# Promoter Prediction on a Genomic Scale – the Adh Experience

Uwe Ohler

Lehrstuhl für Mustererkennung

(Chair for Pattern Recognition, Computer Science V)

University of Erlangen-Nuremberg

Martensstr. 3, D-91058 Erlangen, Germany

and

Berkeley Drosophila Genome Project

Rubin Lab, Rm 539 LSA MC 3200

University of California at Berkeley, Berkeley, CA 94720

ohler@informatik.uni-erlangen.de

### Abstract

We describe our statistical system for promoter recognition in genomic DNA with which we took part in the Genome Annotation Assessment Project (GASP1). We applied two versions of the system; the first uses a region based approach toward transcription start site identification, namely interpolated Markov chains, the second a hybrid approach combining regions and signals within a stochastic segment model. We compare the results of both versions with each other and examine how well the application on a genomic scale compares to the results we previously obtained on smaller datasets.

## Introduction

Withinthenextyear,thecompletegenomesofseveraleukaryotic organismswillbestoredin thedatabases,andwehavetofacethechallengethattheannotat ionprocessisgettingmore andmorecomplicatedforthegenomicsequenceofhighereukaryotessuch *asD. melanogaster.*Thefirstdraftoftheannotationofanewlysequencedgenomeisusua lly limitedtothecodingpartofagene,butacompleteannotationshouldal socontainthe positionsofthetranscriptionstartsites(TSS),asmostofthe regulatoryelementsinvolvedin geneexpressionarelocatedinthepromoterregionupstreamorcloset otheTSS. Theuntranslatedregionbetweentranscriptionandtranslationstartsi te,the5'UTRregion, canspanuptoseveralkilobasesinhighereukaryotes--itisan averageofalmost2,000bases forthetranscriptionstartsitesetcompiledinthepaperbyRee se *etal.* (2000).Therefore,we cannotsimplytakethesequenceupstreamfromthestartcodon.Methods thataimatthe identificationofregulatoryelementsintheupstreamregionsofc o-expressedgenessuchas vanHelden *etal.* (1998)havebeenshowntodeliverpromisingresultsfortheyeastgenom e whichhasveryshortUTRs,buttheywillbehardtoapplywhentheannot ationonlyconsists ofthecodingpartofagene.Ofcourse,TSSidentificationisal leviatedbyfull-lengthcDNA sequencingprojects;butasthesequencingalwaysstartsatthe3' endofagene,weneed additionalmethodstoconfirmthe5'endofthesequences,ortohuntforra relyexpressed genesthatarenotcontainedinthelibrariesatall.Wearei nadesperateneedtoatleastgeta goodguesswheretheTSS(andthusthepromoterregion)islocated,or wewillstartlooking fortheneedleinthewronghaystack.

TheonlyavailablecomparisonofpromoterpredictioningenomicDNAw ascarriedoutby FickettandHatzigeorgiou(1997).Atthistime,noextensiveunstudied genomicsequences wereavailableforcomplexeukaryoticorganisms,andtheauthorsper formedtheirevaluation onasetof18newlyreleasedvertebratesequences,thelongestof whichcomprisedlessthan6

KB. It was therefore a great challenge to see how well a recently developed promoter recognition program performs on a genomic scale, and what we can conclude for the annotation of complex eukaryotic genomes. We will briefly review the two versions of our promoter recognition system that we applied, discuss in detail the results that were described in the paper of Reese *et al.* (2000), and finally draw conclusions on the state of promoter prediction in general.

### *Methods and Data*

*McPromoter* (Ohler *etal*., 1999a) is a statistical method to look for eukaryotic polymerase II transcription start sites in genomic DNA. It consists of a model for promoter sequences, and a mixture model for non-promoter sequences for coding and non-coding sequences. To localize transcription start sites, a window of 300 bases is shifted over the sequence in steps of 10 bases (see figure 1). At every position, the difference between the log-likelihood of the promoter and the non-promoter model is computed. The resulting plot describes the regulatory potential over the sequence, and is smoothed by a median and hysteresis filter (see Niemann, 1990). The program then makes a prediction for each local minimum below a pre-specified threshold (see figure 2 for an example).

We applied two versions of *McPromoter* on the *Adh* sequence (see Ashburner *etal.* (1999) for a comprehensive description of the annotated genes). The difference between the two versions lies in the structure of the promoter model, and we wanted to explore how well our more recent modeling approach improved on the recognition of TSSs. Version 1.1 of *McPromoter* is a content based approach and uses a single interpolated Markov chain (IMC) of 5[th] order to model promoter sequences. As such, the model does not rely on a priori knowledge about the structure of the promoters, but judges the overall composition of the sequence. For the two non-promoter components for coding and non-coding sequences, we also chose interpolated Markov chains. Related methods were described by Audic and Claverie (1997) and

Hutchinson (1996). In the figures of the GASP paper by Reese *et al*. (2000), version 1.1 is denoted by LMEIMC (L _ehrstuhl für M_uster e_rkenung – I_nterpolated M_arkov C_hains). The submodels are trained using the discriminative Maximum Mutual Information (MMI) approach. In contrast to the standard Maximum Likelihood parameter estimation, MMI maximizes the probability of the decision for the correct sequence class, and therefore also takes negative samples into account (Ohler *et al*., 1999b).

In version 2.0, we replaced the single Markov chain promoter model by a more sophisticated *stochastic segment model* which consists of five states for a simplified upstream-TATA-spacer-initiator-downstream structure of eukaryotic promoters (Ohler *et al*., 2000). With this approach, we obtain more accurate statistics for the states, combining region specific states such as the one for the upstream region with states specific for individual signals such as the one for the TATA box. Hybrid approaches that exploit statistics for several regions were previously described by Solovyev and Salamov (1997) and Zhang (1998). Version 2.0 of *McPromoter* is denoted by LMESSM in the GASP overview paper (Reese *et al*., 2000). Both versions were trained on the same representative dataset consisting of *D. melanogaster* promoter and non-promoter sequences of 300 bases length, obtained at http://www.fruitfly.org/sequence/drosophila-datasets.html. Cross-validation classification experiments on this data (described in Ohler *et al*., 2000) gave a recognition rate of 27.9% for version 1.1 and 58.8% for version 2.0 at the very low false positive rate of 1%. We used the system at this threshold for the evaluation of the *Adh* region.

## Results

According to the results described by Reese *et al*. (2000), version 1.1 of *McPromoter* could identify 26 (28.2%) transcription start sites with a false positive rate of 1/2,633 bases, and version 2.0 successfully located 31 promoters (33.6%) with the slightly higher false positive rate of 1/2,437 bases. This compares well with the results described in the comparison of

promoterrecognitionalgorithmsinvertebrateDNA(FickettandH atzigeorgiou,1997), especiallyconsideredthesmalleramountofavailabletrainingda tafortheorganismof *D. melanogaster*.

Anegativelysurprisingfactforuswasthesmallimprovement oftheperformancethatversion 2.0achievedincomparisonwiththeearlierversion.Withtheresul tsfromcross-validation experimentsontherepresentativesetofpromotersandnon-promotersi nmind,weexpected thenewversiontolocalizeapproximately20-30%moreTSSsatthe samerateoffalse predictions.16ofthe26predictionsmadebyversion1.1arecontainedint hesetof31 predictionsfromversion2.0.Consideredthatthemethodsareclosely related,thisnumberis somewhatsmall,andcouldbeduetothedifferenttrainingalgorithms (MMIversusML parameterestimation).

9predictionsfromversion1.1arelocatedwithin+/-40basesofthes tartsite(meandistance 202bases),asopposedto13closepredictionsandameandistanceof166base softhe predictionsobtainedbyversion2.0.Aswedonotknowexactlyhowfarthe trueTSSdiffers fromourcurrentannotation,thisnumberisencouragingtous.Version2.0 isclearlymore successfultoidentifytheexactpositionofthestartsites.

## Discussion

Togetabetterunderstandingwhytheperformanceofversion1.1andve rsion2.0didnot differverymuchfromeachother,welookedatthesystemperforma ncewithoutthe smoothingpost-processingsteps(table1).Whenwelookattheresults withoutthesmoothing post-processingoperations,itbecomesobviousthatthenewversionindee dmakesgreat improvement,andthatmainlythepost-processingisresponsiblethat version2.0worksless wellthanexpected.Thesmoothingwasdesignedspecificallyforar egion-basedapproachlike theMarkovchainsappliedinversion1.1,andworkslesswellonahybri dapproachlike version2.0wherethepromoterregionisdividedinseveraldistinct segments.

A rough extrapolation of the cross-validation results at the currently used threshold (1% false positives) leads to a worst-case rate of 1/2,000 bases false predictions. From the non-smoothed results it becomes clear now that this is obviously not met by reality. A possible explanation is that the available training data is still not representative enough. It certainly contains too little non-coding data, and the available promoter set has a bias towards house-keeping genes.

We already realized a number of plans to improve the model performance of version 2.0. The first idea was to include reverse sequence models for the non-promoter states, as we can both directions of the sequence independently. It is well known that the reverse sequences of genes still resemble the true genes on the opposite strand, and that the statistics of reverse exon and intron sequences are close to the forward sequence -- hence the problem of shadow gene predictions. Nevertheless, we added two new states for reverse exon and intron sequences to have a more accurate model for the non-promoters.

In a second step, we increased the amount of training data. For the *Adh* experiment, we took the model that performed best on three cross-validation experiments and left out one third of the available data to see if our predictions on this set were met by reality. Instead, we took the whole set and determined the 1% false positive threshold by choosing the mean threshold of the three experiments.

Finally, we replaced the median and hysteresis filters by a simple approach to allow only one prediction below the threshold within 300 bases (the model size). A similar smoothing approach is implicitly carried out by the gene finders with integrated promoter predictors; they choose the best prediction in accordance with the model topology which allows for only one prediction before the start codon. But the question remains if some predictions close to the best one might correspond to alternative transcription start sites, and if such a reduction actually filters out useful information.

As a result of these improvements, 20 predictions instead of 13 are now located within +/-40 bases from the putative start site, and we could increase the performance to 34 identified promoters with a false positive rate of 1/3,000 bases.

## ConclusionsandOutlook

The analysis of the *Adh* region showed us clearly that promoter recognition by itself, without context information, still delivers too many false positives to be practically useful on a *genomic* scale. There is still a lot of room for improvement – we think of parallel states for the TATA box region and the downstream region, discriminative training of the segment model, and a non-linear combination of the segment likelihoods. But the overall picture will may be not change in the near future when we exploit only the primary sequence. We will see if the usage of other features such as DNA bendability (Pedersen *et al.*, 1998) can lead to the necessary improvement.

From a different point of view though, the rate of one false positive in three kilobases seems reasonable if one has already an idea where the coding part of the gene is. This information can be provided in both by alignments of cDNA to genomic sequence and abinitio gene finding. We therefore envision a promoter recognition system used within a gene finder that also incorporates EST and cDNA alignment information to extend the coding region on the 5'end. The accuracy of the TSSl localization of *McPromoter* is good enough to then use such a preliminary annotation of the transcription start site for the analysis of upstream regions of co-expressed genes.

Both versions of the McPromoter system can be accessed via the World Wide Web at http://www5.informatik.uni-erlangen.de/HTML/English/Research/Promoter.
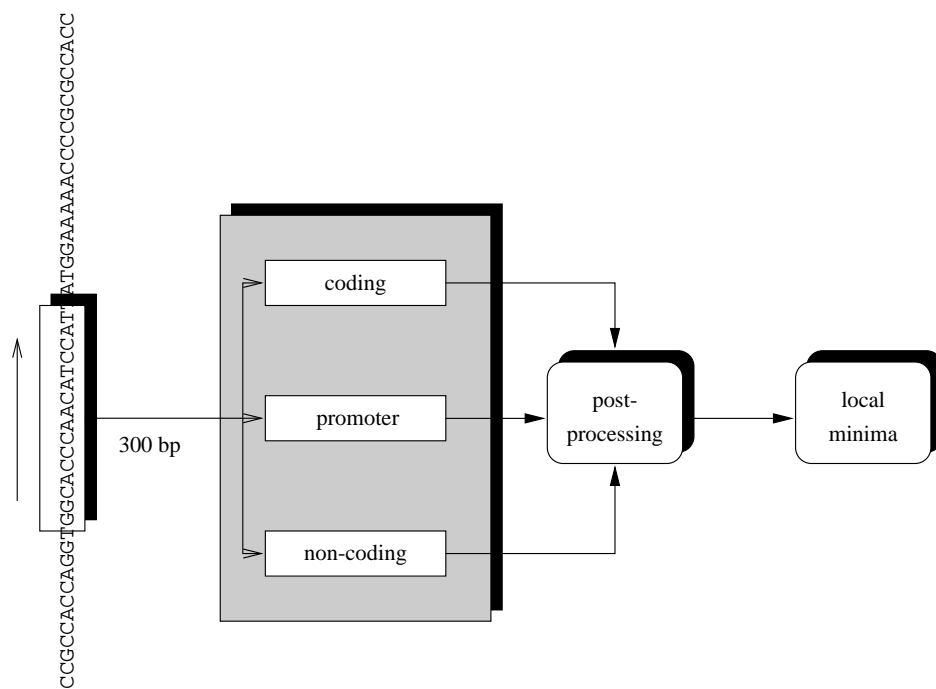
GeorgeHartzellandMartinReesefortheworkonthecollectiona ndevaluationofputative

TSSsinthe *Adh*region,andtoG.Rubin,theheadoftheBerkeleyDrosophilaGenome

Project.

### References

Ashburner,M.,S.Misra,J.Roote,S.E.Lewis,R.Blazej ,T.Davis,C.Doyle,R.Galle,R.

George,N.Harris,G.Hartzell,D.Harvey,L.Hong,K. Houston,R.Hoskins,G.Johnson,C.

Martin,A.Moshrefi,M.Palazzolo,M.G.Reese,A.Spradli ng,G.Tsang,K.Wan,K.

Whitelaw,B.Kimmel,S.CelnikerandG.M.Rubin.1999.Anexplorat ionofthesequenceof

a2.9-MbregionofthegenomeofDrosophilamelanogaster:TheAdhre gion.Genetics

153(1):179-219.

Audic,S.andJ.-M.Claverie.1997.Detectionofeukaryoticpromoter susingMarkov

transitionmatrices.ComputChem.21(4):223-7.

Fickett,J.andA.Hatzigeorgiou.1997.Eukaryoticpromoterrecognition. GenomeRes.7:861-

878.

Hutchinson,G.B.1996.Thepredictionofvertebratepromoterregionsusi ngdifferential

hexamerfrequencyanalysis.ComputApplBiosci.12(5):391-8.

Niemann,H.1990.PatternAnalysisandUnderstanding,2 [nd]edition.Springer,Berlin.

Ohler,U.,S.Harbeck,H.Niemann,E.Nöth,andM.G.Reese. 1999a.InterpolatedMarkov

chainsforeukaryoticpromoterrecognition.Bioinformatics15(5):362-369.

Ohler,U.,S.HarbeckandH.Niemann.1999b.Discriminativetrai ningoflanguagemodel classifiers.Proc.EuropeanConferenceonSpeechandSignalProcess ingTechnology,p. 1607-1610,Budapest.

Ohler,U.,S.Harbeck,G.Stemmer,andH.Niemann.2000.Stochas ticsegmentmodelsof eukaryoticpromoterregions.PacificSymposiumonBiocomputing5:377-388.

Pedersen,A.G.,P.BaldiP,Y.Chauvin,S.Brunak.1998.DNAst ructureinhumanRNA polymeraseIIpromoters.JMolBiol.28;281(4):663-73.

Reese,M.G.,N.Harris,G.Hartzell,U.Ohler,andS. Lewis.2000.Thegenomeannotation assessmentproject.GenomeRes.10,toappear.

Solovyev,V.andA.Salamov.197.TheGene-Findercomputertoolsfor analysisofhuman andmodelorganismsgenomesequences.ProcISMB5:294-302.

VanHelden,J.,B.AndreandJ.Collado-Vides.1998.Extractingregul atorysitesfromthe upstreamregionofyeastgenesbycomputationalanalysisofoligonucl eotidefrequencies.J. Mol.Biol.281(5):827-842.

Zhang,M.Q.1998.Identificationofhumangenecorepromotersinsilic o.GenomeRes8(3): 319-326.

coding

promoter

300 bp

non-coding

post-
processing

local
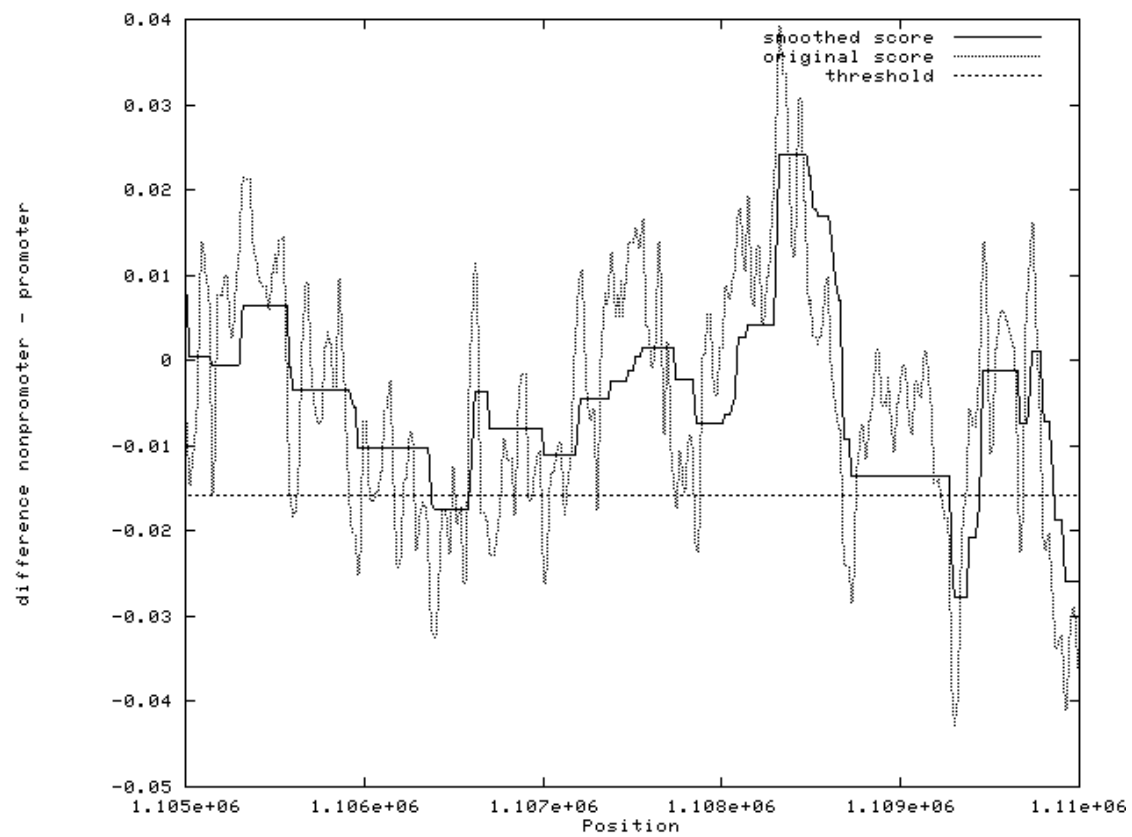minima

CCGCCACCAGGTGGCACCAACATCCATTATGGAAAACCCCGCGCCACC

**Figure 1. Structure of the McPromoter system. A window of 300 bases is shifted over the sequence in steps of 10 bases, and the content is evaluated with the promoter and non-promoter models. The difference between the promoter and the non-promoter log likelihood is stored. After post-processing, the local minima are reported as transcription start site predictions.**

**Figure 2. Application of McPromoter version 2.0 on a 5kB part of the *Adh* region containing the transcription start site for the *Adh* gene. We show the non-smoothed as well as the smoothed output of the system. The strongest local minimm corresponds to the annotated transcription start site of *Adh*.**

| Post-processing | Version1.1 | | Version2.0 | |
|---|---|---|---|---|
| | *Recognized promoters* | *Falsepositive rate* | *Recognized promoters* | *Falsepositive rate* |
| *None* | 47 | 1/450 | 57 | 1/719 |
| *Hysteresis* | 33 | 1/1,833 | 43 | 1/1,653 |
| *Median&Hysteresis* | 26 | 1/2,633 | 31 | 1/2,437 |

**Table1.Comparisonoftheinfluenceofpost-proces        singontheperformanceofthepromoterpredictors.**

**Shownaretheresultswithoutanypost-processing(        i.e.,everylocalminimumisusedasprediction),a        fter**

**hysteresissmoothing,andafterbothmedianandhys        teresissmoothing.Thepost-processingoperations**

**reducethenumberoffalsepositivesforbothvers        ions,butitbecomesclearthattheeffectismuch        better**

**forthepureregion-basedapproachofversion1.1.**