

tRNAscan-SE

by Kevin Gleason

April 28, 2015

What is tRNAscan-SE

- Peter Schattner, Angela N. Brooks and Todd M. Lowe
 - UCSC RNA Center, 2005
- Identifies 99–100% tRNA genes in DNA sequence
- Less than one false positive per 15 gigabases.
- Offered as web server and UNIX package
- Solves the issue of custom programming for RNA identification

Purpose

- Transfer RNA genes make up the single largest gene family
- Conventional gene finders do not target tRNA
 - Lack codon structure and statistical signatures of protein coding genes
- tRNA affects the Codon Bias in protein coding genes

Web Server

- The search mode selection
 - Determines which probabilistic model to be used in
 - Training data from selected species or phylogenetic groups (mammals, yeasts and archaea)
- Query Sequence Selection
 - Sequences to be searched

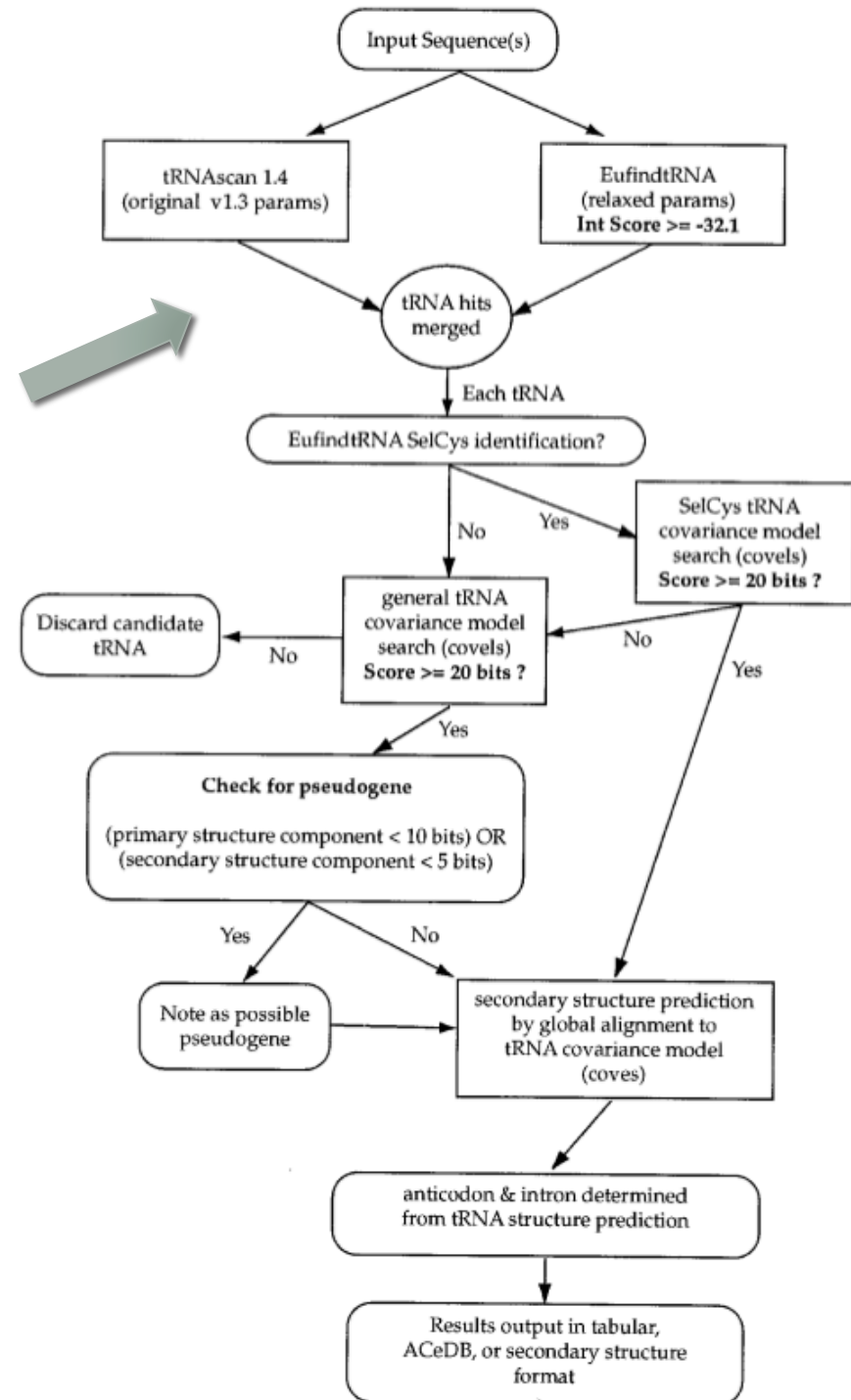
Stochastic CFG

- Watson-Crick and GU Wobble
- $S \rightarrow \lambda \mid AS \mid CS \mid GS \mid US \mid ASU \mid USA$
 $\mid CSG \mid GSC \mid GSU \mid USG \mid SS$
- Probabilities from frequencies of features in databases of RNA structure
- Can predict pairwise and nested structures, but not pseudoknots

The Algorithm

Phase 1

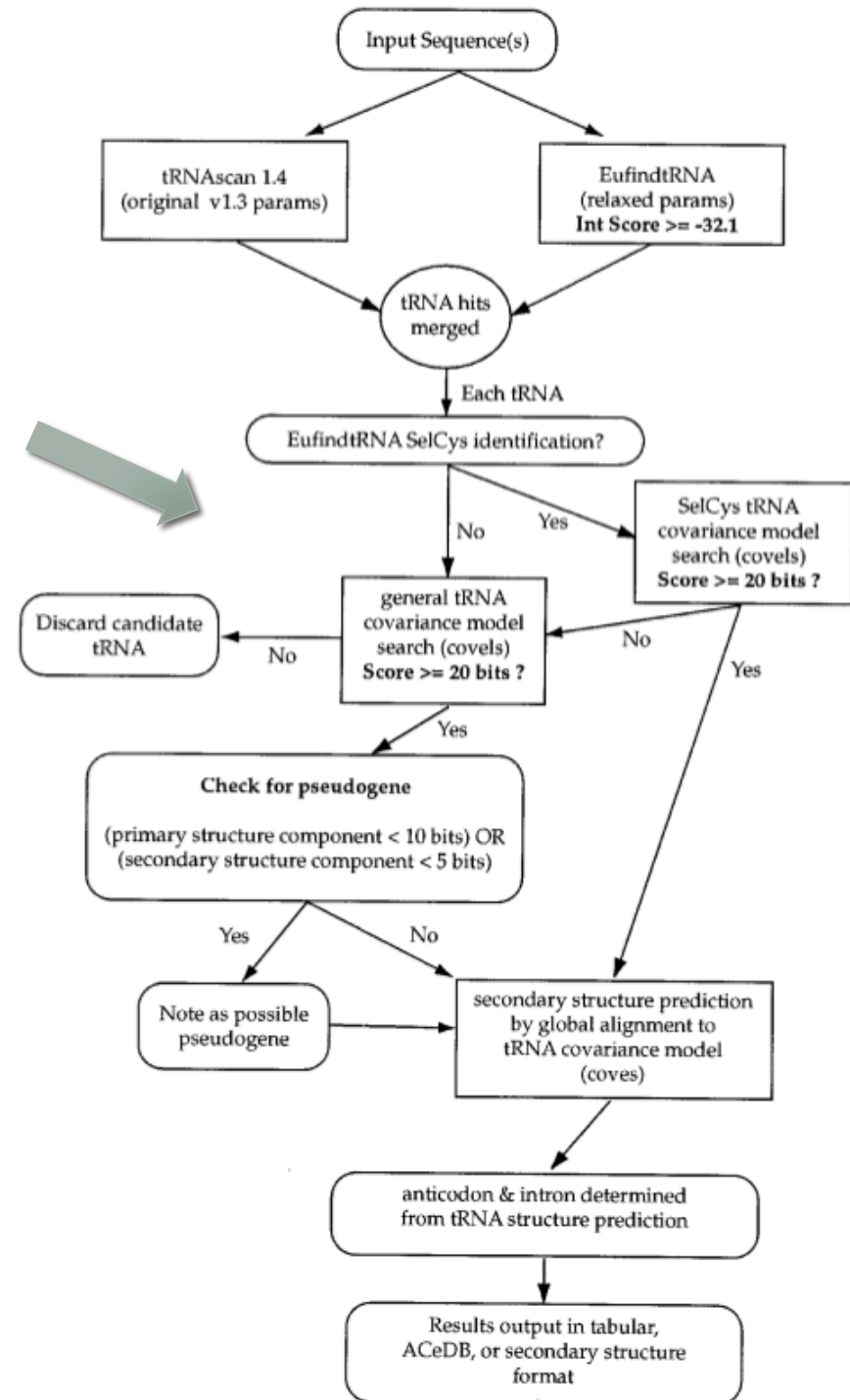
- tRNAscan and Parvesi Algorithm
- Parvesi
 - Linear sequence signals
 - RNA polymerase III promoters and terminators
 - Less specific, but predicts tRNAs that tRNAscan misses
- Results merged



The Algorithm

Phase 2

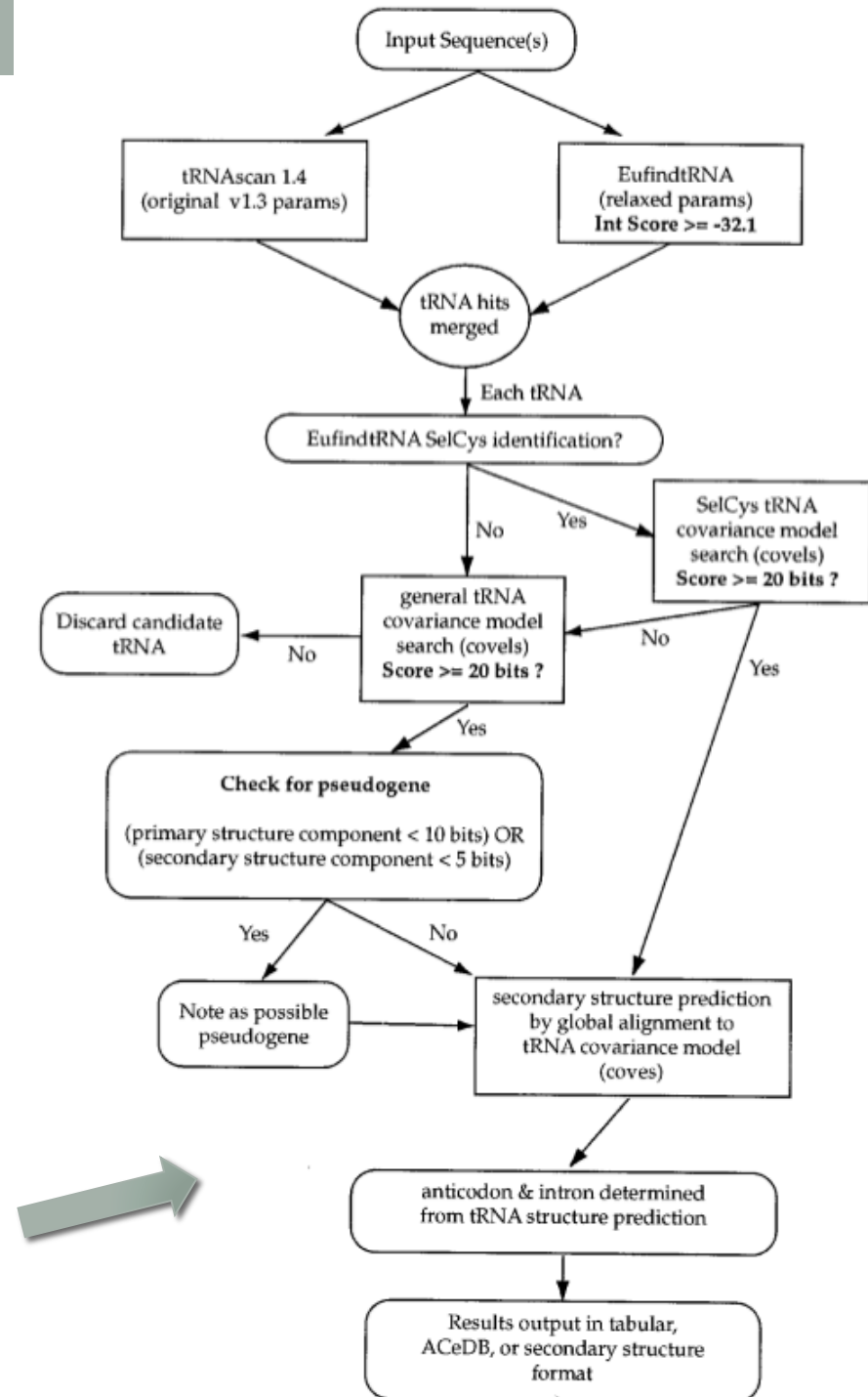
- Extracts the candidate subsequences
 - 7nt padding
- Applies a tRNA covariance model
- Structurally aligned 1415 tRNAs from modified 1993 Sprinzl database
- Discard all low scores



The Algorithm

Phase 3

- tRNA with log odds scores over 20.0 bits
- Trims bounds to covels predictions
- Predicts secondary structure with coves
 - Pseudogenes filtered out
- A second profile created from HMM
 - No secondary structure is conserved
- Predictions formatted and saved



Output

- Final tRNA predictions are saved in tabular, ACeDB or secondary structure output format.
- AceDB - originally "**A C. elegans DataBase**" - is a genome database designed for handling bioinformatic data

(A) tRNAscan-SE Output

Sequence Name	tRNA #	tRNA Begin	Bounds End	tRNA Anti Type	Codon	Intron Begin	Bounds End	Cove Score
chr6.trna18-AlaAGC	1	1	73	Ala	AGC	0	0	40.39

tRNAs decoding Standard 20 AA: 1
 Selenocysteine tRNAs (TCA): 0
 Possible suppressor tRNAs (CTA,TTA): 0
 tRNAs with undetermined/unknown isotypes: 0
 Predicted pseudogenes: 0

Total tRNAs: 1

Isotype / Anticodon Counts:

Ala	: 1	AGC: 1	GGC:	CGC:	TGC:
Gly	: 0	ACC:	GCC:	CCC:	TCC:

Predicted tRNA Secondary Structures:

chr6.trna18-AlaAGC.trna1 (1-73) Length: 73 bp
 Type: Ala Anticodon: AGC at 34-36 (34-36) Score: 40.39
 Seq: GGGGGATTAGCTCAAGCGGTAGGGTGCCTTGCTTAGCATGCAAGAGGtAGCAGGATCGACGCCTGCATTCTCCA
 Str: >>>>>>...>>>>.....<<<<.>>>>.....<<<<.<.....>>>>.....<<<<<<<<<<<<.

Output

Selenocysteine tRNAs (TCA)	4
Possible suppressor tRNAs (CTA,TTA)	0
tRNAs with undetermined or unknown isotypes	10
Predicted pseudogenes	224
Total tRNAs	1516

Intron Summary

tRNAs with introns	Gly	Pro	Val	Val	Arg	Arg	Leu	Asn	Lys
16	GCC	TGG	CAC	TAC	GCG	CCT	CAA	ATT	TTT
	1	1	1	1	1	1	8	1	1

Four Box tRNA Sets

Isotype	tRNA Count by Anticodon				Total
Ala	AGC	GGC	CGC	TGC	81
	51	1	13	16	
Gly	ACC	GCC	CCC	TCC	88
		35	2	51	
Pro	AGG	GGG	CGG	TGG	58
	33		12	13	
Thr	AGT	GGT	CGT	TGT	121
	34	58	18	11	
Val	AAC	GAC	CAC	TAC	56
	28		19	9	

Six Box tRNA Sets

Isotype	tRNA Count by Anticodon								Total
Ser	AGA	GGA	CGA	TGA	ACT	GCT			58
	22	1	10	10	1	14			
Arg	ACG	GCG	CCG	TCG			CCT	TCT	172
	138	1	2	11			7	13	
Leu	AAG	GAG	CAG	TAG			CAA	TAA	66
	33		16	6			8	3	

Two Box tRNA Sets

Isotype	tRNA Count by Anticodon				Total
Phe	AAA	GAA			23
		23			
Asn	ATT	GTT			46
	12	34			
Lys			CTT	TTT	130
			60	70	
Asp	ATC	GTC			49
		49			

Two Box & Other tRNA Sets

Isotype	tRNA Count by Anticodon				Total
Ile	AAT	GAT		TAT	51
	32	1		18	
Met			CAT		37
			37		
Tyr	ATA	GTA			31
	2	29			
Supres			CTA	TTA	0