

NcDNAAlign - a flexible and efficient package for the generation of non- protein coding multiple alignments from genomic sequences



Dominic Rose¹, Jana Hertel², Kristin Reiche¹, Peter F. Stadler^{1,2,3,4}, Jörg Hackermüller^{2,4}

¹Bioinformatics Group, Department of Computer Science, University of Leipzig, Leipzig, Germany; ²Interdisciplinary Center for Bioinformatics, University of Leipzig, Leipzig, Germany; ³The Santa Fe Institute, Santa Fe, New Mexico, USA; ⁴ Fraunhofer Institute for Cell Therapy and Immunology, Leipzig, Germany

At a glance...

Overall goal: Efficiently and customizably generate multiple sequence alignments of non-coding genomic sequences.

Approach:

- Pipeline concept, implemented in PERL, combining custom written and standard bioinformatic software tools
- Only sequence data required as input
- Sequence regions not of interest can be discarded beforehand
- Follows a progressive paradigm starting from pairwise BLAST alignments
- Uses a configurable heuristics for alignment "beautification"
- Proof of principle applications to identification of ultra-conserved elements in nematodes and fishes and ncRNA finding using RNAz in *Gammaproteobacteria*

Results/Conclusions:

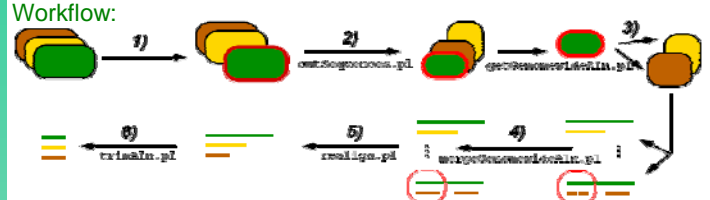
NcDNAAlign...

- Displays comparable sensitivity and specificity as TBA in RNA gene finding using RNAz
- Aligns less DNA than TBA, albeit at higher alignment quality
- Significantly outperforms TBA in computational effort
- Appears to be a reasonable alternative to TBA for pilotstudies, when infrastructure is limited, or when analysis pipelines should be rerun frequently to incorporate newly available data

Availability:

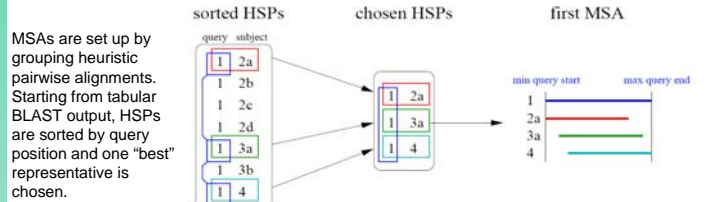
- NcDNAAlign is available under the GNU Public License from <http://www.bioinf.uni-leipzig.de/Software/NcDNAAlign/>.

1) NcDNAAlign Workflow:

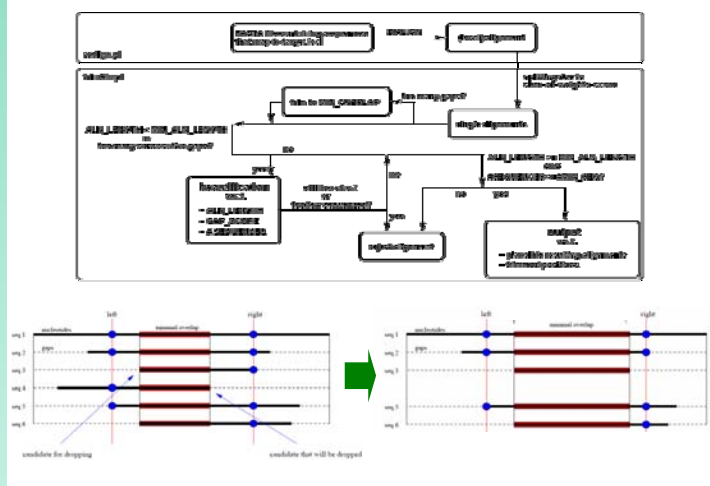


- 1) One species out of all given genomic sequences has to be selected as reference. 2) If selected, sequence data are rided of potentially interfering or uninteresting sequence stretches, reducing the data set to genomic sub sequences. 3) All subsequences of the reference are compared to all subsequences of all other species and local alignments are calculated heuristically (BLAST). 4) Adjacent compatible hits are combined. 5) The best hits (E-value) in each organism for each subsequence of the reference are aligned (DIALIGN). 6) Finally, the alignments are pruned, poorly aligned sequences are removed and the remaining sequences are optionally realigned to obtain an optimal alignment.

Generating multiple alignments:

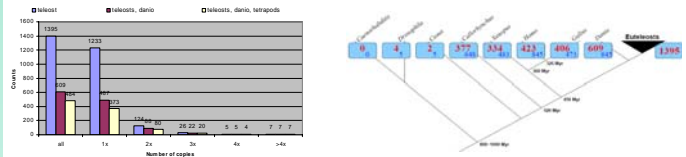


Beautifying Alignments:



2) Results: Teleost UCRs

- Identification of ultra conserved elements (UCE) in teleost genomes using NcDNAAlign
- UCE defined as 100% sequence identity among at least three sequences in alignment and min. 50 nt in length
- Study effect of teleost specific duplication on UCRs

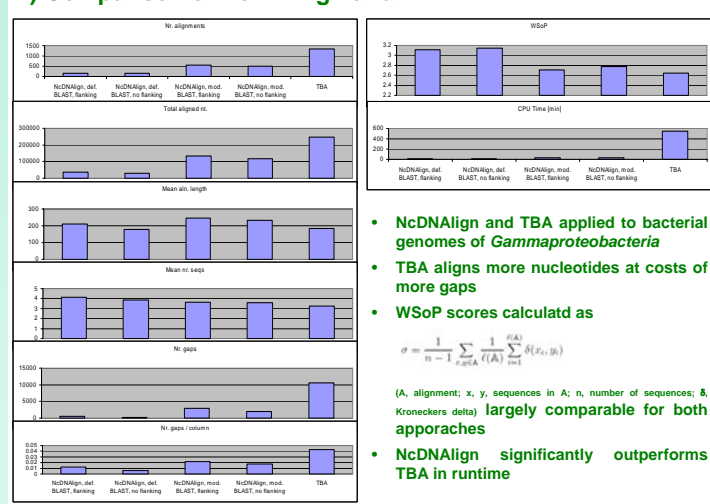


3) Results: Structured ncRNAs in Gammaproteobacteria

Flanking regions	NcDNAAlign default Blast		NcDNAAlign modified Blast		TBA
	+	-	+	-	
	(1)	(2)	(3)	(4)	(5)
(a) CPU time					
Total [min]	15.68	14.35	29.27	27.52	548.36
(b) Alignments					
Nr. alignments	169	169	542	499	1347
Nr. overlapping abs.	155		483		479*
Total align. nucleotides	35 153	29 596	125 907	113 829	235 986
% of E. coli genome	0.76%	0.64%	2.71%	2.45%	5.09%
(c) RNAs					
Nr. hits	126	122	339	300	658
Overlap	99		280		903*
Overall length of hits	25 100	20 618	94 995	80 680	92 888
Mean length of hits	199	169	280	269	141
FDR	0.32	0.33	0.25	0.24	0.26
Nr. annotatable hits	102	80	212	189	469
Nr. non-annot. hits	24	139	24	119	189
(d) Sensitivity					
rRNA, 22 annot.	10	(.46)	9	(.41)	14
tRNA, 86 annot.	55	(.64)	61	(.71)	62
Misc. RNA, 49 annot.	5	(.10)	6	(.12)	18
Overall, 157 annot.	70	(.46)	76	(.50)	94

- Consuming approximately the 20 to 30-fold amount of CPU time, TBA produces roughly 2.5 fold more alignments than NcDNAAlign
- RNAz applied to these alignments yields half the number of hits in NcDNAAlign as compared to TBA at comparable false discovery rates (c).
- Despite the discrepancy in the number of alignments, the sensitivity of the RNAz-based detection of annotated ncRNAs is equal for both approaches.

4) Comparison of NcDNAAlign and TBA



- NcDNAAlign and TBA applied to bacterial genomes of *Gammaproteobacteria*
 - TBA aligns more nucleotides at costs of more gaps
 - WSoP scores calculated as
- $$\sigma = \frac{1}{n-1} \sum_{i=1}^{n-1} \frac{1}{\sum_{j=i+1}^n \delta(r_i, r_j)}$$
- (A, alignment; x, y, sequences in A; n, number of sequences; δ , Kronecker's delta) largely comparable for both approaches
- NcDNAAlign significantly outperforms TBA in runtime