Annotating long non-coding RNAs in model and non-model organisms using a Random Forest strategy FEELnc: FIExible Extraction of LncRNAs

#### Valentin Wucher Roderic Guigó's Lab Centre de Regulació Genòmica (CRG) - Barcelona - Spain

Séminaire MIAT: 22 September 2017



80% of the variants associated with diseases (by GWAS) are localized outside of protein-coding genes (Manolio *et al.*, 2009; Hindorff *et al.*, 2009).

>60% of the human genome is transcribed into RNAs with only 2% corresponding to proteins (Human ENCODE Consortium; Djebali *et al.*, 2012, Mouse ENCODE Consortium, 2015).

#### Need to identify ncRNAs to better annotate genome

Ease the interpretation of genotype to phenotype relationships.

Prospects

### Non-coding RNA

Different type of non-coding RNAs.

From GENCODE v27 annotation (Harrow et al., 2012):

- Total No of Genes: 58,288;
- Protein-coding genes: 19,836 (34%);
- Long non-coding RNA genes: 15,778 (27%);
- Small non-coding RNA genes: 7,569 (13%);
- . . .



## Long non-coding RNA characteristics

Definition:

• Transcripts without coding potential, longer than 200 nt, polyA+/- (Derrien *et al.*, 2012).

Functions (non-exhaustive):

- Can enhance or repress transcription of targeted mRNA(s);
- Can act in *cis* or in *trans*;
- Sponge for microRNAs;
- Make IncRNA protein complexes.

Examples:

- Xist: binds to PRC2 (DNMT3A)  $\rightarrow$  IncRNA-protein complexe;
- LncRNA-protein complexe  $\rightarrow$  DNA hypermethylation;
- $\bullet~\text{DNA}$  hypermethylation  $\rightarrow$  silencing X chromosome.

Prospects

## Standard pipeline for RNA-seq analysis: mRNAs + IncRNAs



Djebali et al., 2017

#### Bottleneck

Lots of novel transcripts.

Introduction

EELnc description

Benchmark

Non-model species

onclusion

Prospects

## How to deal with all assembled transcripts



Novel transcripts = IncRNAs + mRNAs + spurious transcription

Classical pipeline to annotate new transcripts:

- 1. Filter: remove short transcripts, smaller than 200 nt long.
- 2. Discriminate: determine whether the transcript is coding or not.
- 3. Classify: classify the lncRNAs regarding to nearest RNA genes.

Issues:

- The filter (1) and classify (3) steps are made manually;
- Only a minimal classification (3) by one of the tool, none for the others;
- No real guideline for non-model organisms.

But some tools exist to discriminate (2) between coding and non-coding RNAs.

Prospects

## Tools to discriminate mRNAs and IncRNAs

#### Alignment-based:

Advantages:

- High specificity;
- Identify conserved IncRNAs.

Drawbacks:

- Depends on the database;
- Depends on the alignment;
- Slow.

PhyloCSF (Lin, M. et al., 2011), CPC (Kong, L. et al., 2007), ...

#### Alignment-free:

Advantages:

- Usually Fast;
- Independent of alignment;
- Lineage-specific IncRNAs.

CPAT (Wang, L. et al. 2013), CNCI (Sun, L. et al. 2013), PLEK (Li, A. et al. 2014), ...

Drawbacks:

• Designed for model organisms.

FIExible Extraction of LncRNAs (FEELnc):

- All in one: formalized the filtering and classification into modules;
- Stringent set of IncRNAs and mRNAs;
- Classification regarding all RNAs, useful to get potential functional relations;
- LncRNAs detection is genome reference-free, i.e non-model species;
- Available for non-model organisms by replacing lncRNAs.

Introduction

FEELnc description

Benchmark

Non-model species

Conclu

Prospects

## FEELnc: FIExible Extraction of LncRNAs



Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Filter					



Aim:

• Filtering out non IncRNA-like.

Methods:

- Remove short transcripts (< 200 nt);</li>
- Flag transcripts overlapping known mRNAs;
- Keep or discard monoexonic transcripts, antisense or intergenic;

Next step:

• Defined coding and non-coding transcripts.

Prospects

# Coding Potential



#### Aim:

 Defines a protein-coding score and then a cutoff to differentiate mRNAs from IncRNAs.

#### Methods:

• Use features and machine learning, a Random Forest.

Benchmark

Non-model species

Conclusion

Prospects

## Features for the Random Forest

1. RNA size (Cabili *et al.*, 2011; Derrien *et al.*, 2012) (high value  $\rightarrow$  mRNA)

clusion

Prospects

## Features for the Random Forest

- 1. RNA size (Cabili *et al.*, 2011; Derrien *et al.*, 2012) (high value  $\rightarrow$  mRNA)
- 2. ORF coverage (ORF defined with respect to 5 modes):
  - Strict: requires start and stop;
  - Moderates: requires start or stop;
  - Relaxed: total RNA sequence.

(high value  $\rightarrow$  mRNA)

lusion

Prospects

# Features for the Random Forest

- 1. RNA size (Cabili *et al.*, 2011; Derrien *et al.*, 2012) (high value  $\rightarrow$  mRNA)
- 2. ORF coverage (ORF defined with respect to 5 modes):
  - Strict: requires start and stop;
  - Moderates: requires start or stop;
  - Relaxed: total RNA sequence.

(high value  $\rightarrow$  mRNA)

- 3. *k*-mer scores on ORF for multiple *k*-mer sizes:
  - For a specific *k*-mer, the ratio between mRNA frequency and lncRNA frequency;
  - Collaboration with INRIA/Genscale team in Rennes (Fr), KmerInShort developed by Guillaume Rizk from GATB tools (Drezen *et al.*);
  - Very fast and parallel extraction of k-mer profiles up to 15-mer.

(high value  $\rightarrow$  mRNA)

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
<i>k</i> -mer sc	ore calculation	I			

Get the *k*-mer profile for a specific size *k*:

- For all K (e.g. TGC) of size k (e.g. 3);
- Get mRNA  $F_K^m$  and lncRNA  $F_K^{lnc}$  observed frequencies;
- Calculate a score for each K:

$$S_K^k = \frac{F_K^m}{F_K^m + F_K^{lnc}}$$

Get the k-mer profile for a specific size k:

- For all K (e.g. TGC) of size k (e.g. 3);
- Get mRNA  $F_K^m$  and IncRNA  $F_K^{lnc}$  observed frequencies;
- Calculate a score for each K:

$$S_K^k = \frac{F_K^m}{F_K^m + F_K^{lnc}}$$

Get the *k*-mer score for a sequence X:

- For each ORF X;
- Get occurrences  $N_K^X$  of all K of size k;
- Calculate a score for the size k using all K:

$$V_X^k = \frac{\sum_{K=1}^{4^k} S_K^k \times N_K^X}{\sum_{j=1}^{4^k} N_j^X}$$

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects

#### Features: Illustration



Features comparison between 5,000 IncRNA sequences and 5,000 mRNA sequences (GENCODE v24) for the learning and the testing.

Non-model species

s Conc

Prospect

#### Features: Illustration



Use these features to make a Random Forest based model.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Random	Forest				

The **Random Forest** is a machine learning method:

- Forest:
  - A set of decision trees.
- Random:
  - Each tree is done on a sampling of the data;
  - Each node of each tree is done on a subset of the features.
- The model:
  - The **forest** of the trees made on a **sample** of the data.
- The prediction:
  - Each input sequences go through each tree;
  - Each tree vote for a sequence to be coding or non-coding;
  - Each input sequence got a score representing the number of trees which vote for the sequence to be coding.

Non-model species

Conclusio

Prospects

### Random Forest: Illustration



Non-model species

## Random Forest: Illustration



17

Introduction

FEELnc description

Benchmark

Non-model species

s Conclu

Prospects

### Coding potential score distribution



Get a coding potential score for all input sequences, with for the best case: IncRNA scores around 0 (blue) and mRNA scores around 1 (red).

Introduction

FEELnc description

Benchmark

Non-model species

Conclu

Prospects

### Coding potential score distribution



#### Issue

#### Which coding potential cutoff?

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Automati	c cutoff				



Make a 10-fold cross-validation:

• Compute performance on subset of the learning dataset.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Automatio	c cutoff				



Make a 10-fold cross-validation:

- Compute performance on subset of the learning dataset;
- Use to define an optimal cutoff (0.367);
- Automatically defined as Sensitivity = Specificity (0.92).

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Automati	c cutoff				



Make a 10-fold cross-validation:

- Compute performance on subset of the learning dataset;
- Use to define an optimal cutoff (0.367);
- Automatically defined as Sensitivity = Specificity (0.92).

#### lssue:

• Transcripts around the cutoff, not a high confidence in the prediction.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Automati	ic cutoff				



"The (CPS) threshold is (...) somewhat arbitrary, and transcripts that reside in questionable regions of the distribution should be annotated as transcripts of unknown coding potential (TUCPs)"

J.S. Mattick & J.L. Rinn, 2015.

#### Implemented FEELnc solution

Defined 2 cutoffs based on specificity.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
User two	cutoffs				



Implemented FEELnc solution:

- A user defined mRNAs and IncRNAs specificity (e.g. 0.95,0.95);
- Automatically set two cutoffs, one for mRNAs and one for IncRNAs (e.g. 0.225, 0.461).

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
User two	cutoffs				



Implemented FEELnc solution:

- A user defined mRNAs and IncRNAs specificity (e.g. 0.95,0.95);
- Automatically set two cutoffs, one for mRNAs and one for IncRNAs (e.g. 0.225, 0.461).

With two cutoffs, definition of a new class:

• Transcript of Unknown Coding Potential (TUCP).



Aim:

 Predict potential functional relationships between IncRNA transcripts and RNA transcripts.

Method:

- Formalized sub-classes of genomic classification genic and intergenic;
- · Get direction of the relation;
- Use a sliding window around IncRNAs;
- Get the relations for all RNA inside the window.

Conclusion

Prospects

Classes



Introduction

FEELnc description

otion

Benchmark

Non-model species

Conclusion

Prospects

## FEELnc: FIExible Extraction of LncRNAs



FEELnc, three independent modules:

Get all IncRNA-like transcripts.

Use a Random Forest to discriminate between mRNAs and IncRNAs.

Classify the IncRNAs regarding the nearest transcript.

Need to compare FEELnc predictions with state of the art methods.

Prospects

# Benchmarking the Coding Potential module

Compare the FEELnc Coding Potential module against 5 methods:

- CPAT (Wang et al., 2013);
- CNCI (Sun et al., 2013);
- PLEK (Li et al., 2014);
- CPC (Kong et al., 2007);
- PhyloCSF (Lin et al., 2011).

Use 5 performance measures to compare methods:

- Sensitivity;
- Specificity;
- Precision;
- Accuracy;
- Matthews Correlation Coefficient (MCC): summarizes others
  - 1: good predictions; 0: random predictions; -1: opposed predictions.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Benchma	arking				



Data:

- Human GENCODE (v24) training and testing data;
- Only one transcript per gene was extracted, no common genes for learning and testing;
- 5,000 mRNAs and 5,000 IncRNAs on training and testing datasets.

Methods used with default parameters.

Prospects

### Benchmarking results

program	sensitivity	specificity	precision	accuracy	MCC
FEELnc	<u>0.923</u>	0.915	0.916	<u>0.919</u>	<u>0.838</u>
CPAT	0.899	0.924	0.922	0.912	0.823
CNCI	0.829	0.979	0.975	0.904	0.817
PLEK	0.732	<u>0.985</u>	<u>0.981</u>	0.858	0.741
PhyloCSF*	0.906	0.802	0.820	0.854	0.712
CPC*	0.699	0.739	0.728	0.719	0.438

#### Best score

- \*: alignment-based
  - FEELnc performs similarly or better;
  - On this benchmark, alignment-free better than alignment-based.

RNA-seq or model reconstruction can generate truncated/incomplete transcript models (Steijger *et al.*, 2013).

Modified the benchmark dataset by removing percentage of transcripts, either in 5' or 3'.



Good performance regarding other methods even with modified datasets.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
What abou	it non-model s	species?			



Classification tools for IncRNAs work well on species with known IncRNAs.

Prospects

#### What about non-model species?



Classification tools for IncRNAs work well on species with known IncRNAs.

#### Issues

What if no lncRNAs are available, i.e. non model species?

#### Some solutions

- Mimic IncRNAs with other sequences for the learning step;
- 2. Use IncRNA sequences from evolutionary related species.

Prospects

## What about non-model species?



Classification tools for IncRNAs work well on species with known IncRNAs.

#### Issues

What if no lncRNAs are available, i.e. non model species?

#### Some solutions

- Mimic IncRNAs with other sequences for the learning step;
- 2. Use IncRNA sequences from evolutionary related species.

### Mimic IncRNA sequences

How to mimic IncRNA sequences?

- 1. LncRNAs are non-coding  $\rightarrow$  extract non-coding sequences;
- 2. LncRNAs can result from the pseudogenization of protein coding genes (Duret *et al.*, 2006)  $\rightarrow$  modify mRNA sequences.

FEELnc methods to mimic IncRNA sequences:

- 1. Intergenic module: randomly extract genomic intergenic sequences;
- 2. Shuffle module: shuffle mRNA learning sequences while preserving the 7-mer frequencies using Ushuffle (Jiang *et al.*, 2008).

#### Mimic IncRNA sequences

How to mimic IncRNA sequences?

- 1. LncRNAs are non-coding  $\rightarrow$  extract non-coding sequences;
- 2. LncRNAs can result from the pseudogenization of protein coding genes (Duret *et al.*, 2006)  $\rightarrow$  modify mRNA sequences.



#### Mimic IncRNA sequences

How to mimic IncRNA sequences?

- 1. LncRNAs are non-coding  $\rightarrow$  extract non-coding sequences;
- 2. LncRNAs can result from the pseudogenization of protein coding genes (Duret *et al.*, 2006)  $\rightarrow$  modify mRNA sequences.



Introduction

EELnc description

Benchmark

Non-model species

onclusion

Prospects

#### Mimic IncRNA sequences: Results



Learning method Performance on the human testing dataset using as learning IncRNAs either the GENCODE IncRNAs, the shuffle module or the intergenic module.

Prospects

### What about non-model species?



Classification tools for IncRNAs work well on species with known IncRNAs.

#### Issues

What if no lncRNAs are available, i.e. non model species?

#### Some solutions

- Mimic IncRNAs with other sequences for the learning step;
- 2. Use IncRNA sequences from evolutionary related species.

Prospects

## Evolutionary related IncRNA sequences



Train the Random Forest:

- Use mRNAs of the species;
- Use IncRNAs from the evolutionary related species.

#### Apply the model:

• On transcript models of the species.

#### How to test?

Learning IncRNAs sequences from the NONCODE 2016 database (Zhao *et al.*, 2016)

Introduction

EELnc descriptior

Bend

Non-model species

Conclusion

Prospects

#### Evolutionary related IncRNA sequences: Results



FEELnc performance and times of speciation are anti-correlated.

Introduction

EELnc description

Benchmark

Non-model species

Conclusion

Prospects

#### Evolutionary related IncRNA sequences: Results



Shuffle module: 0.748 MCC  $\rightarrow$   ${\sim}100$  Myr.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Applicatio	n: Dog				

Dog (Wucher et al., 2017):

- Collaboration with the European LUPA Consortium (Michel Georges) and the BROAD institute (Kerstin Lindblad-Toh);
- 16 tissues;
- 20 RNA-seq;
- ~2,500 new IncRNA genes;
- $\bullet ~\sim \!\! 10,000$  new IncRNA transcripts.



Dog (Chris Barber, Wikipedia)

Introduction

EELnc description

Benchma<u>rk</u>

Non-model species

Conclus

Prospects

### Applications: Others



Chicken (Andrei Niemimäki, Wikipedia)



Ectocarpus (Akirapeters, Wikipedia)

Chicken (Muret et al., 2017):

- With Sandrine Lagarrigue (INRA, Agrocampus, FR);
- Adipose and liver tissues;
- 16 RNA-seq;
- ~2,200 new IncRNA genes.

Ectocarpus (algae) (Cormier et al., 2016):

- With Mark Cock (CNRS,Roscoff, FR);
- Male and femelle gametophytes;
- 10 RNA-seq;
- ~700 new IncRNA genes;
- First IncRNA catalogue in algae.

Conclusion

Prospects

#### One tool: three applications

#### FEELnc: from transcript models to IncRNA classifications.



Conclusion	Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
	Conclusi	on				

User friendly:

- Automatic threshold;
- Easy to use.

Flexible:

- Set two specificity thresholds for stringent predictions;
- Five ORF type definitions;
- Coding potential can be used on FASTA or GTF.

Non-model species compatible:

- Mimic IncRNAs by shuffling mRNA sequences;
- Coding potential module is alignment- and genome reference-free;
- Guideline for species without annotated IncRNAs.

Performs similar or better than other tools.

FEELnc on github: https://github.com/tderrien/FEELnc

A nextflow/docker implementation of STAR/Cufflinks/FEELnc pipeline has been developed by Evan FLODEN: https://github.com/skptic/IncRNA-Annotation-nf

Published in Nucleic Acids Research, along with a dog extended annotation: Wucher *et al.*, 2017.

All data (included benchmarking data), command lines and scripts to make figures are available through Supplementary:

http://nar.oxford journals.org/content/early/2017/01/03/nar.gkw1306/suppl/DC1

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Prospects					

Method:

- Add new features, e.g transcript expressions or exons number;
- Modified the 12-mer score, e.g. translate amino acids;
- Predict pseudogenes (can be still reference-free?).

Bioanalyses:

- FEELnc used in the FAANG consortium, i.e. farming animals;
- Improve IncRNA classification by adding new data, as chromosome configuration capture (Hi-C);
- Use and integrate IncRNA classification with multi-omics data.

Benchmark

Non-model species

s Conclu

Prospect

#### Acknowledgments

CNRS - IGDR - Dog genetics team (Rennes):

**Thomas Derrien**, **Christophe Hitte**, Catherine André, Céline Le Béguec

INRA/INRIA - IGEPP/IRISA -EGI/Genscale team (Rennes): Fabrice Legeai, Guillaume Rizk

INRA - PEGASE (Rennes): Kévin Muret, Sandrine Lagarrigue

CNRS - Algal Genetics Group (Roscoff): Alexandre Cormier, Erwan Corre, Susana M. Coelho, Mark Cock

Nathalie Villa-Vialaneix

And all of you for your attention!

INRA - Toulouse: Sarah Djebali, Sylvain Foissac

CRG - Comparative Bioinformatics (Barcelona): Evan Floden

CRG - Computational Biology of RNA Processing (Barcelona): Roderic Guigó and all lab's members.

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects
Annexes					

Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects





Introduction	FEELnc description	Benchmark	Non-model species	Conclusion	Prospects







Original sequence: AGACTTAGCA Original count:

AC	AG	CA	СТ	TA	TT	GA	GC
1	2	1	1	1	1	1	1

Ushuffle with *k*-mer size = 2

Permuted sequence: ACTTAGCAGA Permuted count:

AC	AG	CA	СТ	TA	TT	GA	GC
1	2	1	1	1	1	1	1

Mechanisms of IncRNA-mediated regulation	n Associated examples of lncRNAs in cancer
A Chromatin remodelling	ANRIL: PCR1- mediated repression of INK4A-ARF-INK4 tumour suppressor locus, upregulated in prostate cancer hotspot in various GWAS (Kotake et al, 2011; Pasmant et al 2011)
Heterochromatin	XIST: Involved in X-chromosomal inactivation downregulated in female breast, ovarian and cervical cancer cell lines (Kawakami et al., 2004), suppresses haematologie cancer în vivo în mice (Yildirim et al., 2013)
/ / / / / / / / / / / / / / / / / / / /	KCNQ1071: Loss of imprinting in colorectal cancer (Nakano et al, 2006)
	HOTAIR: Overexpressed in breast cancer, promotes cancer metastasis (Gupta et al, 2010)
3 Transcriptional co-activation and -repression Activation	LincRNA-p21: Regulation of p53 response upon DNA damage; upregulated in various cancer cell lines (Huarte et al 2010)
or inhibition	H19: Upregulated in gastric cancer; ectopic expression promotes cell proliferation (Yang et al. 2012)
	SRA: Transcriptional coactivator of stereoid receptors upregulated in breast tumorigenesis (Leygue et al, 1999)
Protein inhibition Telomorase 3' TERRA	<b>TERRA:</b> Facilitates telomeric heterochromatin formation and inhibits telomerase by direct binding: expression significantly reduced in many human cancer cell lines (Redon <i>et al</i> , 2010)
Post-transcriptional modifications Pre-	MALATI: Control of alternative splicing by regulating the distribution of serine/arginine splicing factors (SR) and their protein levels in nuclear speckles, upregulated in various eaacer tissues, promotes cell motility and proliferation (Schmidt et al. 2011: Tripathi et al. 2010; Xu et al. 2011)
Decoy incRNA incRNA in cRNA in steady state overspin steady state	A Med PTENP1: Pseudogene of the tumour suppressor gene PTEN controls PTEN expression levels by competing for microRNA binding with PTEN; lost in many human cancers (Poliseno e al, 2010)

From: Cheetham, S. W., Gruhl, F., Mattick, J. S., & Dinger, M. E. (2013). Long noncoding RNAs and the genetics of cancer. British journal of cancer, 108(12), 2419.