

SUPPLEMENTAL FIGURES AND TABLES

Identification of LINE retrotransposons and long non-coding RNAs expressed in the octopus brain

Giuseppe Petrosino^{1,*}, Giovanna Ponte^{1,*}, Massimiliano Volpe^{1,*}, Ilaria Zarrella¹, Federico Ansaloni^{8,2}, Concetta Langella¹, Giulia Di Cristina¹, Sara Finaurini², Monia T. Russo³, Swaraj Basu¹, Francesco Musacchia¹, Filomena Ristoratore¹, Dinko Pavlinic⁴, Vladimir Benes⁴, Maria I. Ferrante³, Caroline Albertin⁵, Oleg Simakov^{6,7}, Stefano Gustincich^{8,2}, Graziano Fiorito^{1,#}, Remo Sanges^{1,2,8,#}

¹ Department of Biology and Evolution of Marine Organisms, Stazione Zoologica Anton Dohrn, Villa Comunale, 80121 Napoli (SZN), Italy.

² Neurobiology Sector, Scuola Internazionale Superiore di Studi Avanzati (SISSA), Via Bonomea 265, 34136 Trieste, Italy.

³ Department of Integrative Marine Ecology, Stazione Zoologica Anton Dohrn, Villa Comunale, 80121 Napoli (SZN), Italy.

⁴ Scientific Core Facilities & Technologies, GeneCore, European Molecular Biology Laboratory (EMBL), Meyerhofstrasse 1, 69117 Heidelberg, Germany.

⁵ Marine Biological Laboratory (MBL), Woods Hole, Massachusetts, USA.

⁶ Okinawa Institute of Science and Technology Graduate University, Onna, Okinawa 9040495, Japan.

⁷ Department of Molecular Evolution and Development, Wien University, Althanstraße 14 (UZA I), 1090 Wien, Austria.

⁸ Central RNA Laboratory, Istituto Italiano di Tecnologia (IIT), Via Enrico Melen 83, 16152 Genova, Italy.

* These authors contributed equally to this work

corresponding authors:

Remo Sanges, e-mail: remo.sanges@gmail.com

Graziano Fiorito, e-mail: graziano.fiorito@szn.it

Table S1: *Octopus vulgaris* samples utilized for sequencing. Brain: supra-oesophageal (SEM) and sub-oesophageal (SUB) masses, and optic lobe (OL); arm: only anterior second left arm (ARM).

Animal	Sex	Weight	Part	Partname	Samplename
13/01	F	296	supra-oesophageal mass	SEM	IZ10
13/01	F	296	sub-oesophageal mass	SUB	IZ11
13/01	F	296	optic lobe	OL	IZ12
13/01	F	296	anterior arm	ARM	IZ13
13/02	M	336	supra-oesophageal mass	SEM	IZ15
13/02	M	336	sub-oesophageal mass	SUB	IZ16
13/02	M	336	optic lobe	OL	IZ17
13/02	M	336	anterior arm	ARM	IZ18
13/03	M	288	supra-oesophageal mass	SEM	IZ20
13/03	M	288	sub-oesophageal mass	SUB	IZ21
13/03	M	288	optic lobe	OL	IZ22
13/03	M	288	anterior arm	ARM	IZ23

Table S2: Counts of the assembled *Octopus vulgaris* and *O. bimaculoides* transcriptomes.

	<i>O. vulgaris</i>	<i>O. bimaculoides</i>
Assembled and filtered transcripts	64477	92820
Total bases	84399088	102485827
GC content (%)	37.9	37
Contig N50	2087	1573
Median transcript length (bp)	795	744
Average transcript length (bp)	1308	1104
Minimum length (bp)	201	201
Maximum length (bp)	20031	32440
248 CEGs Complete (%)	97.20	98.79
248 CEGs Complete + Partial (%)	98.4	100

Table S3: Counts of the repeats composition for the assembled *Octopus vulgaris* transcriptome.

	Nucleotides	Percentage	Transcripts	Percentage
Bases Masked	6584938	7.8	46944	72.8
Retroelements	2102784	2.5	16926	26.2
DNA transposons	1501995	1.8	11913	18.5
Unclassified	169576	0.2	250	0.4
Total interspersed repeats	3774355	4.5	22915	35.5
Satellites	150431	0.2	1326	2.1
Simple repeats	2368596	2.8	34833	54
Low complexity	470441	0.6	7704	11.9

Table S4: LINEs used for phylogenetic analysis (see also: Ohshima and Okada³⁶ and Jurka et al.⁷¹).

Clade	Name	Species	Sequence ID
CRE	SLACS	<i>Trypanosoma brucei</i>	CAA34931
	CZAR	<i>Trypanosoma cruzi</i>	AAA30239
	CRE1	<i>Crithidia fasciculata</i>	AAA75435
	CRE2	<i>Crithidia fasciculata</i>	AAB40036
Genie/Gil	GilM	<i>Giardia intestinalis</i>	AAI47180
R4	R4A1	<i>Ascaris lumbricoides</i>	AAA97394
R1	R1Dm	<i>Drosophila melanogaster</i>	CAA36227
	R1	<i>Bradysia coprophila</i>	AAA29813
	RT1	<i>Anopheles gambiae</i>	AAA29363
	RT2	<i>Anopheles gambiae</i>	AAA29365
	R1Bm	<i>Bombyx mori</i>	AAC13649
	TRAS1	<i>Bombyx mori</i>	BAA07467
	SART1	<i>Bombyx mori</i>	BAA19776
	LOA	LOA	<i>Drosophila silvestris</i>
LOA	BAGGINS1	<i>Drosophila melanogaster</i>	Repbase
	Bilbo	<i>Drosophila subobscura</i>	AAB92389
	Lian-Aa1	<i>Aedes aegypti</i>	AAB65093
Tad1	Tad1	<i>Neurospora crassa</i>	AAA21781
	MgL	<i>Magnaporthe grisea</i>	AAB71689
	CgT1	<i>Glomerella cingulata</i>	AAA85636
Jockey	BMC1	<i>Bombyx mori</i>	BAB21761
	amy	<i>Bombyx mori</i>	AAA17752
	Juan-A	<i>Aedes aegypti</i>	AAA29354
	Juan-C	<i>Culex pipiens</i>	AAA28291
	NLR1Cth	<i>Chironomus thummi</i>	AAB26437
	Doc6	<i>Drosophila melanogaster</i>	Repbase
	G5	<i>Drosophila melanogaster</i>	Repbase
	G5A	<i>Drosophila melanogaster</i>	Repbase
	Jockey	<i>Drosophila melanogaster</i>	AAA28675
	Doc	<i>Drosophila melanogaster</i>	CAA35587
	Fw	<i>Drosophila melanogaster</i>	AAA28508
	Fw2	<i>Drosophila melanogaster</i>	Repbase
	G4	<i>Drosophila melanogaster</i>	Repbase
	Helena	<i>Drosophila yakuba</i>	AAC24972
	BS	<i>Drosophila melanogaster</i>	Repbase
BS3	<i>Drosophila melanogaster</i>	Repbase	
TART-B1	<i>Drosophila melanogaster</i>	AAC46494	
I	I-1_DR	<i>Danio rerio</i>	Repbase
	IVK	<i>Drosophila melanogaster</i>	Repbase
	I	<i>Drosophila melanogaster</i>	AAA70222
	ingi	<i>Trypanosoma brucei</i>	CAA29181
	L1Te	<i>Trypanosoma cruzi</i>	CAB41693
Rex1	Rex1-1_DR	<i>Danio rerio</i>	Repbase
CR1	Sam3	<i>Caenorhabditis elegans</i>	AAA93347
	Sam1	<i>Caenorhabditis elegans</i>	AAA21080
	Q	<i>Anopheles gambiae</i>	AAA53489
	CR1-3_AG	<i>Anopheles gambiae</i>	Repbase
	CR1-5_AG	<i>Anopheles gambiae</i>	Repbase
	T1	<i>Anopheles gambiae</i>	AAA29367
	CR1-2_AG	<i>Anopheles gambiae</i>	Repbase
	CR1-4_AG	<i>Anopheles gambiae</i>	Repbase
	DMCR1A	<i>Drosophila melanogaster</i>	Repbase
	CR1	<i>Gallus gallus</i>	AAC60281
	PsCR1	<i>Platymys spixii</i>	BAA88337
	L3	<i>Homo sapiens</i>	Repbase
	SR1	<i>Schistosoma mansoni</i>	AAC06263
	L2	UnaL2	<i>Anguilla japonica</i>
Maui		<i>Takifugu rubripes</i>	AAD19348
CR1-3_DR		<i>Danio rerio</i>	Repbase
CR1-1_AG		<i>Anopheles gambiae</i>	Repbase
R2	R2Bm	<i>Bombyx mori</i>	AAB59214

Clade	Name	Species	Sequence ID
	R2Nv	<i>Nasonia vitripennis</i>	AAC34927
	R2Fa	<i>Forficula auricularia</i>	AAC34906
	R2Dm	<i>Drosophila melanogaster</i>	CAA36225
	R2Am	<i>Anurida maritima</i>	AAC34903
	R2Lp	<i>Limulus polyphemus</i>	AAC34904
NeSL-1	NeSL-1	<i>Caenorhabditis elegans</i>	CAB04870
RTE	BovB	<i>Bos taurus</i>	Repbase
	BovB_VA	<i>Vipera ammodytes</i>	Repbase
	RTE-1	<i>Caenorhabditis elegans</i>	AAA50641
	RTE-2	<i>Caenorhabditis elegans</i>	AAB00700
	RTE-1_AG	<i>Anopheles gambiae</i>	Repbase
	JAM1	<i>Aedes aegypti</i>	Repbase
	Rex3	<i>Xiphophorus maculatus</i>	Repbase
	SR2	<i>Schistosoma mansoni</i>	Repbase
	RTE-3_AG	<i>Anopheles gambiae</i>	Repbase
	RTE-2_CPB	<i>Chrysemys picta bellii</i>	Repbase
	RTE-4_AMi	<i>Crocodylidae</i>	Repbase
	RTE-5_AMi	<i>Crocodylidae</i>	Repbase
	RTE-6_AMi	<i>Crocodylidae</i>	Repbase
	RTE-7_AMi	<i>Crocodylidae</i>	Repbase
	RTE-8_AMi	<i>Crocodylidae</i>	Repbase
L1	L1Hs	<i>Homo sapiens</i>	AAA51622
	L1Md	<i>Mus musculus</i>	AAA66024
	L1-1_DR	<i>Danio rerio</i>	Repbase
	L1-10_DR	<i>Danio rerio</i>	Repbase
	L1-6_DR	<i>Danio rerio</i>	Repbase
	L1-8_DR	<i>Danio rerio</i>	Repbase
	L1-3_DR	<i>Danio rerio</i>	Repbase
	L1-5_DR	<i>Danio rerio</i>	Repbase
	L1-2_DR	<i>Danio rerio</i>	Repbase
	L1-4_DR	<i>Danio rerio</i>	Repbase
	DRE	<i>Dictyostelium discoideum</i>	Repbase
	ATLINE1_4	<i>Arabidopsis thaliana</i>	Repbase
	ATLINE1_5	<i>Arabidopsis thaliana</i>	Repbase
	Tal1-1	<i>Arabidopsis thaliana</i>	AAA75254
	ATLINE1_6	<i>Arabidopsis thaliana</i>	Repbase
	Cin4	<i>Zea mays</i>	Repbase
	Tx1	<i>Xenopus laevis</i>	AAA49976
	Zepp	<i>Chlorella vulgaris</i>	BAA25763

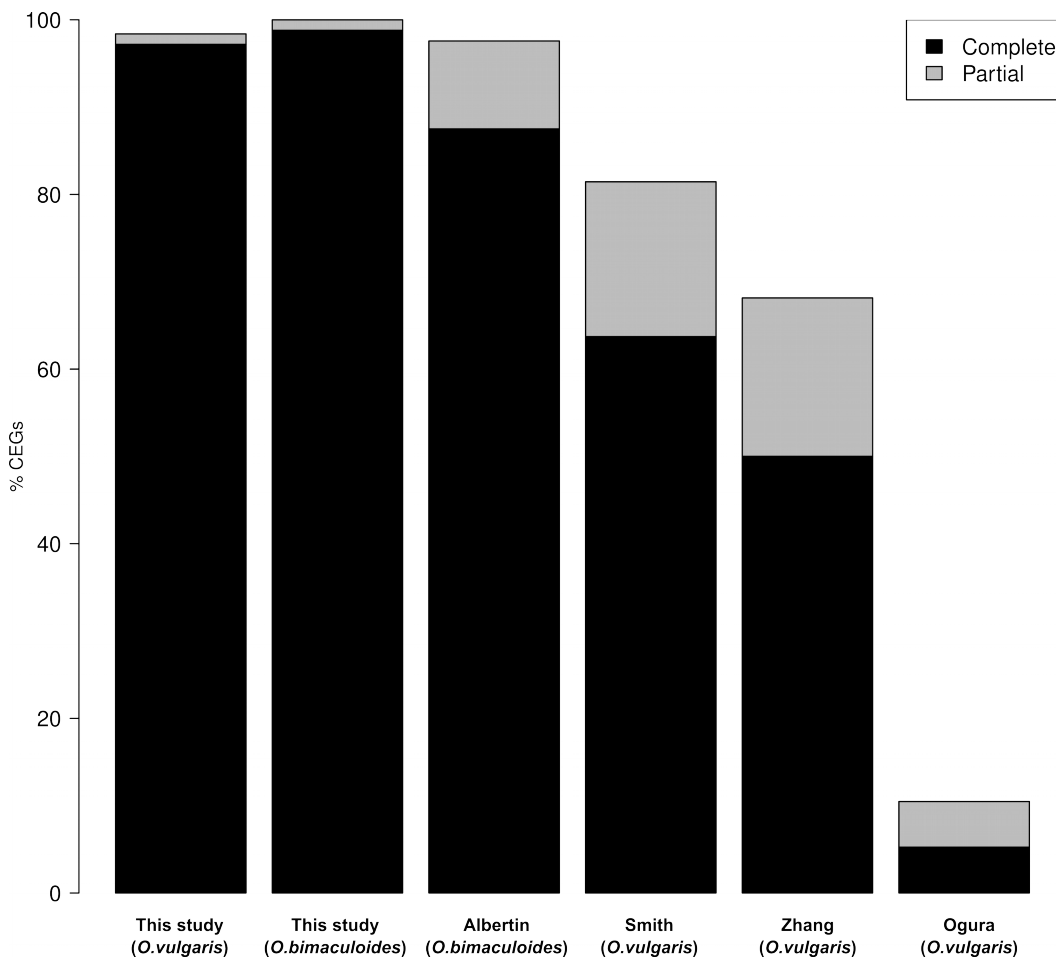


Figure S1: Completeness of *Octopus* transcriptomes. Percentages of core eukaryotic genes (CEGs) represented into this and published assemblies of octopus transcriptomes (Albertin et al.⁵; Smith et al.⁶²; Zhang et al.⁸⁴; Ogura et al.⁸⁵). The barplots indicate the percentages of CEGs present into every published transcriptome. The portion defined as “complete” identify all those transcripts whose assembled sequence are predicted to be full-length while “partial” indicates the reconstruction of only a fragment of specific CEGs. All the transcriptomes originate from *Octopus vulgaris* with the exception of the transcriptome from Albertin et al.⁵ which originate from *Octopus bimaculoides*. Two transcriptomes for *O. bimaculoides* are included here, the original one assembled by Albertin et al.⁵ and the same assembled for the aims of this study as described in Methods. The transcriptomes assembled in this work should be considered the most complete available to date for the genus *Octopus*, since they contain the highest percentages of core eukaryotic genes according to CEGMA (98% for *O. vulgaris* and 100% for *O. bimaculoides*).

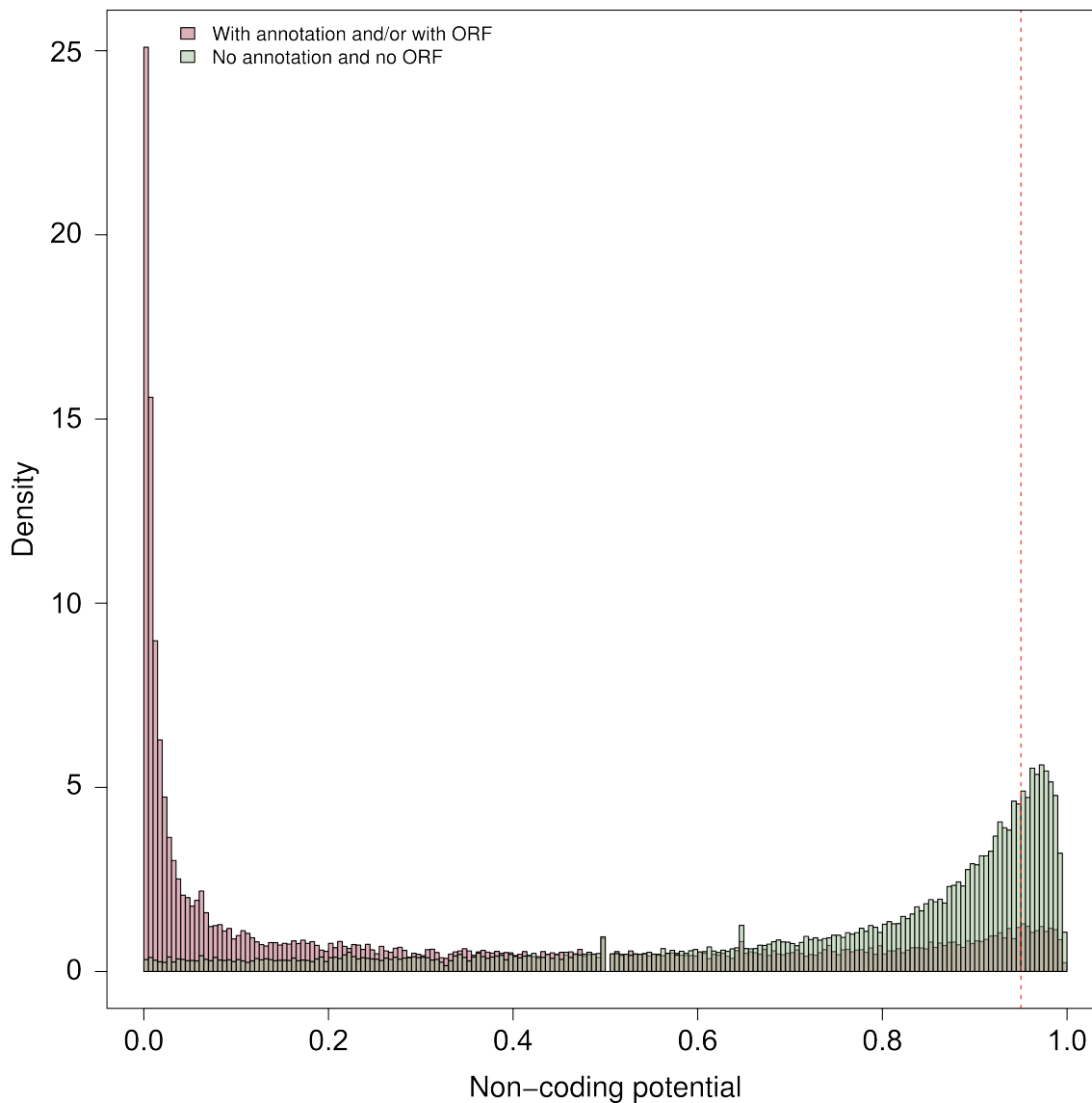


Figure S2: Stringent classification of lncRNAs. Non-coding potential score plot as measured by Portrait⁵⁷ for the assembled transcripts. The transcriptome has been divided in two groups: the first group of genes showing at least one BLAST match against a protein or a domain or a ribosomal or a small RNA in the annotation analysis and/or an ORF bigger than 100 aa (light red); and the second group of genes not showing any BLAST match and whose longest ORF results shorter than 100 aa (light green). Only transcripts without any match, with an ORF smaller than 100 aa and a non-coding potential bigger than 0.95 have been classified as non-coding. The vertical red dotted line represents the 0.95 non-coding potential cut-off used. The transcripts classified as non-coding are those plotted in the green bars at the right of the vertical red dotted line. According to Portrait recommendations, a non-coding potential score bigger than 0.5 is sufficient to classify a transcript as non-coding. We applied more stringent conditions and despite the combination of multiple parameters and the application of these stringent cutoffs we were still able to discover that a high proportion of transcripts likely represent lncRNAs. Here an underestimation of the true proportion of lncRNAs is possible because the RNAseq was not conducted using a strand specific protocol and therefore there is the possibility that annotation for “coding” has been utilized also for cases of “antisense non-coding” overlapping sense coding regions.

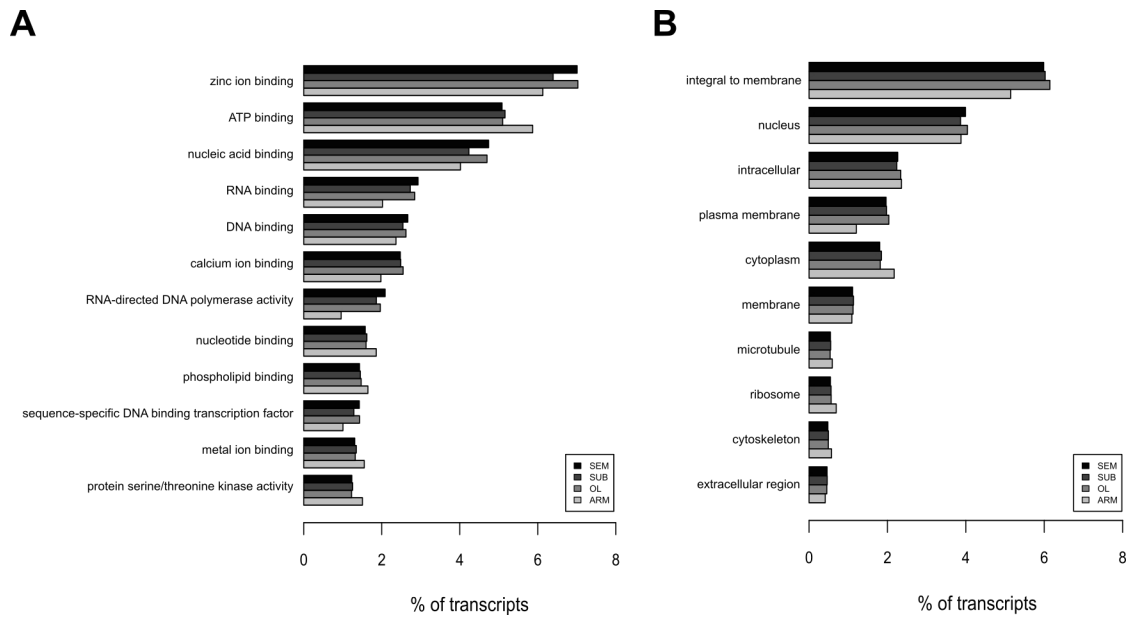


Figure S3: Most represented GO Molecular Function and Cellular Component classes. Barplots showing the percentage of top represented GO classes in the transcripts expressed for every tissue considered. The top 10 represented classes for every tissue were selected and the percentage of transcripts expressed associated to the given class is reported. **a**, Molecular function classification confirms the findings obtained with the biological process classification showing a higher rate of RNA-directed DNA polymerase activity in samples deriving from the brain. It is also important to underline that the top represented class in all the parts is the zinc ion binding which confirm the expansion of zinc fingers protein in octopus⁵. **b**, Transcriptome classification according to cellular component division results to be generally similar among the sampled parts. The octopus brain (SEM, SUB, OL) appears to contain a higher number of transcripts whose protein product is localized to the plasma membrane; higher representation of transcripts localized into the cytoplasm are found in ARM.

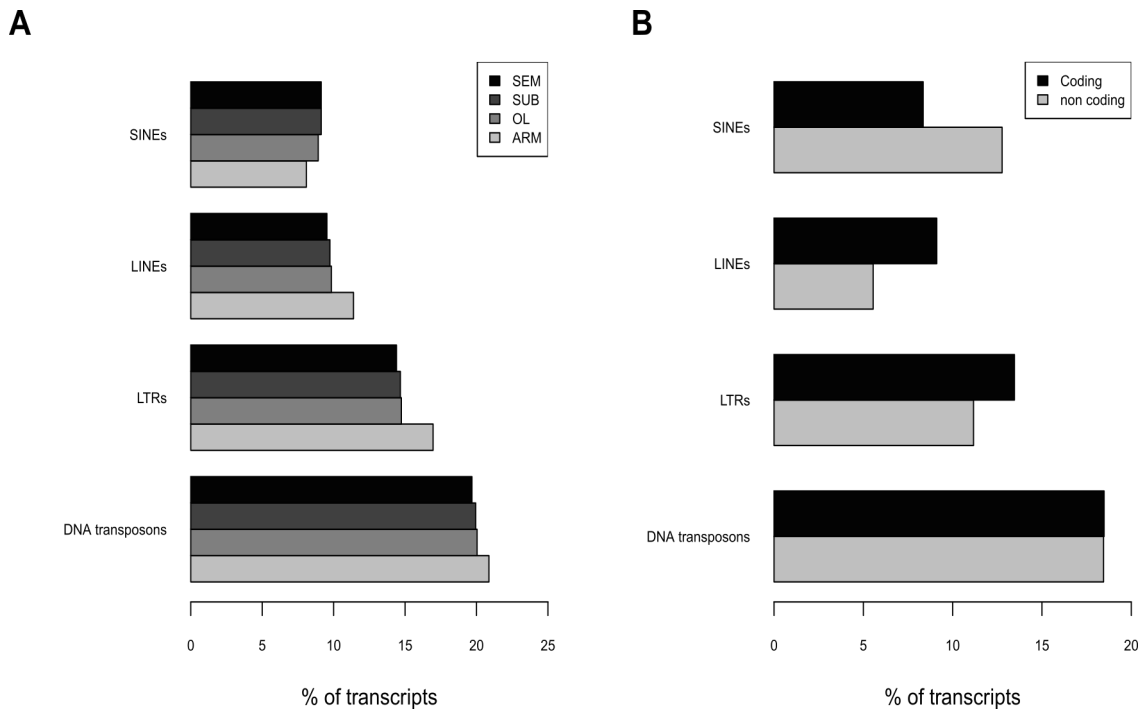


Figure S4: Percentages of transcripts associated to transposable elements. a, Barplots representing the percentage of expressed transcripts containing a fragment from a transposable element. SINE elements are most frequently embedded in brain-expressed transcripts while LINES, LTRs and DNA transposons are associated to a higher number of arm-expressed transcripts (brain: SEM, SUB, OL; arm: ARM). **b**, Barplots representing the percentage of expressed coding and non-coding transcripts containing fragments from the different transposons. SINEs are enriched in non-coding transcripts while LINES and LTRs fragments are more frequently embedded in coding transcripts. DNA transposons results to be equally distributed.

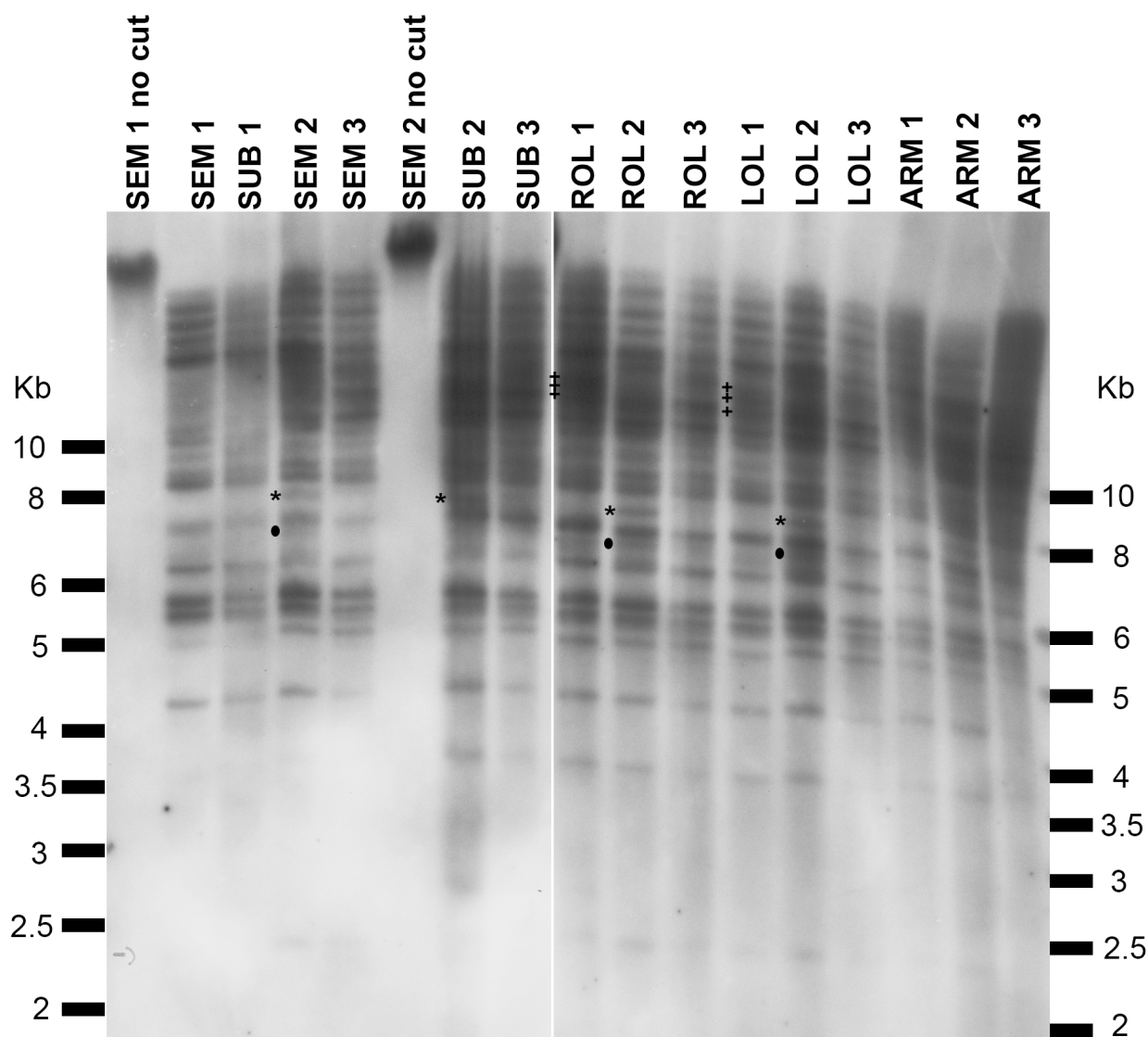


Figure S5: Southern blot analysis of the *Octopus vulgaris* LINE element. Genomic DNA extracted from SEM, SUB, OL and ARM of three different animals (#: 1, 2 and 3) were digested by restriction enzyme EcoRI and analyzed by Southern blotting. The fragment indicated by an asterisk (*) is specific of the individual #2. A plus symbol (+) highlights fragments present in OL and probably ARM of the individual #1, but not in SEM and SUB of the same individual. A dot (•) marks a fragment found in all the tissues of octopus #2 except in the SUB. DNA molecular markers are reported on both sides of the panel.

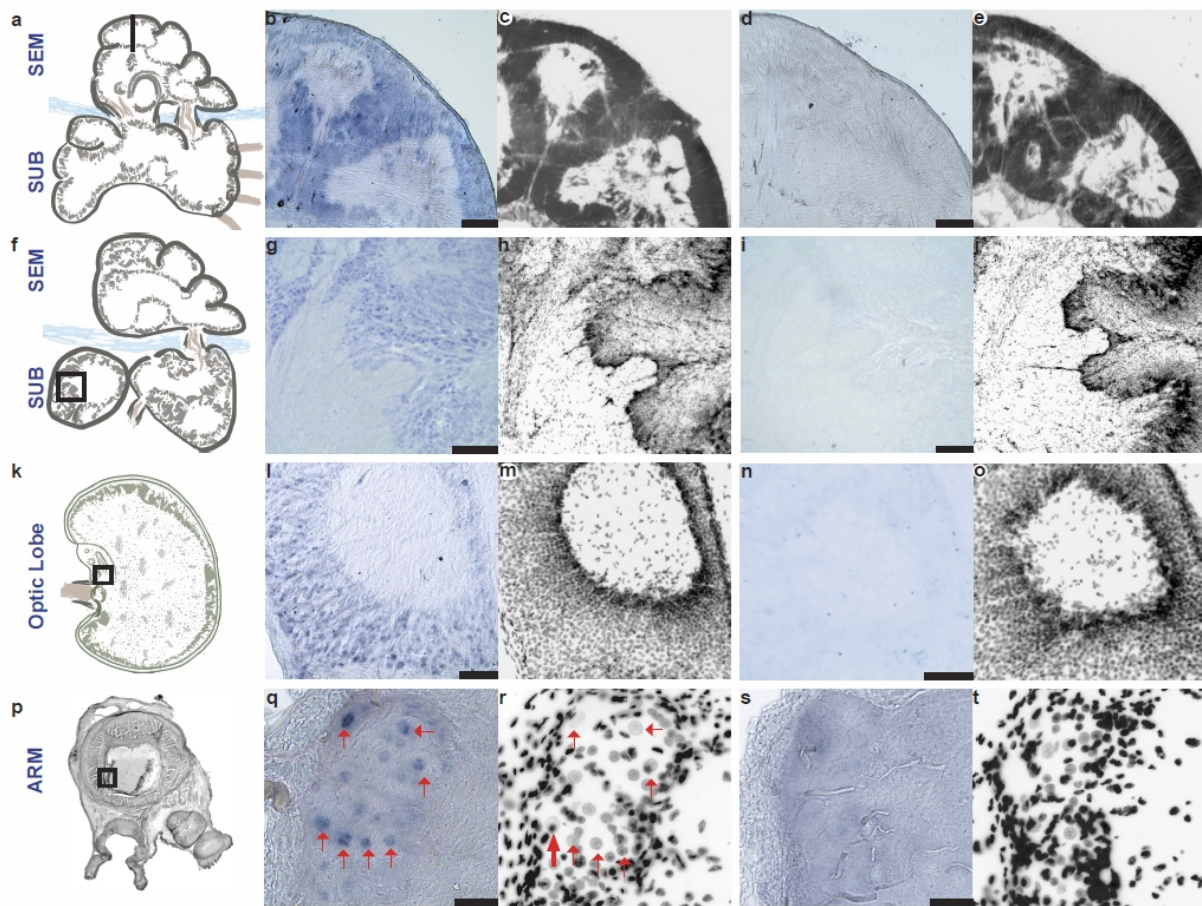


Figure S6: *In situ* hybridization for LINE in the brain (SEM, SUB and OL) and arm of adult *Octopus vulgaris*, and relative controls. Schematic diagrams (**a**, **f**, **k**, **p**) of the organization of the octopus brain (SEM and SUB) at the sagittal (**a**) and parasagittal (**f**) planes, of the optic lobe (OL, horizontal section across the midline: **k**), and of the arm (**p**) with the typical distribution of different muscular bundles surrounding the arm nerve cord, in the middle; suckers appear on the ventral side (bottom). The diagram of the octopus arm has been drawn by superimposing tracings (after Milligan staining) of a typical transverse section at the medial length of an arm. **b**, RTE-2_OV mRNA is detected in numerous cells of the vertical lobe mostly at the cortical layer of each girus. **c**, the corresponding section after DAPI staining where only nuclei appear marked. **d,e**, A similar section at the level of the vertical lobe after staining with sense RTE-2_OV probe serve as control (the same section is shown also after DAPI staining, **e**). **g-l**, Sections of the SUB at the level of the posterior pedal lobe, with positive cells marked after *in situ* hybridization (**g**) and the corresponding DAPI stained cells (**h**) to reveal the intricate patterns of neurons. Control staining (sense) where no positive cells are revealed are shown in (**i**) again with the same section stained to show nuclei (DAPI, **j**). **l**, An area of the peduncle lobe (OL) at the level of the spine where RTE-2_OV mRNA appear localized (the same section after DAPI, **m**). **n,o**, A nearby proximate section at the same level (OL) after *in situ* hybridization with sense RTE-2_OV probe (**n**) and DAPI staining (**o**). **q**, RTE-2_OV mRNA is seen at the octopus arm only in few large motor neurons (arrows) of the nerve cord; note the corresponding section after hybridization with sense probe (**s**). DAPI staining of the same sections is shown in (**r**, **t**). Sections stained with DAPI are presented to show the cellular populations of the respective areas. Scale bars: 100 μ m.