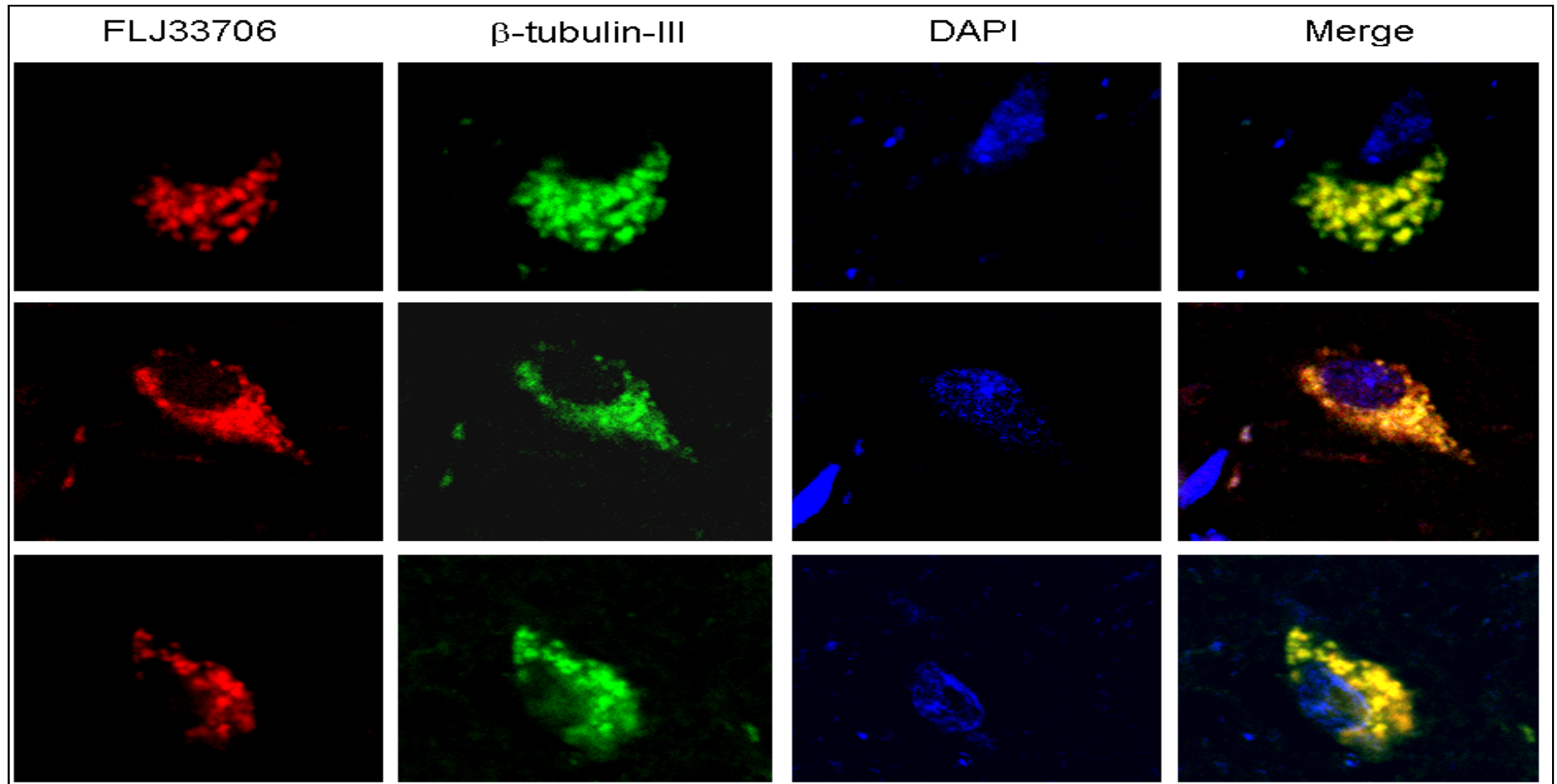


Western blot assay using an antibody designed against a 17-amino-acid peptide confirmed expression of FLJ33706 protein

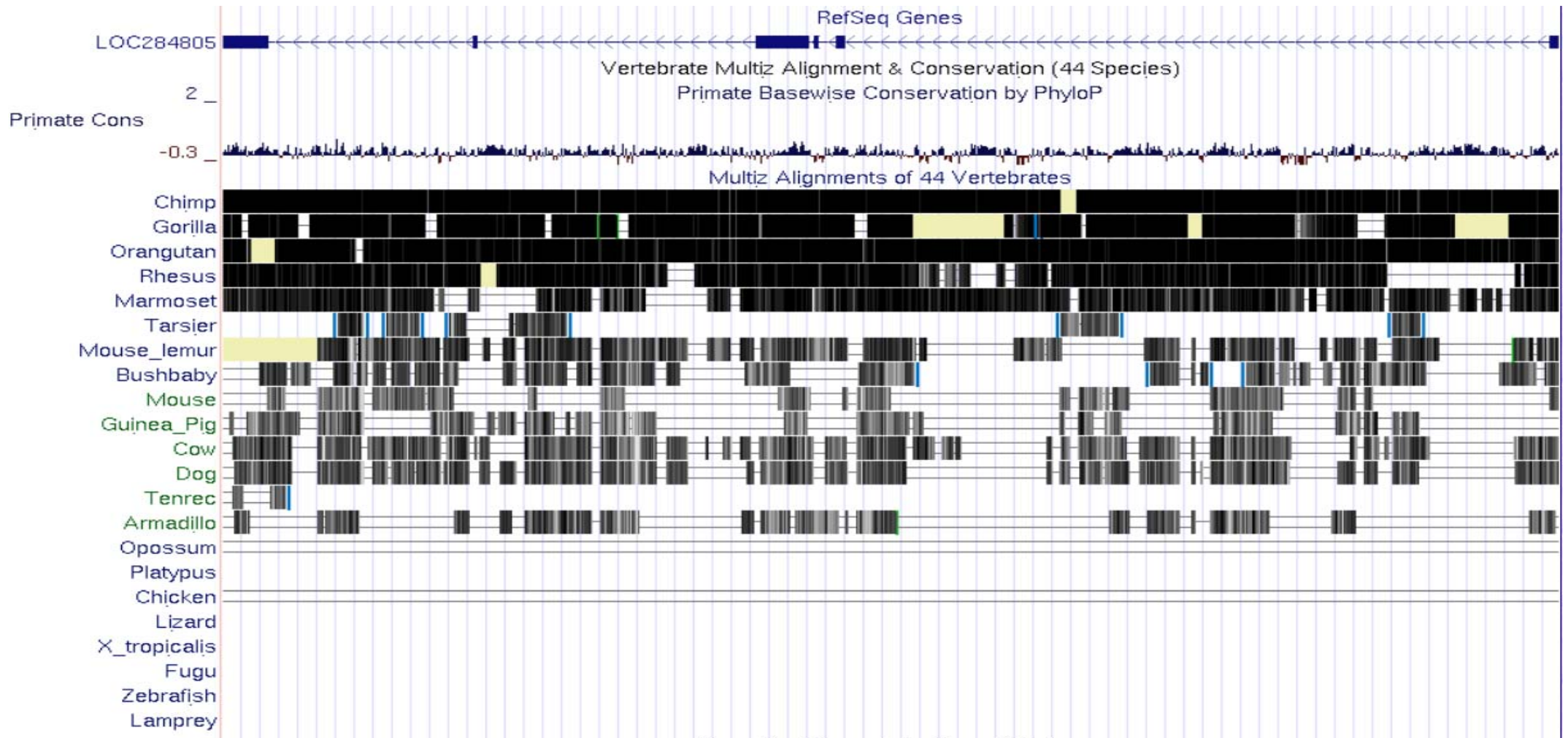


- (A) The band was not detected in pre-immune serum or in the presence of excess synthetic antigenic peptides
- (B) The band was detected only after transformation of FLJ33706 recombination plasmids in *E. coli* (a) His-tag specific antibody and (b) anti-FLJ33706.
- (C) The band was detected in human cortex, midbrain, and cerebellum, but not in mouse.
- (D) FLJ33706 expression can be detected in the cortex of seven different human individuals.

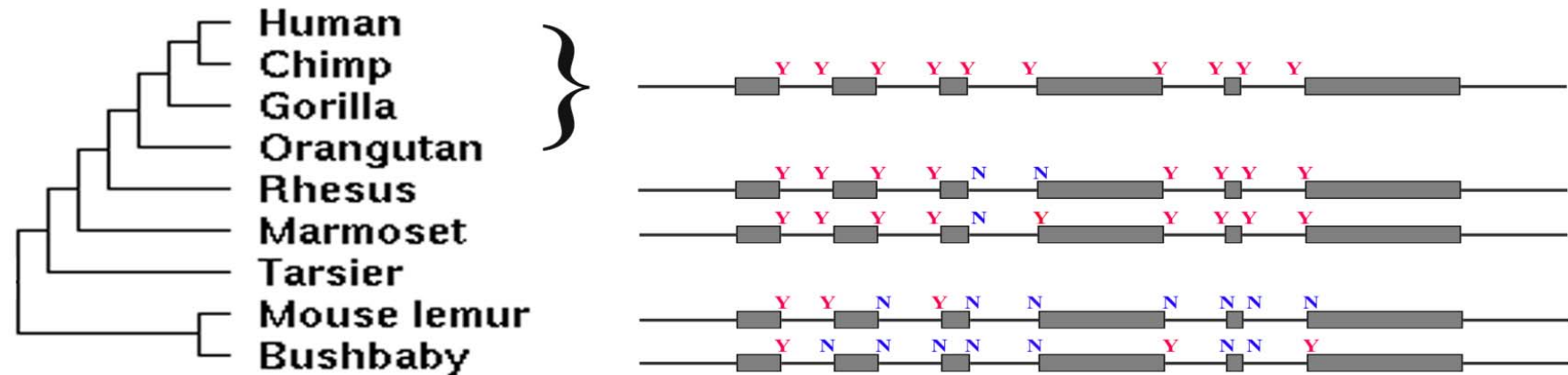
Immunohistochemistry studies of human cortex slides showed enrichment in cytoplasm of neuronal cells



The DNA segment emerged in eutherian mammals



Insertion of *Alu* elements generated splicing sites



		Intron 1						Intron 2						Intron 3						Intron 4						Intron 5							
Human	...CAG	G	T	ggg...tcc	A	G	ACT...GAG	G	T	aag...ttt	A	G	AGA...CGG	G	T	aag...aac	A	G	GCC...CTG	G	T	agg...gac	A	G	AGT...CAG	G	T	acg...tcc	A	G	ACT...		
Chimp	...CAG	G	T	ggg...tcc	A	G	ACT...GAG	G	T	aag...---	A	G	AGA...CCG	G	T	aag...aac	A	G	GTC...CTG	G	T	agg...gac	A	G	AGT...CAG	G	T	acg...tcc	A	G	ACT...		
Gorilla	...CAG	G	T	ggg...tcc	A	G	ACT...GAG	G	T	aag...---	A	G	AGA...CCG	G	T	aag...aac	A	G	GCC...CTG	G	T	agg...gac	A	G	AGT...CAG	G	T	aca...tcc	A	G	ACT...		
Orangutan	...CAG	G	T	ggg...tcc	A	G	ACT...GAG	G	T	aag...ttt	A	G	AGA...CCA	G	T	aag...aac	A	G	GCC...CTG	G	T	agc...gac	A	G	--T...CAG	G	T	acg...nnn	N	N	NNN...		
Rhesus	...CAG	G	T	ggg...tet	A	G	ACT...GAG	G	T	aag...---	A	G	AGA...CCG	G	A	aag...aat	A	A	GTC...CTG	G	T	agc...gac	A	G	AGT...TAG	G	T	acg...tcc	A	G	ACT...		
Marmoset	...CAG	G	T	gga...tcc	A	G	ACT...GAG	G	T	aag...---	A	G	AGA...CGG	C	T	aag...aac	A	G	GCC...CTG	G	T	agc...tag	G	G	AGT...CAG	G	T	acg...tcc	A	G	ACT...		
Mouse_lemur	...CAG	G	T	ggg...tte	A	G	ACT...GGG	G	T	aag...---	A	A	----CCA	G	A	gag...aac	T	G	CCA...---	-	-	---,===	=	=	===,===	=	=	===,nnn	N	N	NNN...		
Bushbaby	...CAG	A	G	agg...===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,CAA	G	T	agc...===	=	=	===,===	=	=	===,tcc	A	G	ATT...		
Mouse	...CAG	=	=	===...tcc	A	G	ACT...===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===...tg	G	A	GCC...		
Guinea_Pig	...===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===,===	=	=	===...cc	G	C	ACT...		
Cow	...GAG	-	C	tgg...tgc	A	G	ACC...GAG	G	T	cac...===	=	=	===,===	=	=	===...age	A	G	CCA...===	=	=	===,===	=	=	===,===	=	=	===...tgc	A	G	ATG...		
Dog	...CAG	-	T	egg...===	=	=	===...GAG	G	T	cac...===	=	=	===,===	=	=	===...a-c	A	G	CCA...---	-	-	---,---	-	-	---,---	-	-	---...tcc	A	G	ATG...		
Armadillo	...CAG	G	T	ggg...===	=	=	===,===	G	C	etg...---	=	=	---...AAG	G	A	agg...aac	A	G	CCA...---	-	-	---,===	=	=	===,===	=	=	===,===	=	=	===,---		

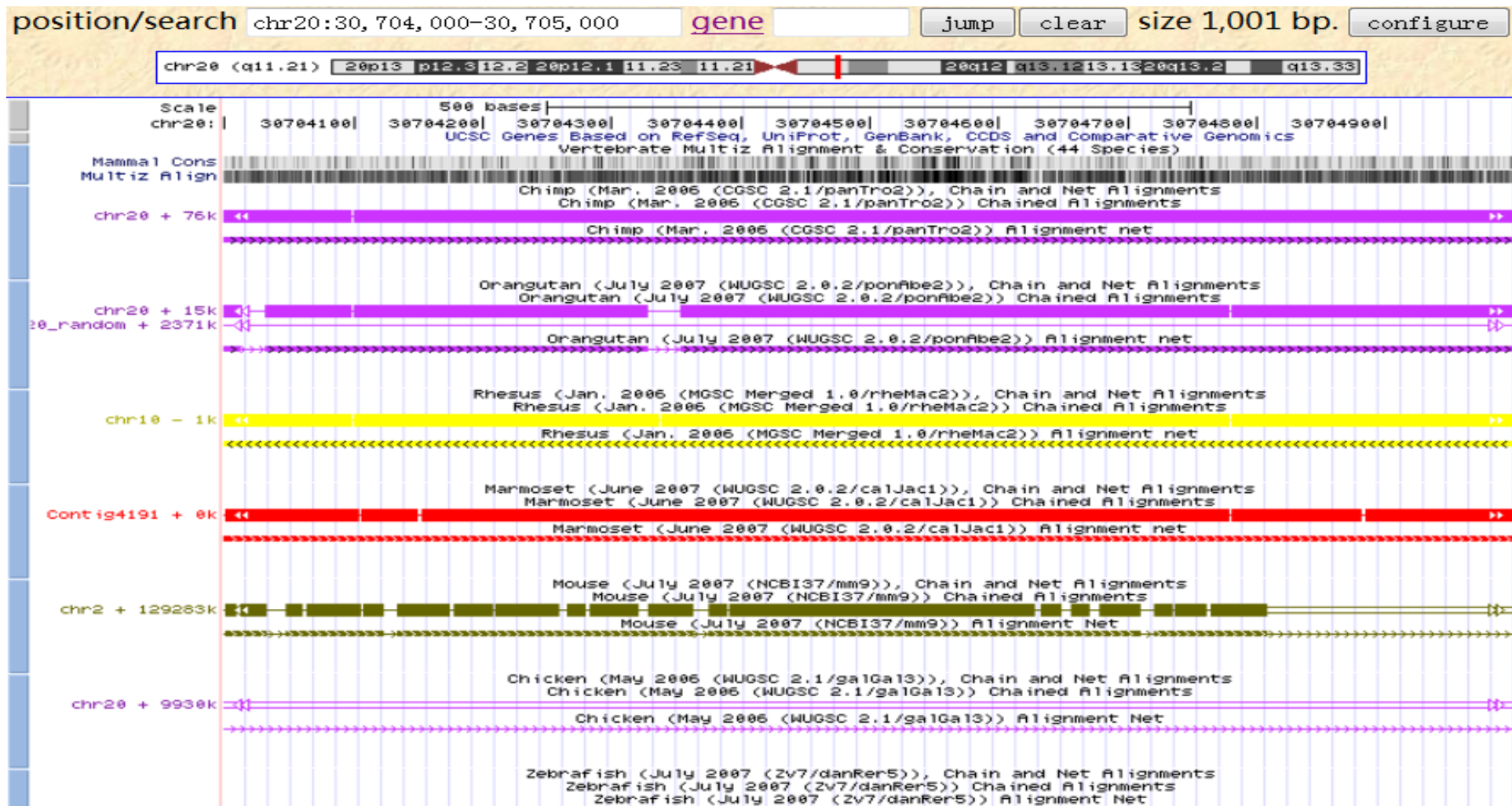
Two changes in human escaped two stop codons

	<u>M</u>	<u>V</u>	<u>R</u>	<u>A</u>	<u>I</u>	<u>N</u>	<u>D</u>	<u>W</u>	<u>R</u>	<u>F</u>	<u>K</u>	<u>G</u>	<u>L</u>		
Human	A T G	G T C	C G G	G C G A T T	A A C	G A T	T G G	C G C	T T T	A A A	G G A	C T G			
Chimp	.	.	A	.	.	.	A			
Gorilla	.	.	A	.	.	.	A			
Orangutan	A	.	G	.	AA	.			
Rhesus	G	C	A	T	G	G	A	T	.	.	G	T			
Position	1	6	10	13	14	21	24	28	31	35	39	41	43		
	<u>R</u>	<u>A</u>	<u>T</u>	<u>V</u>	<u>A</u>	<u>G</u>	<u>L</u>	<u>G</u>	<u>R</u>	<u>A</u>	<u>P</u>	<u>Q</u>	<u>R</u>	<u>P</u>	
Human	C G G G C C A C A	G T C	G C T G G A	C T T G G C	G C G	A G G	G C T	C C C C A G	C G C	C C T					
Chimp	C	T		
Gorilla	C	G	.	.	A	T	.	.	.		
Orangutan	C A	T	A	.		
Rhesus	C	T	G	C	A	.	T	T	.	.	T	A	T		
Position	45	46	47	49	51	52	60	61	64	66	71	76	77	82	84
	<u>P</u>	<u>W</u>	<u>E</u>	<u>V</u>	<u>L</u>	<u>L</u>	<u>S</u>	<u>R</u>	<u>R</u>	<u>R</u>	<u>M</u>	<u>T</u>	<u>V</u>	<u>D</u>	
Human	C C T ... T G G	G - - A A	G T T	C T C	C T C A G C	C G G	C G G	A G G	A T G	A C G G T G	G A C				
Chimp	.	C	G	C	C	.	.	T		
Gorilla	.	C A	G G	C	C		
Orangutan	T	C	G	C	C	G	T	.	A	A	T T	C	.		
Rhesus	.	C A	G	.	C	T G	.	T	T G		
Position	92	104	106	107	110	112	113	127	132	134	139	144	145	147	
	<u>L</u>	<u>S</u>	<u>L</u>	<u>T</u>	<u>C</u>	<u>F</u>	<u>L</u>	<u>Q</u>	<u>S</u>	<u>N</u>	<u>R</u>	<u>STOP</u>			
Human	C T G	T C G C T G	A C C	T G T	T T C	C T C	C A G	T C C A A T	C G G	T A G					
Chimp	.	T	.	.	.	C	.	.	G	.	.	.			
Gorilla	.	A	C			
Orangutan	.	A	.	C			
Rhesus	T	.	C	.	G	.	.	C	G	.	.	.			
Position	152	154	155	158	161	165	167	183	186	187	190	195			

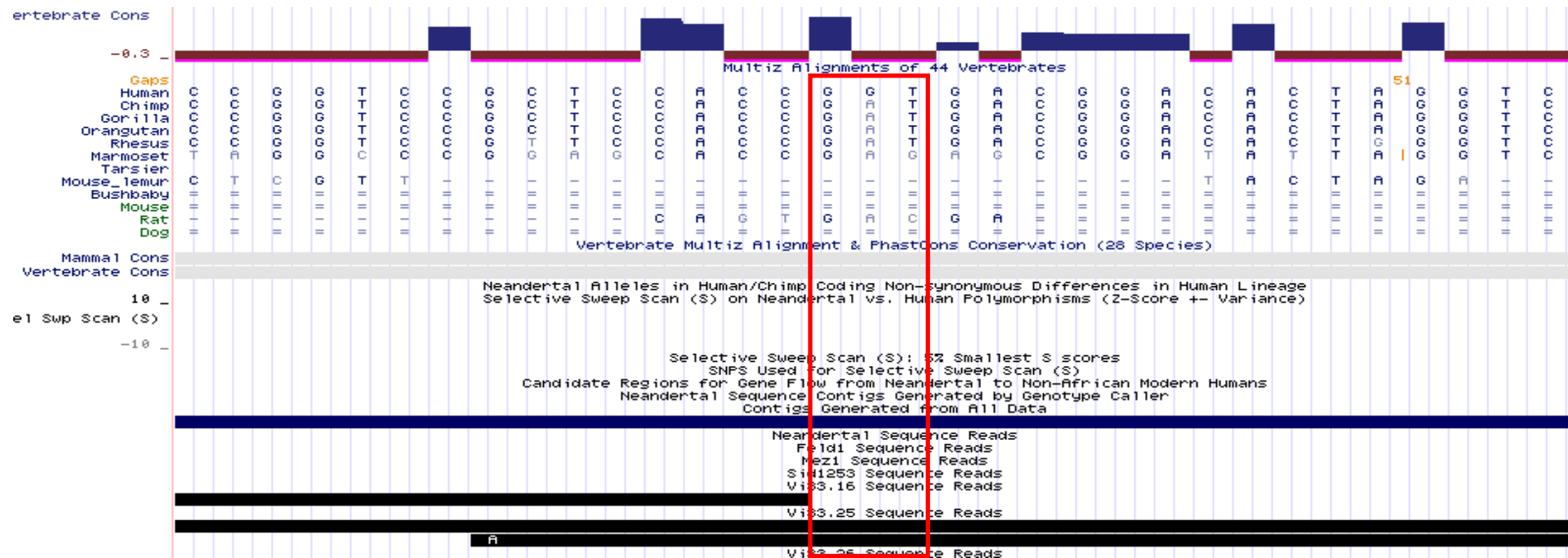
There are signals of enhancer and transcription factor binding sites in the 5kb upstream regions



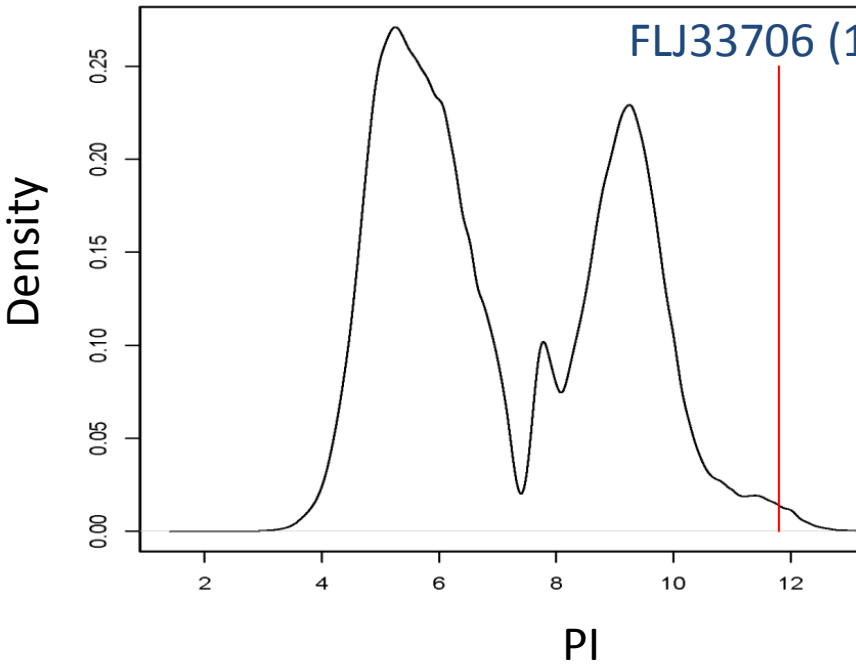
Promoter region is absent in chicken/zebrafish,
emerged in mouse, and is similar in rhesus and chimpanzee



The Open Reading Frame is intact in Neadertal genome



FLJ33706 has high PI



GO Term	FDR q-value
RNA binding	5.50E-08
cytosolic ribosome	3.68E-07
macromolecular complex	1.63E-06
cytosolic large ribosomal subunit	4.61E-05
RNA splicing	6.71E-05
cytosolic part	7.73E-05
ribosomal subunit	4.54E-04
large ribosomal subunit	7.89E-04
intracellular organelle part	9.99E-04
organelle part	0.001136772
ribonucleoprotein complex	0.003187642
cellular biosynthetic process	0.007101674
MHC class II receptor activity	0.009220135
translation	0.010595406
mRNA processing	0.012153244
RNA processing	0.012167141
structural constituent of ribosome	0.017365179
mRNA metabolic process	0.020473341
macromolecule metabolic process	0.021017467
intracellular non-membrane-bound organelle	0.024935299
non-membrane-bound organelle	0.024935299
ribosome	0.036638186

Unpublished

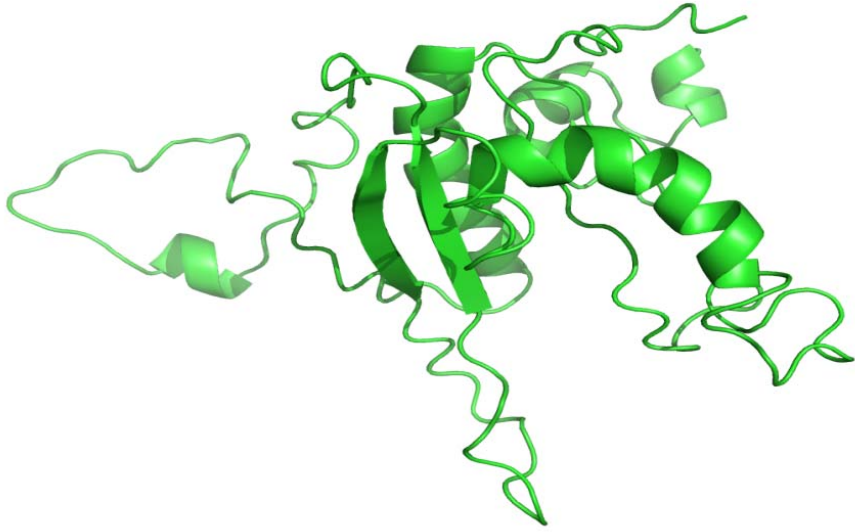
GO enrichment of proteins with PI > 11 (FDR < 0.05)

Predicted Secondary Structure include four helices & one beta strand

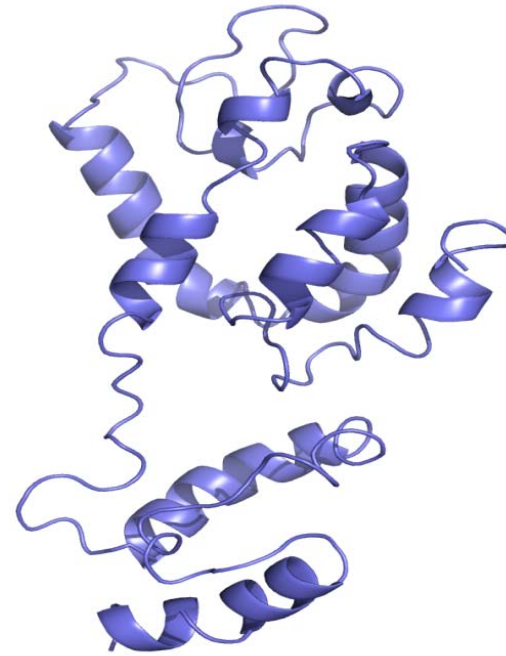


Unpublished

Predicted 3D Structure (probably not reliable)



Scored as best by I-TASSER



Another putative conformation

How many other human-specific *de novo* genes are there?

Where did they originate from?

生物信息学：导论与方法

Bioinformatics: Introduction and Methods

Ge Gao 高歌 & Liping Wei 魏丽萍

Center for Bioinformatics, Peking University



<https://www.coursera.org/course/pkubioinfo>