

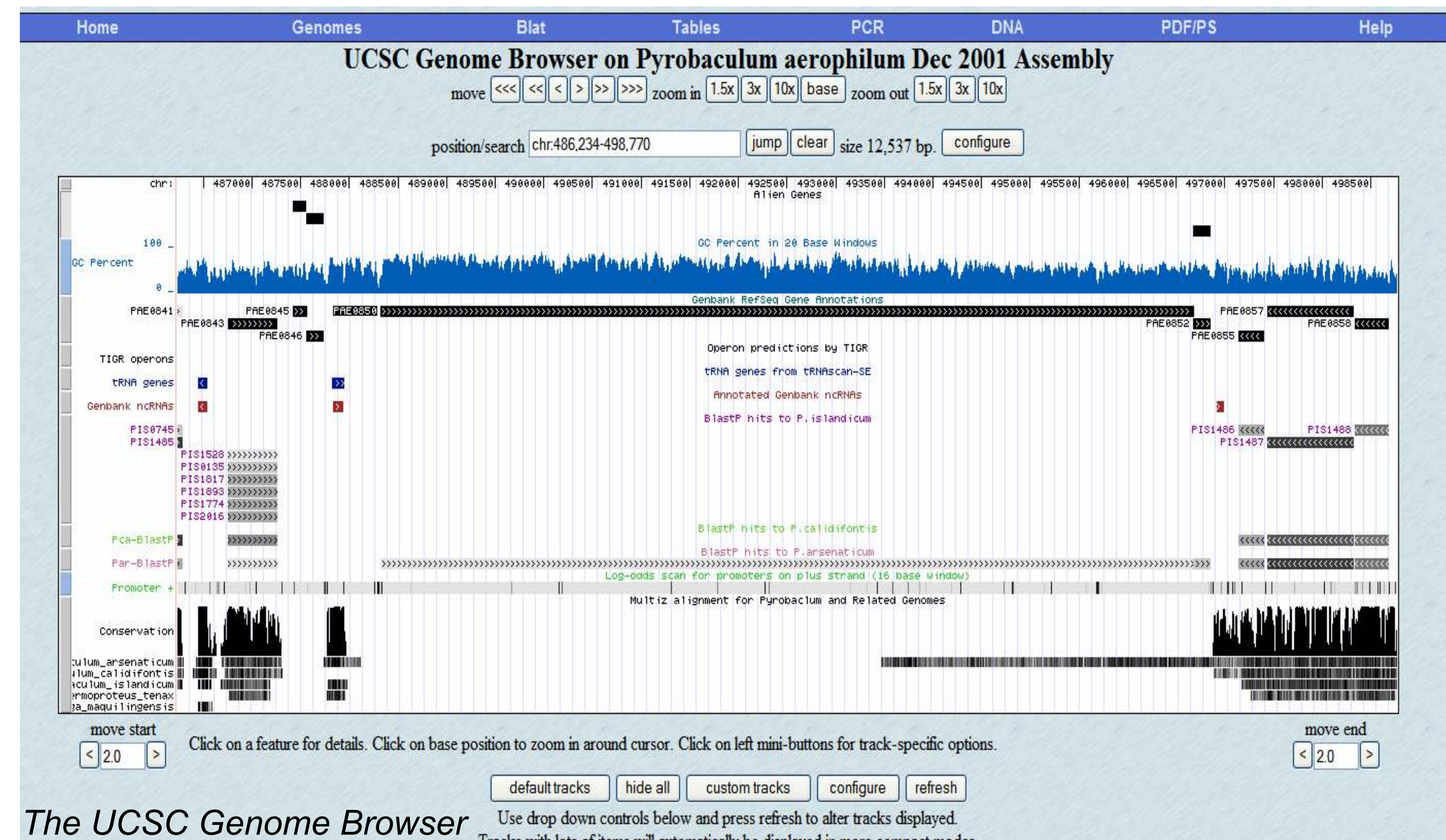
# Using Conserved Elements to Assist in Gene Annotation

Christoph Rau

Todd Lowe, Advisor

## Why is it important?

- Attempting to discover/annotate every important gene by hand would be nearly impossible.
- Early algorithms had lower accuracies, and returned many false positives.
- Better algorithms result in less 'busy' work for the researcher and more quality time investigating the results.



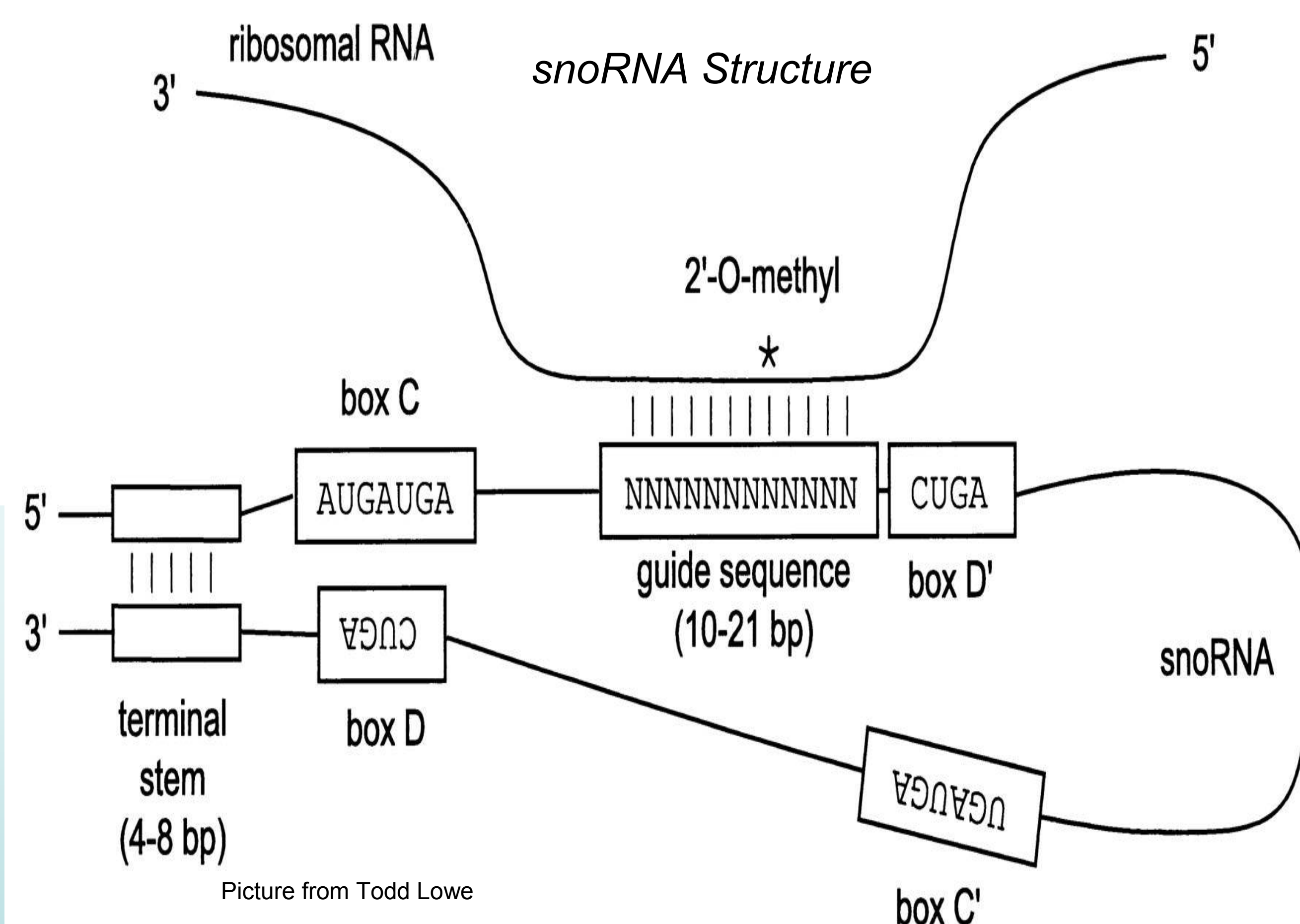
The UCSC Genome Browser Picture from Todd Lowe

## Finding Highly Expressed and 'Alien' Genes

- Not all genes in a genome are equally interesting to a researcher.
- Two interesting categories are highly expressed genes, whose proteins appear in higher amounts in the organism, and 'alien' genes, whose genes appear to be foreign to the genome.
- Highly expressed gene predictions can be used to determine metabolic pathways, while alien gene predictions can be used to examine evolutionary questions.
- The algorithm uses a 'codon usage' analysis. It is postulated that similarly expressed proteins will have similar codon frequencies.
- Each codon is a sequence of three nucleotides that codes for an amino acid.
- If a gene shares similar codon frequencies with known highly expressed genes and shows differences from the rest of the genome, it is marked as "highly expressed."
- If a gene's codon usage differs from both the whole genome AND the known highly expressed genes, then it is marked as an 'alien' gene.

## Finding SnoRNAs

- RNA molecules are an intermediate step between DNA and proteins
- However, not all genes code for proteins. Some form RNA molecules that perform other functions
- Small nucleolar RNAs (snoRNAs) modify other non-coding RNAs like rRNAs and tRNAs
- There is some evidence that suggests they can act on mRNAs as well, but this is not well documented



Picture from Todd Lowe

- The original algorithm uses the highly conserved structure of snoRNAs to probabilistically determine potential snoRNAs based on rRNA modifications.
- A basic tenant of comparative genomics is that important sequences are conserved between genomes.
- The new algorithm compares the snoRNAs of different genomes, increasing the overall accuracy.
- This process also allows multiple genomes to be analyzed at speeding up the annotation process

Thank you

David Bernick Todd Lowe



This work was completed as part of UCSC's SURF-IT summer undergraduate research program, an NSF CISE REU Site. This material is based upon work supported by the National Science Foundation under Grant No. CCF-0552688

Pyrococcus horikoshii sRNAs				snoRNA results					
Name	C Box	D' Box	C' Box	D Box					
sR1	AAAGAAGGG	ATGATGA	AGCCTCCGCAC	CTGA	ATGA	TGAGGA	GTGGACGGCTC	CTGA	GCTACTCCT
sR2	AAAAAGAGG	ATGATGA	GTTTTCCCTCACT	CTGA	GGAG	TGATGA	GGAGCCGATCA	CTGA	CCTGATCAT
sR3	AATTGTGGC	ATGATGA	ATAGCAAGCCAG	CTGA	AGAG	TGATGA	AGTGAACACCC	CTGA	GCTTACATAA
sR4	CCGATTTGG	ATGATGA	GGGAGATTTCGG	CCGA	GTGG	TGAGGA	GACTCGATGGG	CTGA	CTTTCTAAG
sR5	ATAAGGTGTG	ATGATGA	ACGCCATCGATA	CTGA	GATA	TGATGA	CCGATTCCTGG	CTGA	TTTCTTTAT
sR6	TGGAATGGG	ATGATGA	AGTTTGTACCC	CTGA	AGAA	TGATGA	ACCCTGCCCTTA	CTGA	CCATTAGCT
sR7	GTAATCCGG	ATGATGA	ACCTCATCCCAA	CTGA	ATAAA	TGATGA	ATGCATCAGG	CTGA	CATTACCTT
sR8	CGAATAGCG	ATGATGA	GCTCCATCCCTAC	CTGA	GTG	TGATGA	ATGTAGCCGCT	CTGA	GCAACCCTTA
sR9	CCGACAGGG	ATGATGA	GCTTTGCTTTG	CTGA	GCAGA	TGATGA	CCAGCCCTTGG	CTGA	CCTGCTAATT
sR10	TTAGGACTC	ATGATGA	ACTACTCCGGG	CTGA	GGGG	TGATGA	ACCACCTACCGG	CTGA	GGTGAAGCA
sR11	CTCTAGCGG	ATGATGA	CTTTCCGAGTG	CTGA	GCTGG	TGATGA	GTAAACAGCTGT	CTGA	CTTCCCTTT
sR12	ATCGGCTTG	ATGATGA	GCGTTACCGGT	CTGA	GCTG	TGATGA	TATCGGACTGT	CTGA	CTAGTATCT
sR13	TCATTGTGG	ATGATGA	GATGGCGGATTG	CTGA	GAGA	TGATGA	GACCTTAGGG	CTGA	TTAATTTCG
sR14	CTATTACCA	ATGATGA	CGGATCAACGG	CTGA	TCGAA	TGATGA	CCTCCGATCAC	CTGA	GGGTTCAAG
sR15	AGATGAGAG	ATGATGA	GTAACCCGTTG	CTGA	GGGG	TGATGA	GAGGATCGACTAG	CTGA	ACAACTCTT
sR16	ATCAAGTCTG	ATGATGA	ACCTCCCTCAC	CTGA	AAGG	TGATGA	GCACACCGTAGG	CTGA	GGGTGATAAT
sR17	ACGCTTGGC	ATGATGA	GAGCGACTGCA	CTGA	AAAG	TGATGA	CAGGCGCTTG	CTGA	CGGGTATCG
sR18	TTTATTTTA	ATGATGA	AACAGCCAGACC	CTGA	TGGGA	TGATGA	GTGGTGGCTTAG	CTGA	TGTTGCGGTA
sR19	TTGCGGGCTC	ATGATGA	GCTTCCCTACGGC	CCGA	GCTTAGG	CGATGA	GGAATACAGCAGGG	CTGA	TTTTGGTAT
sR20	ATGAATGGC	ATGATGA	GGCTCGATTGG	CTGA	AT	CGATGA	TTGAGAGGACTTG	CTGA	CGGGTGATTA
sR21	TACCATCCG	ATGATGA	GACCGTACTGG	CCGA	AGTA	TGATGA	GCACTCGGTAG	CTGA	GGCTGAAAA
sR22	GGAATCCG	ATGATGA	GAAACGGGTACTG	CCGA	GT	CGATGA	GGAAGAGAA	CTGA	GGAAGAAAA

Picture from Todd Lowe