



¹ Department of BioHealth Informatics, Indiana University School of Informatics and Computing, IUPUI. ² Centre for Computational Sciences, Indianapolis. ³ Department of Medical and Molecular Genetics, Indiana University School of Medicine, Indianapolis.

ABSTRACT

With advancements in in omics technologies, the range of biological processes long non-coding RNAs (lncRNAs) are involved in is expanding extensively [1, 2]. The accelerating rate of evidence discovery of lncRNAs' role in various critical biochemical, cellular, and physiological processes is necessitating a robust platform of lncRNA annotation resources. Available resources with lncRNA ontology annotations are rare: despite a plethora of resources for annotating genes, and an extensive body of lncRNA literature. Here, we present a lncRNA annotation extractor and repository (Lantern), was developed using PubMed's abstract retrieval engine and NCBO's recommender annotation system [3]. Between 1-150 abstracts were extracted per lncRNA, which were subsequently used for extracting annotations with respect to each ontology by querying NCBO's recommender system via Application Programming Interface (API). To evaluate the quality of annotations in Lantern, benchmarking analysis was performed by deploying Lantern's pipeline over 182 lncRNAs from lncRNAdb [4] and compared the extracted annotations against annotations mapped onto the lncRNAdb's manually curated free text. Benchmarking analysis suggested that Lantern has a recall of 0.62 against lncRNAdb for 182 lncRNAs and precision of 0.8 based on manual verification of ontology annotations for 50 lncRNAs. Additionally, IncRNAs were also annotated with multiple omics information, like: RBP-interaction, tissue specific expression, protein c-expression, coding potential, sub-cellular localization and SNPs for around 11000 lncRNAs; retrieved and analyzed by robust NGS tools and pipelines. The extracted annotations for 11000 lncRNAs are available at http://www.iupui.edu/~sysbio/lantern/.

OBJECTIVES

- Develop an automated pipeline for annotating lncRNAs with multiple ontologies.
- Extract up to date annotations for multiple ontologies using latest literature.
- Develop a lncRNA annotation resource with a web interface.
- Annotate lncRNAs with high quality controlled multiple omics annotations for improving the annotation of the noncoding transcriptome.

MATERIALS & METHODS

- LncRNA relevant literature extracted using PubMed API.
- Ontology annotations from the literature were extracted using NCBO ontology recommender.
- Extracted ontologies were benchmarked against annotations extracted from lncRNAdb.
- The following formulae were used to calculate precision and recall:

Recall = <u>Number of annotations common across extracted annotations and gold standard annotations</u> Total number of annotations extracted

Precision was calculated based on manual verification of GO annotations for 50 randomly selected lncRNAs. The formula used to calculate precision is as follows:



Figure 1: Lantern workflow overview. Flow chart show various NGS and ontology extraction pipelines merging to annotate lncRNAs with multi-level omics and ontology information. All the annotation pipelines were integrated, and the extracted annotations are hosted on the online resource. The ontology annotation extraction pipeline showing two part work flow, first part benchmarks the pipeline developed, and the annotations extracted, second part shows the implementation of the developed pipeline over all the human lncRNAs from GENCODE. Respective flowcharts on either side showing data resources and steps involved in extracting diverse set of omics data for annotating lncRNAs.

Lantern: a semi-automated pipeline and repository for annotating lncRNAs with ontologies Swapna Vidhur Daulatabad¹, Rajneesh Srivastava¹, Sarath Chandra Janga^{1, 2, 3}

Lantern currently hosts annotation information for 11,290 lncRNAs. pipeline were benchmarked.

• where a lncRNA localizes in the cell,

- which tissue is more representative of which lncRNA,
- which RBPs interact with which lncRNAs,
- SNPs on lncRNA specific to a tissue and phenotype and

	Number of IncRNAs with this annotation	Annotation
Gene ontology, hum ontology and SNON 69	769	Ontology
Acro	7942	RBP-interactions
Acr	9982	Tissue specific expression
Acı	5331	Protein co-expression
For 2	9898	Coding Potential
For 10 cellular com	11290	Sub-cellular localization
Acr	6714	GTEx eQTL SNPs
1421 phen	2569	GWAS SNPs

the domain feature.

recall and precision

tations both sources (PubMed and lncRNAdb)



100 to 150 abstracts extracted per lncRNA.





0251562(MALAT1)	
Brain - Anrygdala Pancreas Brain - Anrygdala Pancreas Skin - Sun Exposed (Lower leg) Spleen Spleen Skin - Sun Exposed (Lower leg) Spleen Spleen Skin - Sun Exposed (Lower leg) Spleen Skin - Sun Exposed (Lower leg) Spleen Skin - Substand Underson (Lower leg) Brain - Hippocampus Underson (Lower leg) Brain - Nucleus accumbers (basal ganglia) Userus Brain - Caudate (basal ganglia) Brain - Cortex Brain - Cucket (basal ganglia) Brain - Cortex Brain - Caudate (basal ganglia) Heart - Left Ventricle Minor Salivary Gland Lung Colon - Fransverse Colon - Fransverse Adrenal Gand Lung Cells - Transformed fibroblasts Artere - Tibial Rrein - Substantia nigra Lung Cells - EBV-transformed fymphocytes Lung	
Brain - Amygdala I Pancreas Pancreas Adipose - Subcutaneous 5pleen Skin - Sun Exposed (Luwer leg) Fpleen Brain - Hippocampus Nunde Blood Brain - Hippocampus Nunde Blood Brain - Nucleus accumbers (basal ganglia) Uterus Brain - Nucleus accumbers (basal ganglia) Heart - Left Vertricle Brain - Fontal Cortex (BA9) Stomach Brain - Fontal Cortex (BA9) Heart - Left Vertricle Colon - Transverse Lung Brain - Fontal Cortex (BA9) Heart - Left Vertricle Colon - Transverse Lung Brain - Fontal Cortex (Ba9) Heart - Left Vertricle Colon - Transverse Lung Colon - Transverse Heart - Left Vertricle Colon - Transformed fibroblasts Hung Brain - Substantia nigra Heart - Left Vertricle Cells - EBV-transformed fibroblasts Lung Cells - EBV-transformed fymphocytes Long	
boxplots showing global expression of putative lncRNA effic expression of lncRNA HULC, observed to highly up-	
A ontology annotation extractor and repository.	
LUSION hcRNA with ontology information using litera- hchmarked against a manually curated posited onto a web interface to easily navigate NAs was put together, along with tissue specific zing contemporary NGS data analysis pipelines. hality-controlled ontology annotations and com- for improving the annotation of the noncoding	
ENCES alth and disease—size and shape matter. Brief- 15-129. s of long non-coding RNA biogenesis and func- p. 47. ogy Recommender 2.0: an enhanced approach urnal of biomedical semantics, 2017. 8(1): p. 21. ding the reference database for functional long 14. 43(D1): p. D168-D173.	
thank members of the Janga lab for helpful dis-	