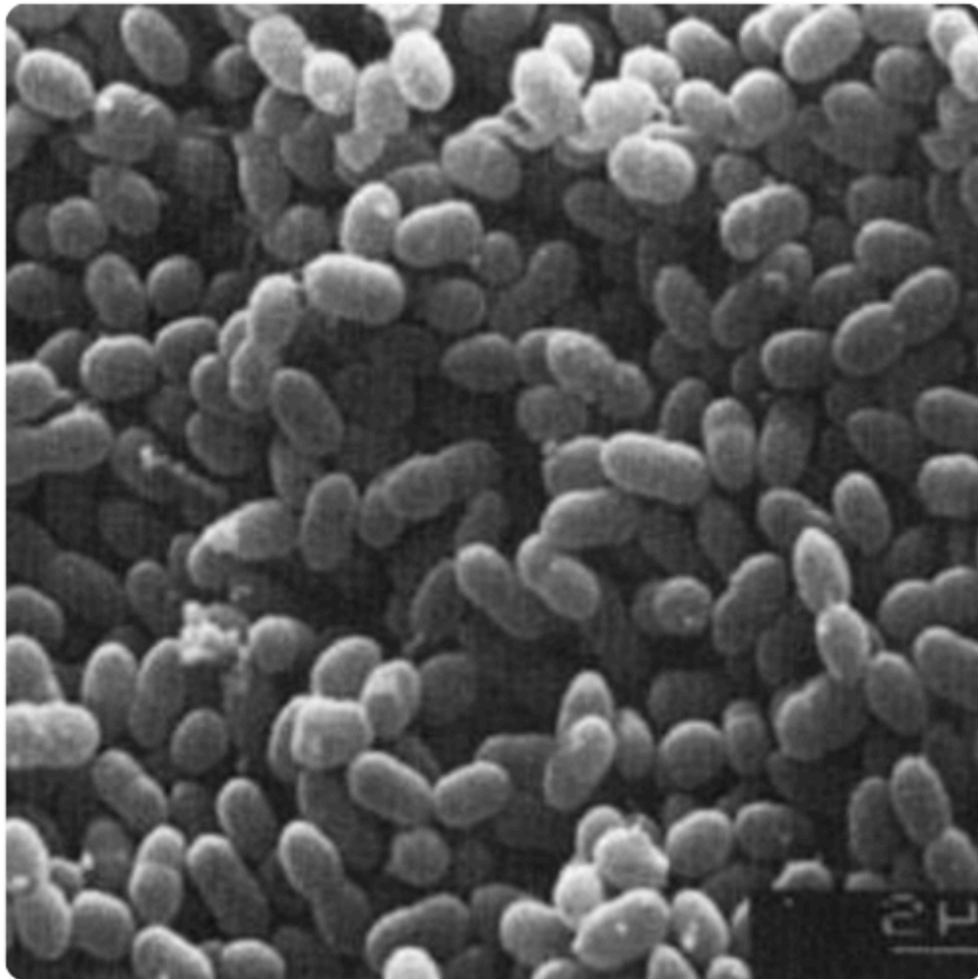


Annotation du génome de *Lactococcus lactis*

Bioinformatique pour la génomique

Rapport de projet
Effectué dans le cadre du parcours
Master 1 Bioinformatique et Biologie des Systèmes
Année universitaire 2023-2024



Morphologie en microscopie électronique de *Lactococcus lactis* subsp. (1)

RÉSUMÉ

Lactococcus lactis est une bactérie gram-positive jouant un rôle central dans les fermentations laitières et la production de levains, et constitue un modèle clé en recherche scientifique grâce à son génome séquencé. Cette étude se concentre sur la réannotation d'un fragment de son génome, en combinant des outils bioinformatiques avancés tels qu'ORFfinder, GeneMark, GeneMarkHMM, ScanForMatches, BlastP, SignalP et DeepTMHMM. L'objectif est de caractériser les régions codantes (CDS) en identifiant les cadres de lecture ouverts (ORF), les codons d'initiation et de terminaison, ainsi que les motifs fonctionnels associés, incluant les sites de fixation du ribosome (RBS), les promoteurs sigma, et les terminateurs rho-dépendants.

Les analyses ont permis d'identifier sept protéines hypothétiques, dont deux transmembranaires, et de prédire leurs fonctions par homologie de séquences et leur localisation subcellulaire. Les résultats mettent en évidence la robustesse de l'approche multi-outil pour une caractérisation précise des gènes et de leurs produits, malgré certaines limites méthodologiques, comme l'utilisation d'algorithmes initialement optimisés pour d'autres génomes, notamment celui de *E. coli*.

Des perspectives d'amélioration incluent l'entraînement spécifique des algorithmes aux données génomiques de *L. lactis* et l'intégration de méthodes élargies pour explorer les mécanismes de terminaison alternatifs. Ces travaux apportent une contribution significative à la compréhension moléculaire de cet organisme d'intérêt industriel et scientifique, consolidant son rôle central dans la production alimentaire et la biotechnologie.

Points clés

- *Lactococcus lactis* est crucial pour les fermentations laitières et constitue un modèle en recherche scientifique.
- Cette étude porte sur l'annotation bioinformatique d'un fragment du génome de *L. lactis*.
- Les outils bioinformatiques utilisés incluent ORFfinder, GeneMark, GeneMarkHMM, ScanForMatches, BlastP, SignalP et DeepTMHMM.
- Identification des CDS via l'analyse des ORF, des codons start/stop, et des motifs fonctionnels (RBS, sites sigma, terminateurs rho-dépendants).
- Sept protéines hypothétiques, dont deux transmembranaires, ont été identifiées et caractérisées.
- Des limites méthodologiques sont relevées, notamment l'utilisation d'algorithmes entraînés sur *E. coli*.
- Des améliorations futures prévoient un entraînement spécifique aux données de *L. lactis* et une exploration étendue des mécanismes de terminaison.
- Ces travaux enrichissent la compréhension moléculaire et appliquée de *L. lactis*.

Mots clés

Lactococcus lactis, bactérie gram-positive, fermentations laitières, levains, annotation génomique, bioinformatique, ORFfinder, GeneMark, GeneMarkHMM, ScanForMatches, BlastP, SignalP, DeepTMHMM, régions codantes (CDS), ORF (Open Reading Frames), protéines hypothétiques, homologie de séquences, localisation subcellulaire, améliorations méthodologiques

Les scripts utilisés et les analyses ont été reportés dans le dossier GitHub :
https://github.com/CamilleAstrid/Annotation_du_genome_de_Lactococcus_lactis

INTRODUCTION

Lactococcus lactis est une bactérie gram-positive largement répandue dans les environnements laitiers, notamment dans le lait cru, où elle joue un rôle clé dans la fermentation. Utilisée depuis des siècles dans la fabrication de fromages et autres produits laitiers, cette bactérie se distingue par son métabolisme anaérobie aérotoleérant et sa capacité à croître de manière optimale à une température d'environ 30 °C. Deux sous-espèces principales de *L. lactis* sont couramment employées dans la production de levains, ce qui en fait un micro-organisme d'intérêt industriel majeur.

Au-delà de son utilisation industrielle, *L. lactis* constitue également un modèle pour la recherche scientifique. Elle est la première bactérie lactique dont le génome a été entièrement séquencé, ouvrant la voie à une meilleure compréhension des mécanismes moléculaires et des fonctions biologiques associées à ce groupe de micro-organismes.

Au cours de cette étude, nous allons reproduire l'annotation d'une partie du génome de *L. lactis* (Annexe1) en utilisant des outils bioinformatiques tels qu'ORFfinder (2), GeneMark (GM) (3), GeneMarkHMM (GMH) (4) et ScanForMatches (SFM) (5). Les objectifs de cette étude incluent l'analyse des cadres de lecture ouverts (ORF) à laquelle viennent s'ajouter la recherche de codons d'initiation (codon START) et de signal de terminaison (codon STOP). L'ensemble de ces analyses nous amèneront à considérer des régions codantes (CDS). Nous y adosserons l'identification de la présence de séquences spécifiques telles que les sites de fixation du ribosome (RBS), du facteurs sigma et la présence de terminateurs rho-dépendants. Une fois les CDS identifiés, nous déterminerons la fonction des protéines issues de ces gènes par analogie de séquence (BlastP), ainsi que leur localisation subcellulaire (SignalIP et DeepTMHMM). Afin de permettre une meilleure analyse des résultats obtenus par les différents logiciels, nous développerons un programme permettant la conversion des fichiers de sortie de GM, GMH et SFM au format GFF (general feature format). Il s'agit d'un format de fichier utilisé pour décrire les gènes et d'autres éléments de séquences d'ADN, d'ARN et de protéines.

MATÉRIELS ET MÉTHODES

Quels sont les logiciels et langages utilisés pour cette étude ?

ORFfinder (Open Reading Frame Finder) (2) recherche des ORF dans la séquence d'ADN saisie, renvoyant les coordonnées de chaque ORF ainsi que la traduction de leur séquence en protéines. Ce programme est utilisé pour identifier dans un ADN nouvellement séquencé des segments potentiellement codants pour des protéines.

GeneMark (3) est un algorithme bayésien basé sur les modèles de chaînes de Markov dépendant du cadre (non homogènes) des régions codantes pour les protéines. Il est capable d'identifier le brin d'ADN codant et le cadre de lecture correct, tout en générant un signal distinctif pour localiser un gène. Pour obtenir de bons résultats, la séquence à analyser doit être extraite de la même population statistique que l'ensemble de formation. Ainsi, on ne peut pas s'attendre à ce que l'algorithme entraîné sur le génome d'*E. coli* fonctionne correctement sur les séquences provenant du génome d'autres espèces (telles que *L. lactis*).

GeneMarkHMM (4) est un algorithme utilisé pour améliorer la qualité de la prédiction des gènes en identifiant avec précision les limites des gènes. Il intègre les modèles de GM dans un cadre basé sur les modèles de Markov cachés (HMM), avec les limites des gènes modélisées comme des transitions entre états cachés. Pour affiner les prédictions des codons d'initiation de la traduction, il utilise des motifs spécifiques des sites de liaison au ribosome. GMH est significativement plus précis que GM pour la prédiction exacte des gènes.

ScanForMatches (5) est un outil en C conçu pour rechercher des motifs dans des séquences d'ADN ou de protéines au format FASTA. Il inclut des fonctionnalités pour identifier des motifs complexes, gérer des appariements non standards, et tolérer des variations comme des décalages, des insertions ou des suppressions. Les motifs sont définis par des règles et permettent des recherches basées sur des plages, des appariements inversés ou des boucles (comme les structures en épingle à cheveux). La recherche peut s'effectuer sur les deux brins. Cet algorithme permet également la recherche de répétitions, de motifs spécifiques mais aussi l'utilisation de matrices de poids position.

Python (6) est un langage de programmation créé par Guido van Rossum (1991). Le nom du langage vient de la série humoristique "Monty Python's Flying Circus". Ce langage met l'accent sur la lisibilité du code, en utilisant une syntaxe simple et minimaliste qui permet aux développeurs de se concentrer sur la logique de leur programme. Langage interprété et multiplateforme, Python supporte plusieurs paradigmes de programmation, notamment procédural, orienté objet et fonctionnel, et s'accompagne d'une riche bibliothèque standard pour répondre à divers besoins. Il permet ainsi un large éventail d'applications telles que le développement web, l'analyse des données, l'intelligence artificielle et l'automatisation.

Perl (Practical Extraction and Report Language) (7) est un langage de programmation créé par Larry Wall (1987). Ses principales caractéristiques sont sa flexibilité, sa capacité à gérer des expressions régulières puissantes, sa syntaxe concise et son interprétation dynamique. Perl est principalement conçu pour le traitement de texte et l'extraction de données. Il offre également une portabilité multiplateforme et bénéficie d'une grande quantité de modules. Il est donc largement utilisé dans l'administration système, la bioinformatique, le développement web et l'automatisation de tâches. Cependant, Perl peut être difficile à lire en raison de sa syntaxe flexible, et bien qu'il soit rapide pour certaines tâches, il peut être moins performant que des langages compilés. Bien qu'il ait perdu de sa popularité au profit de langages comme Python, il reste pertinent pour des tâches spécifiques dans des environnements Unix/Linux.

Quelles sont les étapes mises en place pour cette étude ?

Identifier les régions codantes

Nous avons d'abord identifié les **ORF** éventuels avec ORFfinder. Nous avons sélectionné le code génétique Standard et restreint l'analyse à des tailles d'ORF de 300 nucléotides. " 'ATG' et codon d'initiation alternatif" a été sélectionné. En effet, chez les procaryotes notamment, il y a différents codons start possibles.

Pour mieux définir les **CDS** présents sur notre fragment génomique d'intérêt, nous avons utilisé deux algorithmes : GM et GMH. GM nous permet de localiser les **codons start et stop** éventuels ainsi que les **régions** probablement **codantes**. GMH permet de vérifier la présence de **gènes typiques ou atypiques** (transfert horizontal) en plus des informations fournies par GM.

Les paramètres utilisés pour GM sont respectivement pour l'espèce, la taille de la fenêtre, le pas d'étude et le seuil : *Lactococcus lactis* Il1403 ; 120 ; 12 et 0,5. Une seconde étude avec GM a été menée en modifiant le seuil à 0,4. Concernant, GMH, nous avons sélectionné la même espèce que GM afin d'avoir des résultats comparables.

L'utilisation des deux algorithmes vient conforter la présence de CDS éventuels. Cependant, ces deux algorithmes ont été entraînés sur le génome de *E. coli*, il nous faut donc de nouvelles approches pour venir appuyer la présence d'un gène codant pour une protéine.

Pour cela, nous avons choisi de rechercher la présence de **RBS**, de **site de fixation du facteur sigma** et des **sites de terminaison rho-dépendant**. Ces trois analyses ont été effectuées en ligne de commande avec l'outil SFM (Annexe2) et les matrices et pattern correspondants (Annexe3).

Définir la fonction et la localisation des protéines issues des régions codantes identifiées

Une fois les CDS identifiés, nous avons récupéré les séquences protéiques à l'aide de la suite EMBOSS. Nous avons **découpé la séquence** pour en extraire les CDS avec Extractseq et **traduit les séquences** génomiques en séquences protéiques avec Transeq. Les séquences protéiques hypothétiques obtenues ont ensuite été utilisées comme query d'un BlastP. Cette étape permet d'identifier les **fonctions des protéines** par analogie de séquences.

La localisation des protéines a été étudiée par la recherche de **peptide signal** sur la protéine qui indiquerait un adressage à la membrane (protéine exportée ou sécrétée) avec l'outil SignalIP. Ensuite, une analyse menée avec DeepTMHMM a prédit la **position subcellulaire** de la protéine (intracellulaire, membranaire ou extracellulaire) ainsi que sa **structure tertiaire** venant appuyer les informations obtenues sur sa localisation subcellulaire.

RÉSULTATS

Les résultats de ORFfinder (Annexe4), GM (Annexe5), GMH (Annexe6), SFM (Annexe7), SignalIP (Annexe8) et DeepTMHMM (Annexe9) ont été retranscrits dans les tableaux ci-dessous.

ORF	ORFfinder					
	Label	Strand	Frame	Start	Stop	Length (nt aa)
0	/	/	/	/	/	/
3	ORF3	+	2	455	3436	2982 993
4	ORF4	+	3	3447	5249	1803 600
1	ORF1	+	1	5230	6486	1257 418
5	ORF5	+	3	6483	7220	738 245
2	ORF2	+	1	7261	9270	2010 669
6	/	/	/	/	/	/

Tableau A : Résultats obtenus avec ORFfinder

Figure 1 Analyse des régions codantes et des protéines issues.
Identification des régions codantes par la présence de codon initiateur de la transcription, codon stop, RBS, site de fixation du facteur sigma et site de terminaison rho-dépendant. La probabilité qu'une région soit codante ajoutée aux données précédentes permet la délimitation de gènes. L'étude est ensuite poursuivie par l'analyse fonctionnelle des protéines hypothétiques générées par homologie de séquences.

ORF	GeneMark (seuil 0.5)				GeneMark (seuil 0.4)				GeneMarkHMM				
	LEnd	REnd	Strand	Frame	LEnd	REnd	Strand	Frame	Strand	LeftEnd	RightEnd	GeneLength	Class
0	/	/	/	/	/	/	/	/	+	174	347	174	2
3	425	3436	direct	fr 2	425	3436	direct	fr 2	+	455	3436	2982	2
4	3444	5249	direct	fr 3	3444	5249	direct	fr 3	+	3447	5249	1803	2
1	5185	6486	direct	fr 1	5185	6486	direct	fr 1	+	5242	6486	1245	2
5	6453	7220	direct	fr 3	6453	7220	direct	fr 3	+	6486	7220	738	2
2	7210	9270	direct	fr 1	7210	9270	direct	fr 1	+	7222	9270	2049	1
6	9388	9558	direct	fr 1	/	/	/	/	+	9475	>9555	81	1

Tableau B : Résultats obtenus avec GeneMark et GeneMarkHMM

<i>ORF</i>	<i>RBS</i>	<i>sigma</i>	<i>rho</i>	<i>Blastp</i>
0	161-176	111-151	557-582	<i>gallidermin/nisin_family_lantibiotic</i>
3	282-298	399-435	3986-4009	<i>ABC_transporter_ATP-binding_protein</i>
4	2531-2549	3294-3333	6683-6717	<i>lantibiotic_dehydratase</i>
1	4581-4599	4317-4348	6683-6717	<i>Nisl/Spal_family_lantibiotic_immunity_lipoprotein</i>
5	6469-6485	6344-6389	7224-7248	<i>lanthionine_synthetase_C_family_protein</i>
2	7205-7224	6873-6916	9277-9310	<i>Nisl/Spal_family_lantibiotic_immunity_lipoprotein</i>
6	9397-9411	9307-9353	/	<i>response_regulator_transcription_factor</i>

Tableau C : Résultats obtenus avec ScanForMatches et BlastP (NCBI)

<i>ORF</i>	<i>SignalIP</i>					<i>DeepTMHMM</i>	
	<i>Prediction</i>	<i>SP(Sec/SPI)</i>	<i>TAT(Tat/SPI)</i>	<i>LIPO(Sec/SPII)</i>	<i>OTHER</i>	<i>Type</i>	<i>Localisation</i>
0	OTHER	0.217604	0.004405	0.009520	0.768471	Globulaire	Inside
3	OTHER	0.125556	0.002389	0.016899	0.855156	alpha-TM	Transmembranaire
4	OTHER	0.003432	0.000448	0.001506	0.994614	Globulaire	Inside
1	OTHER	0.019379	0.002161	0.003979	0.974481	alpha-TM	Transmembranaire
5	OTHER	0.017380	0.000367	0.021509	0.960745	Globulaire	Inside
2	OTHER	0.004834	0.000481	0.000758	0.993927	Globulaire	Inside
6	OTHER	0.001969	0.000097	0.000325	0.997610	Globulaire	Inside

Tableau D : Résultats obtenus avec SignalIP et DeepTMHMM

Afin de mieux comparer les résultats, un programme pour convertir les fichiers a été créé et utilisé (documents complémentaires 4 et 5). Les analyses montrent la présence supposée de sept protéines dont deux transmembranaires. Les protéines 0 et 6 ne sont pas identifiées avec ORFfinder puisque leur taille est inférieure à celle de la recherche. En effet, elles seraient tronquées et donc les analyses sur leur promoteur et terminateur seraient approximatives. Tous les gènes sont dans le même sens (direct). La protéine 3 est une ABC transporter et est bien transmembranaire à trois domaines transmembranaires et une partie extracellulaire. La protéine 1 est une lipoprotéine responsable d'une réponse immunitaire et est membranaire. Cependant, ces résultats ne se retrouvent pas dans l'analyse faite par SignalIP.

CONCLUSION & DISCUSSION

L'étude menée sur le génome de *Lactococcus lactis* a permis d'explorer diverses méthodes de prédiction des régions codantes et d'analyse des fonctions protéiques. L'utilisation combinée de plusieurs outils bioinformatiques, tels qu'ORFfinder, GeneMark, GeneMarkHMM, ScanForMatches, SignalIP et DeepTMHMM, a fourni une vue d'ensemble des gènes potentiels et de leurs caractéristiques. Ces approches ont permis d'identifier sept protéines hypothétiques, dont deux transmembranaires.

Cependant, des limites ont été constatées, notamment l'utilisation de logiciels initialement développés pour d'autres organismes, comme *E. coli*, ce qui peut altérer la précision des prédictions pour *L. lactis*. De même, la recherche s'est concentrée uniquement sur les terminateurs rho-dépendants, négligeant d'autres

types de terminaison, et certains promoteurs sigma n'ont pu être analysés en raison de problèmes techniques.

Pour les études futures, l'amélioration de la pertinence des outils bioinformatiques via un entraînement spécifique au génome de *L. lactis*, ainsi que l'intégration d'analyses complémentaires pour détecter d'autres types de terminateurs, seront nécessaires. Ces ajustements permettraient de renforcer la fiabilité et la robustesse des résultats, contribuant ainsi à une compréhension plus approfondie de cet organisme d'intérêt industriel et scientifique.

RÉFÉRENCES

- [1] Mémoire de Magistère : Isolement et sélection des souches de bactéries lactiques productrices des métabolites antibactériennes, par BELARBI Fatima (2010-2011) doi:10.13140/RG.2.2.13373.82405
- [2] Wheeler, D L et al. "Database resources of the National Center for Biotechnology Information." Nucleic acids research vol. 28,1 (2000): 10-4. doi:10.1093/nar/28.1.10
- [3] Borodovsky M. and McIninch J. "GeneMark: parallel gene recognition for both DNA strands." Computers & Chemistry, 1993, Vol. 17, No. 19, pp. 123-133
- [4] Lukashin, A V, and M Borodovsky. "GeneMark.hmm: new solutions for gene finding." Nucleic acids research vol. 26,4 (1998): 1107-15. doi:10.1093/nar/26.4.1107
- [5] <https://blog.theseed.org/servers/2010/07/scan-for-matches.html> By The SEED Team on July 16, 2010. The utility was written by Ross Overbeek; David Joerg and Morgan Price wrote sections of an earlier version. It is worth noting that it was strongly influenced by the elegant tools developed and distributed by David Searls.
- [6] Van Rossum, G., & Drake, F. L. (2009). Python 3 Reference Manual. Scotts Valley, CA: CreateSpace.
- [7] Wall, L., & others. (1994). The Perl programming language. Prentice Hall Software Series.

ANNEXES

>seq L_lactis_lactis

GAAACATTAACAAATCTAAAAACAGTCTTAACCCCGGATGAGAATTCGGGGGTTTTTTCTATTGGTAATATAACGCGAGCATAATAAACGGCTCTGATT
AAATTCGTAAGTTGCCATCAATGCTCAGGACGAATATTGCATGATCTATCATAAATTATAAGGAGGCACTCAAAATGAGTACAAAAGATTTTAACTTGG
ATTTGGTATCTGTTTCGAAGAAAGATTCAAGGTGCATCACCACGCATTACAAGTATTTTCGTATGTACACCCGGTTGTAACAGGAGCTCTGATGGGT
TGTAACATGAAAACAGCAACTTGTCAATTGTAGTATTCACGTAAGCAAATAACCAAAATCAAAGGATAGTATTTTGTAGTTCAGACATGGATACTATCCT
ATTTTTATAAGTTATTTAGGGTTGCTAAATAGCTTATAAAAAATAAAGAGAGGAAAAACATGATAAAAAAGTTTCATTTAAAGCTCAACCGTTTTAGTAA
GAAATACAATTTTATCTCCAAACGATAAACGGAGTTTACTGAATATACTCAAGTCATTGAGACTGTAAGTAAAAATAAAGTTTTTTTGGAACAGTTAC
TACTAGCTAATCCTAAACTCTATGATGTTATGCAGAAATATAATGCTGGTCTGTAAAGAAGAAAAAGGGTTAAAAAATTATTTGAATCTATTTACAAGT
ATTATAAGAGAAGTTATTTACGATCAACTCCATTTGGATTATTTAGTGAACCTCAATTGGTGTTTTTCGAAAAAGTTCACAGTACAAGTTAATGGGAA
AGACTACAAAGGGTATAAGATTGGATACTCAGTGGTTGATTGCGCTAGTTCATAAAATGGAAGTAGATTCTCAAAAAAGTTATCATTTACTAGAAATA
ATGCAAATTATAAGTTTGGAGATCGAGTTTTTCAAGTTTATACCATAAATAGTAGTGAGCTTGAAGAAGTAAATATTAAATATACGAATGTTTATCAAAT
TATTTCTGAATTTTGTGAGAATGACTATCAAAAATATGAAGATATTTGTGAAACTGTAACGCTTTGCTATGGAGACGAATATAGAACTATCGGAACA
ATATCTTGGCAGTCTGATAGTTAATCATTTATGATCTCTAATTTACAAAAGATTGTGTGTCAGATTTTTCTTGGAACACTTTTTGACTAAAAGTTGAA
GCAATAGATGAAGATAAAAAATATATAATCTCTGAAAAAAGTTCAAAAAGTTTATTCAAGAATACTCAGAAATAGAAATTTGGTGAAGGTATTGAGAA
ACTGAAAGAAATATATCAGGAAATGTCACAAATCTTGAGAATGATAATTATATTTCAAATTTGATTAAATAGTGATAGTGAAATAAATTTGATGTTAAA
CAAAAGCAACAATTAGAACATTTAGCTGAGTTTTTAGGAAATACGACAAAATCTGTAAAGAAACATATTTGGATGACTATAAGGATAAATTTATCGA
AAAATATGGTGTAGATCAAGAAGTACAAATAACAGAATTATTTGATTCTACATTTGGCATAGGAGCTCCATATAATTATAATCATCCTCGAAATGACTTT
CTGATGTCGACCAAGTATGATATCTCAGAAAGGAGAGAGAGAAAAAGTTCAAAAAGTTTATTCAAGAATACTCAGAAATAGAAATTTGGTGAAGGTATTGATA
TCTTGACGACTTAGAGTCTCATTATCAAAAAATGGACTTAGAAAAGAAAAAGTGAACCTCAAGGGTTAGAATTATTTTTGAATTTGGCAAAGGAGTAT
GAAAAAGATATTTTTATTTTAGGGGATATCGTTGGAAATAATAATTTGGGAGGGGCATCAGGTAGATTTTCTGCACTCTCTCCGGAGTTAACAGTTAT
CATAGAACGATAGTAGATTCTGTGCGAAAGAGAAAAATGAGAATAAAGAAATACATCGTGTGAAATAGTATTTCTCCAGAAAAATACAGACATGCTAA
CGTTAGTACATACATAAATTATGAGAGGAAAGTACTTCCATTTTTACAAGTCAAGTACAAATGAAGTTCTGTAACTAATATCTATATTGGAATAGAC
GAAAAAGAAAAATTTATGCACGAGACATTTCAACTCAAGAGGTATTGAAATTTCTACATTACAAGCATGTACAATAAAACGTTATTCAGTAATGAGCT
AAGATTTCTTTACGAAATTTCAATAGATGACAAAGTTTGGTAATTTACCTTGGGAACCTATTTACAGAGACTTTGATTATATTCCAGTTTAGTATTGAC
GAAATAGTAAATATCTCTGTAAATGGAATTTGGGGAAGGGATGTAATAGTACAATAGACAATAAGAGAACCTTATTTCAAAGCAAAAGAAATTTCCCA
AAGAGTTTATATTGTCAATGGAGATAATAAGTTTATTTATCAGAGAAAAACCCATTTGGATGGAATTTAGAGTCGGCGTAAAGAAGAGCTCA
AAAAAGAAAAAGATTTTATAGAGCTACAAGAATATTTTGAAGATGAAAAATATCATAAATAAAGGAGAAAAAGGGGAGAGTTGCCGATGTTGTAGTGCCTT
TTATTAGAACGAGAGCATTAGGTAATGAAGGGAGAGCATTATAAGAGAGAAAAAGAGTTTCGGTTGAACGGCGTGAAAAATTGCCCTTTAACGAGT
GGCTTTATCTAAAGTTGTACATTTCTATAAATCGTCAAAATGAATTTTACTGTCTGATCTTCCAGATATTCAGAAAAATAGTAGCAAAACCTGGGTGGAA
ATCTATTTCTCTAAGATATATGATCCTAAACCACATATTAGATTGCGTATAAATGTTTCAGATTTATTTTTAGCTTACGGATCTATTTCTGAAATCTTAA
AAAGGAGTCGGAATAATAGGATAATGTCAACTTTTGATATTTCTATTTATGATCAAGAGTAGAAAAGATATGGTGGATTGTGATACTTTAGAGTTATCCG
AAGCAATATTTGTGCCGATTCTAAAATTTATCCAAATTTGCTTACATTGATAAAAGATACTAATAATGATTGGAAGTCGATGATGTATCAATCTTGGT
GAATATTTATATCTGAAATGCTTCTTTCAGAAATGATAACAAAAAGATTCTTAATTTTTGAATTTAGTTAGTACTAAAAAGGTTAAAGAAAAATGTCAT
GAAAAAGATTGAACATTATCTAAGCTTCTGAAAGTTAATACTAGGTGACCAAAATTTTTATGACAAGAATTTTAAAGAAATTAAGCATGCCATAAA
AATTTATTTTTAAAAATGATGCTCAAGATTTTGAACCTTCAGAAAGTTTATTCAAATTTAGTCAAGTATGATTCATGTCCTAATCAATCAATAGTTGGTAT
TGAACGAGATAAAAGAGAAATTAATTTATACACACTTCAAAGGTTGTTTGTTCGGAAGAATACATGAAATGAGGACTAATAGATGGATGAAGTGAA
AGAATTCACATCAAAACAATTTTTAATACTTTACTTACTCTTCCAAGCACCTTGAAGTTAATTTTCAGTTGGAAAAACGTTATGCAATTTATTTAATT
GTGCTAAATGCTATCAGAGCTTTTGTTCGGTTGGCTAGTCTTTTATTTATCAAGATTAAATAAATCTGTGCTAGGTTACAGGAGACATCTTATCAATA
TATTATCATCTATTTTATTTGTTTGAAGTATAACAACAGCTTGGGACAGCTGGAAGTTATGTTAGTGGAATAATTTGATAGTGGAAATTTCTTCAAGTAT
CAATATGCGCCTCATGAGGACTACCTCATCTCTTGAATTAAGTGATTATGAGCAGGCTGATATGTATAATATCATAGAAAAAGTTACTCAAGACAGCAC
TTACAAGCCTTTTCAGCTATTTAATGCTATCATTGTTGTGCTTTCATCGTTATCTCATTTGTTATCTAGTCTATTTTTATTGGAACATGGAACATTGGGGT
AGCAATTTTACTCTTATTTGTTCCAGTATTATCTTTGGTACTTTTTCTCAGAGTGGGACAAATTAGAGTTTAAATCCAGTGGCAGAGCAAGTCTGTA
AAGAGAAACATGGTATATGTATATTTTATGACTCATGATTTTCTTAAAGAAATCAAGTAAATAATATAGCAATATTAATCAATCAATTAATTTGGAA
AATTAAGAAAAAGGATTTATCAACCAAGATTTAGCTATTGCTCGTAAGAAGACATATTTCAATATTTTTCTTGATTTTCATTTTGAATTTGATAAATATTCTT
ACGATATTTGCTATGATCTTTTCGGTAAGAGCAGGAAAACTTCTTATAGGTAATTTGGTAAGTCTCATACAAGCTATTTCTAAAATCAATACTTATTCTC
AAACAATGATTCAAAATATTTACATCTATTATAATACTAGTTTGTATTGGAACAACCTTTTGAGTTTAAAGAGAGAAAAAGTGATGTCACAAAAAAA
TAGAAGATACTGAAATATGCAATCAACATATAGGAAGTTGAAAGTAATTAATTTATCATATGTTTACCCTAATTCGAATGCCTTTGCATAAAGAATAT
CAATTTATCTCTTTGAAAAAGGAGAAATTAAGTCTGATTGTGAGAAAAATGGTTCAAGGAAAAAGTACACTAGTAAAGATAATTTACGATTATATCAAC
CAACTATGGGAATAATCCAATACGACAAAATGAGAAGTAGTTTGATGCCTGAGGAGTTTATCAGAAAAACATATCGGTGCTGTTCCAAGATTTTGTG
AAGTATGAGTTAACGATAAGAGAGAATATAGGATTGAGTGATTGTCTTCTCAATGGGAAGATGAGAAAAATTAATAAGTACTAGATAATTTAGGACTC
GATTTTTTGAACAAATAATCAATATGTAAGTGTATGATACGAGTTAGGAAATTTGGTTTCAAGAAGGGCATCAACTTTACAGGAGGTGAGTGGCAAAAAAT
TGCAATTAGCAAGGACATCTTTAAGAAAGCTTCAATTTATATTTTATGATGAACCAAGTGCCTGCACTGCTGTAGCTGAGGAAAAATTTGATTA
TTTTGTGCTCTTTTCGGAATAAATATTTCAATTTTCTCATAGTTTGAATGCTGCCAGAAAAGCAATAAAATCGTGGTTATGAAAGATGGACA
GGTCGAAGATGTTGGAAGTCATGATGCTCTCTGAGAAGATGTCATACTATCAAGAACTTTATTTACAGAGCAATATGAGGATAATGATGAATAAAA
AAAAATATAAAAAAGAAATGTTGAAAAAATTTATGCTCAATGGGATGAGAGAACTAGAAAAAATAAAGAAAACTTCGATTTCCGGAGAGTTGACTCTCTC
TACAGGATTTGCCCTGGTATATTTAATGTTAGCGGAGTTAAAAAATAAAGAAATCAAGATATATCAGAAAAAGATGACAAATTTATTAATTAAT
GTTAGCAAACCTTTCAACATATGGGCTTTTAAACAGGATCACTTTATTCGGGAGCAGCTGGCATTGCAATTAAGTATCTACATTTACGAGAAGATGACGA
AAAATATAAGAATCTTCTTGATAGCCTAAATAGATATATCGAATATTTTCGTCAGAGAAAAAATGGAAGGATTTAATTTGAAAAACATTACTCCTCTGAT
TATGACGTGATTGAAGGTTTATCTGGGATCTTTCTCTATTTATTAATCAACGACGAGCAATATGATGATTGAAAAATCTCATTATCAATTTTTATC
AAATCTGACTAAAGAAAAACAATGGACTAATATCGCTTTACATACGGAATCGGAGATCAGATGTTCAATCAGAAAGTGAGATTGATCCACTAGGCTGTT
TGAATATGGGATTAGCACATGGACTTGCTGGAGTGGGCTGTATCTTAGCTTATGCCACATAAAAGGATATAGTAATGAAGCCTCGTTGTCAGCTTTGC
AAAAAATTTATTTTATTTATGAAAAGTTTGAACCTTGAAGGAAAAAACAGTTTCTATGGAAGATGGAAGTTGTAGCAGATGAATAAAAAAGAGAA
AGTAATTAGGGAAGCAAGTTTCATTAGAGATGATGGTGCTATGGAGGTCAGGATATAGTCTGCTATACTTATACGGAGGATTAGCACTGGATAATGA
CTATTTTGTAGATAAAGCAGAAAAATATTAGAGTCAGCTATGCAAAAGGAACTTGGTATTGATTATATATGATTGCCATGGCTATTCTGGTTAATA
GAAATTTGTTCTTTTATTTAAGCGGCTATTAATAACAAAAAAGTTTGAATCATACATGGAAGAATTTAATGTTAGTGAAGAACTTTATCAAACTTTGAAGAATAC
GGAGATGAAAGTGGCACGGGTTTTCTTGAAGGAATAAGTGGCTGTATACTGGTATTATCGAAATTTGAATATTCAATCAATTTTACTTATTGGAGACAA
GCACTGTTACTTTTTGACGATTTTTTGAAGGAGGGAAGAGGAAATGAGAAGATATTTAATACTTATTGTGGCCTTAATAGGGATAACAGGTTTATCA
GGGTGTTATCAAAACAGTCAATAAAGGTTGAGGTTTGACGAAGGAAGTTATACTAATTTTATTTATGATAATAAATCGTATTTTCGTAAGTATAAGGAG
ATTCTCAGGAGAACGTTAAACAAATCCAAAGTAAATTTTATAGCTGTTGATTGTTGACATGGAAGAATGAAAGTGAGAAAGTGAAGAACTTTGAAAGATAAC
ATAGTGTGACTTTGGTCTTAAATAATATTATGAGGCTTCTGACAAGTCGCTATGTATGGGTATTAACGACAGATACTATAAGATACTTCCAGAAAGTGA
TAAGGGGGCGGTCAAAGCTTTGAGATTACAAAACCTTTGATGTGACAAGCGATATTTCTGATGATAATTTTGTATTGATAAAAAATGATTACGAAAAA
TTGACTATATGGGAAATATTACAGTATATCGGACACCACCGTATCTGATGAAGAATTTGGGAGAATATCAGGATGTTTTAGCTGAAGTACGTGTGTTTG


```

ATTCAGTTAGTGGCAAAAGTATCCCGAGGTCTGAATGGGGGAGAATTGATAAGGATGGTTCAAATTCCAAACAGAGTAGGACGGAATGGGATTATGG
CGAAATCCATTCTATTAGAGGAAAAATCTCTTACTGAAGCATTGCGGTTGAGATAAATGATGATTTTAAGCTTGCAACGAAGGTAGGAACTAGAGTG
AAAAAATACAGTGGTTTCCTTTTATCGTTTGTTCGTTGGGTTTATCAGCAACTGTCATGGGGAGACAACAATTCACAACAGTTACTCTCAAATAA
TATTAATACGGAATTAATTAATCATAATTCTAATGCAATTTTATCTTCAACAGAGGGATCAACGACTGATTTCGATTAATCTAGGGGCGCAGTCACCTGCA
GTAATCGACAACAAGGACTGAATTGGATGTAAGTGGTCTGCTAAAACCTTATTACAGACATCAGCTGTTCAAAAAGAAATGAAAGTTTCGTTGC
AAGAACTCAAGTTAGTTCTGAATTCAGTAAGAGAGATAGCGTTACAAATAAAGAAGCAGTTCAGTATCTAAGGATGAGCTACTTGAGCAAAGTG
AAGTAGTCGTTTCAACATCATCGATTCAAAAAATAAAATCCTCGATAATAAGAAGAAAAAGAGCTAACTTCGTTACTTCTCTCCGCTTATTAAGGAA
AAACCATCAAATTCATAAGATGCATCTGGTGTAATTGATAATTCTGCTTCTCCTCTATCTTATCGTAAAGCTAAGGAAGTGGTATCTCTTAGACAACCTT
TAAAAAATCAAAAAGTAGAGGCACAACCTCTATTGATAAGTAATTCTTCTGAAAAGAAAGCAAGTGTTTATACAAATTCACATGATTTTGGGATTAT
CAGTGGGATATGAAATATGTGACAAAATAATGGAGAAAGCTATGCGCTCTACCAGCCCTCAAAGAAAATTTCTGTTGGAATTATTGATTCAGGAATCAT
GGAAGAACATCCTGATTTGTCAAATAGTTTAGGAAATTATTTAAAAATCTTGTTCCTAAGGGAGGGTTTGATAATGAAGAACCTGATGAACTGGAA
ATCCAAGTGATATTGTGACAAAATGGGACACGGGACGGAAGTCGAGGTCAGATTACAGCAAATGGTAATTTTAGGAGTAGCACCAGGGATTAC
TGTAATATATACAGAGTATTTGGTGAATCTTTGAAAATCGGAATGGGTAGCTAGAGCAATAAGAAGAGCTGCGGATGATGGGAACAAGGTCATC
AATATAAGTGCTGGACAGTATCTTATGATTTTCAGGATCGTATGATGATGGAACAAATGATTATCAAGAGTATCTTAATTATAAGTCAGCAATAAATTATG
CAACAGCAAAAAGGAAGTATTGTTGTCGAGCTCTTGGTAATGATAGTTTAAACATACAAGATAACCAAAACAATGATAAACTTTCTTAAGCGTTTCAGA
AGTATAAAGGTTCTTGGAAAAGTTGTAGATGCACCGAGTGATTGAGGATGTAATAGCCGTAGGTGGAATAGATGGTTATGGTAATTTTCTGATTTT
AGTAATATTGGAGCGGATGCAATTTATGCTCCTGCTGGCACAACGGCCAATTTTAAAAAATATGGGCAAGATAAATTTGTCAGTCAGGGTTATTATTTG
AAAGATTGGCTTTTTACAATACTAATACTGGCTGGTACCAATATGTTATGGCAACTCATTTGCTACTCTAAAGTATCTGGGGCACTGGCATTAGTA
GTTGATAAATATGGAATAAAGAATCCTAACCACTAAAAAGGTTTCTTCTAATGAATTTCTCCAGAAGTTAATGGGAATAGAGTATTGAATATTGTTGAT
TTATTGAATGGGAAAAATAAAGCTTTTAGCTTAGATACAGATAAAGGTCAGGATGATGCTATTAACCATAAATCGATGGAGAATCTTAAAGAGTCTAG
GGATACAATGAAACAGGAACAAGATAAAGAAATTCAAAGAAATACAATAACAATTTTCTATCAAAAATGATTTTCATAACATTTCAAAAAGAAGTA
ATTTCAAGTTGATTATAATTAATCAAAAAATGGCTAATAATCGAAATTCGAGAGGTGCTGTTTCTGTACGAAGTCAAGAAATTTTACCTGTTACTGGA
GATGGAGAAGATTTTTACCGGCTTAGGTATAGTGTGATCTCAATCCTTGGTATATTGAAAAGAAAAGACTAAAAATTGATAGGCTGAAAAAGGGTG
GAAATCCACTCTTTTTCTTTTTCAATCCTTGACTTCCGCTTGAAATCACTTGAGCATTCTACAATTACCGGCTTTAGGTATAGTGTGTATCTTGAGTAC
TAAACAATCGGAGGTAAAGTGGTGATAAAATTTTAATAGTTGATGATGATCAGGAAATTTAAAAATTAATGAAGACAGCATTAGAAATGAGAACTA
TGAAGTTGCGACGCATCAAAACATTTCACTTCCCTTGGATATTACTGATTTTCAGGGATTTGATTTGATTTT

```

Annexe 1 : séquence génomique au format fasta du fragment d'intérêt du génome de *Lactococcus lactis*

```

tar -xvf scan_for_matches.tgz
# dézippe le dossier
gcc -O -o scan_for_matches ggpunit.c scan_for_matches.c
# compile le programme
./run_tests tmp
# lance le programme avec un fichier test (tmp)
diff tmp test_output
# vérifie que rien n'est retourné et donc que le programme fonctionne correctement

scan_for_matches matrice_file.txt < seq1_L.lactis.txt > output_file.txt
# exécute le programme avec la matrice poids position du RBS et la séquence d'intérêt
wc -l output_file.txt
# compte le nombre de ligne et donc retourne le nombre de séquence RBS identifiées (multiplié par 2)

scan_for_matches sigma_matrice_file.txt < seq1_L.lactis.txt > output_sigma_file.txt
# exécute le programme avec la matrice poids position du site du facteur sigma et la séquence d'intérêt

scan_for_matches rho_pattern_file.txt < seq1_L.lactis.txt > output_rho_file.txt
# exécute le programme avec le pattern du site de terminaison rho-dépendant et la séquence d'intérêt

```

Annexe 2 : commande terminal (bash) utilisée pour analyser le fragment génomique d'intérêt avec le programme ScanForMatches

```

{(-23,-53,20,-33),(-22,-34,20,-46),(14,-46,-7,-15),(-27,-52,19,-17),(-4,-17,14,-10)}>44 5...12 DTG
# matrice_file.txt
# matrice poids position utilisée pour rechercher les RBS dans notre séquence d'intérêt
# les poids sont attribués respectivement à A, C, G et T

```

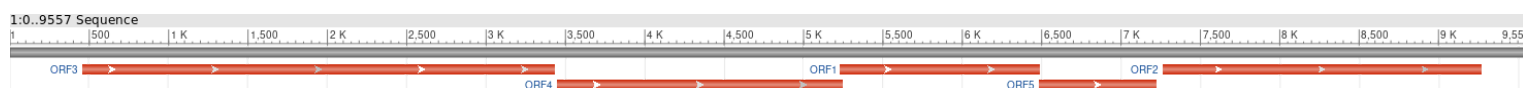
```

{(-22,-29,-25,16),(-18,-25,-21,16),(-28,-17,18,-11),(12,-9,-19,-6), (-2,12,-15,-7),(9,-13,-10,-1)}> 30
16...35 {(-32,-29,-45,17),(17,-45,-45,-28),(-3,-9,-21,11),(14,-7,-19,-21),
(14,-7,-19,-20),(-28,-35,-45,17)}> 30
# sigma_matrice_file.txt
# matrice poids position utilisée pour rechercher les sites de fixation du facteur  $\sigma$  dans notre séquence d'intérêt
# les poids sont attribués respectivement à A, C, G et T

r1={AU,UA,GC,CG,GU,UG,GA,AG}
p1 = 3...5 p2 = 3...10 3...9 r1~p2 ~p1 TTTTTT[1,0,0]
# rho_pattern_file.txt
# pattern utilisé pour rechercher les terminateurs rho-dépendants dans notre séquence d'intérêt

```

Annexe 3 : matrices et pattern utilisés avec l'outil ScanForMatches



Annexe 4 : résultats obtenus avec ORFfinder et le fragment génomique d'intérêt

GENEMARK PREDICTIONS						
Sequence: seq_L_lactis_lactis						
Sequence file: seq.fna						
Sequence length: 9557						
GC Content: 31.57%						
Window length: 120						
Window step: 12						
Threshold value: 0.500						

Matrix: Lactococcus_lactis_II1403						
Matrix author: -						
Matrix order: 4						
List of Open reading frames predicted as CDSs, shown with alternate starts						
(regions from start to stop codon w/ coding function >0.50)						
Left end	Right end	DNA Strand	Coding Frame	Avg Prob	Start Prob	

455	3436	direct	fr 2	0.55	0.02	
623	3436	direct	fr 2	0.56	0.03	
785	3436	direct	fr 2	0.57	0.03	
848	3436	direct	fr 2	0.57	0.05	
1310	3436	direct	fr 2	0.53	0.02	
3447	5249	direct	fr 3	0.63	0.75	
3456	5249	direct	fr 3	0.63	0.10	
3561	5249	direct	fr 3	0.65	0.37	
3630	5249	direct	fr 3	0.63	0.14	
3687	5249	direct	fr 3	0.63	0.31	
5230	6486	direct	fr 1	0.67	0.09	
5239	6486	direct	fr 1	0.68	0.90	
5242	6486	direct	fr 1	0.68	0.73	

5374	6486	direct	fr 1	0.71	0.00
5653	6486	direct	fr 1	0.69	0.21
5806	6486	direct	fr 1	0.69	0.92
5827	6486	direct	fr 1	0.69	0.84

6483	7220	direct	fr 3	0.53	0.02
6507	7220	direct	fr 3	0.54	0.63
6564	7220	direct	fr 3	0.58	0.01
6696	7220	direct	fr 3	0.71	0.35
6738	7220	direct	fr 3	0.72	0.25
6789	7220	direct	fr 3	0.75	0.50
6873	7220	direct	fr 3	0.78	0.18
6939	7220	direct	fr 3	0.73	0.03

7222	9270	direct	fr 1	0.89	0.15
7279	9270	direct	fr 1	0.91	0.75
7504	9270	direct	fr 1	0.91	0.00
7792	9270	direct	fr 1	0.90	0.00

9406	9558	direct	fr 1	0.51	0.45
9409	9558	direct	fr 1	0.57	0.64

List of Regions of interest
(regions from stop to stop codon w/ a signal in between)

LEnd	REnd	Strand	Frame
425	3436	direct	fr 2
3444	5249	direct	fr 3
5185	6486	direct	fr 1
6453	7220	direct	fr 3
7210	9270	direct	fr 1

ABOUT THE MATRIX USED:

Training set derived by GeneMarkS, 4.27 September 2014
Tue Sep 23 15:01:07 2014

Annexe 5.1: résultats obtenus avec GeneMark et le fragment génomique d'intérêt (seuil à 0,5). Les éléments en surbrillance indiquent les différences observées avec les résultats obtenus pour un seuil à 0,4 (Annexe 5.2). Le graphe associé est renseigné dans les documents complémentaires (document1).

GENEMARK PREDICTIONS

Sequence: seq_L_lactis_lactis
Sequence file: seq.fna
Sequence length: 9557
GC Content: 31.57%
Window length: 120
Window step: 12
Threshold value: 0.400

Matrix: Lactococcus_lactis_II1403
Matrix author: -
Matrix order: 4

List of Open reading frames predicted as CDSs, shown with alternate starts
(regions from start to stop codon w/ coding function >0.40)

Left end	Right end	DNA Strand	Coding Frame	Avg Prob	Start Prob

	455	3436	direct	fr 2	0.55 0.02
	623	3436	direct	fr 2	0.56 0.03
	785	3436	direct	fr 2	0.57 0.03
	848	3436	direct	fr 2	0.57 0.05
	1310	3436	direct	fr 2	0.53 0.02
	3447	5249	direct	fr 3	0.63 0.75
	3456	5249	direct	fr 3	0.63 0.10
	3561	5249	direct	fr 3	0.65 0.37
	3630	5249	direct	fr 3	0.63 0.14
	3687	5249	direct	fr 3	0.63 0.31
	5230	6486	direct	fr 1	0.67 0.09
	5239	6486	direct	fr 1	0.68 0.90
	5242	6486	direct	fr 1	0.68 0.73
	5374	6486	direct	fr 1	0.71 0.00
	5653	6486	direct	fr 1	0.69 0.21
	5806	6486	direct	fr 1	0.69 0.92
	5827	6486	direct	fr 1	0.69 0.84
	6483	7220	direct	fr 3	0.53 0.02
	6507	7220	direct	fr 3	0.54 0.63
	6564	7220	direct	fr 3	0.58 0.01
	6696	7220	direct	fr 3	0.71 0.35
	6738	7220	direct	fr 3	0.72 0.25
	6789	7220	direct	fr 3	0.75 0.50
	6873	7220	direct	fr 3	0.78 0.18
	7222	9270	direct	fr 1	0.89 0.15
	7279	9270	direct	fr 1	0.91 0.75
	7504	9270	direct	fr 1	0.91 0.00
	7792	9270	direct	fr 1	0.90 0.00
	9406	9558	direct	fr 1	0.51 0.45
	9409	9558	direct	fr 1	0.57 0.64
List of Regions of interest					
(regions from stop to stop codon w/ a signal in between)					
LEnd	REnd	Strand	Frame		

	425	3436	direct	fr 2	
	3444	5249	direct	fr 3	
	5185	6486	direct	fr 1	
	6453	7220	direct	fr 3	
	7210	9270	direct	fr 1	
	9388	9558	direct	fr 1	

ABOUT THE MATRIX USED:					
Training set derived by GeneMarkS, 4.27 September 2014					
Tue Sep 23 15:01:07 2014					

Annexe 5.2: résultats obtenus avec GeneMark et le fragment génomique d'intérêt (seuil à 0,4). Les éléments en surbrillance indiquent les différences observées avec les résultats obtenus pour un seuil à 0,5 (Annexe 5.1). Le graphe associé est renseigné dans les documents complémentaires (document2).

GeneMark.hmm PROKARYOTIC (Version 3.42)
Date: Fri Nov 29 04:52:54 2024
Sequence file name: seq.fna
Model file name:
/home/genemark/parameters/prokaryotic/Lactococcus_lactis_II1403/GeneMark_hmm_combined.mod
RBS: true
Model information: Lactococcus_lactis_II1403

FASTA definition line: seq_L_lactis_lactis

Predicted genes

Gene	Strand	LeftEnd	RightEnd	Gene	Class
#			Length		
1	+	174	347	174	2
2	+	455	3436	2982	2
3	+	3447	5249	1803	2
4	+	5242	6486	1245	2
5	+	6483	7220	738	2
6	+	7222	9270	2049	1
7	+	9475	>9555	81	1

Annexe 6 : résultats obtenus avec GeneMarkHMM et le fragment génomique d'intérêt

```
>seq_L_lactis_lactis:[161,176]
GGAGG CACTCAA ATG
>seq_L_lactis_lactis:[282,298]
GGAGC TCTGATGGG TTG
>seq_L_lactis_lactis:[1802,1814]
GGGGA TATCG TTG
>seq_L_lactis_lactis:[2317,2335]
GGATG TAAATAGTAAG ATG
>seq_L_lactis_lactis:[2531,2549]
GGAGA AAAGGGGAGAG TTG
>seq_L_lactis_lactis:[4581,4599]
GGAGA ATTAAGTCTA TTG
>seq_L_lactis_lactis:[5497,5512]
GGAGC AGCTGGCA TTG
>seq_L_lactis_lactis:[5611,5625]
GAAGG ATTTAAT TTG
>seq_L_lactis_lactis:[5796,5808]
GGAGA ATCAG ATG
>seq_L_lactis_lactis:[6340,6352]
GGAGA TGAAA GTG
>seq_L_lactis_lactis:[6367,6379]
GAAGG AATAA GTG
>seq_L_lactis_lactis:[6469,6485]
GGAGG GAAGAGGAA ATG
>seq_L_lactis_lactis:[6836,6854]
GGGGG CGGTCAAAGCT TTG
>seq_L_lactis_lactis:[6990,7005]
GGAGA ATATCAGG ATG
>seq_L_lactis_lactis:[7067,7086]
GGGGA GAATTGATAAGG ATG
>seq_L_lactis_lactis:[7205,7224]
GAAGG TAGGAACTAGA GTG
>seq_L_lactis_lactis:[7449,7462]
GGATG TAACTG GTG
>seq_L_lactis_lactis:[7593,7606]
GGATG AGCTAC TTG
>seq_L_lactis_lactis:[7942,7954]
GGAGA AAGCT ATG
>seq_L_lactis_lactis:[8071,8086]
```

```
GGAGG GTTTGATA ATG
>seq_L_lactis_lactis:[8547,8566]
GGATG TAATAGCCGTAG GTG
>seq_L_lactis_lactis:[8608,8626]
GGAGC GGATGCAATTT ATG
>seq_L_lactis_lactis:[9397,9411]
GGAGG TAAAGTG GTG
```

Annexe 7.1 : résultats obtenus avec ScanForMatches et le fragment génomique d'intérêt pour la recherche des RBS

```
>seq_L_lactis_lactis:[62,103]
TTGGTA ATATAACGCGAGCATAATAAACGGCTCTGA TTAAAT
>seq_L_lactis_lactis:[111,151]
TTGCCA TCAATGCTCAGGACGAATATTGCATGATC TATCAT
>seq_L_lactis_lactis:[272,304]
TTGTAA AACAGGAGCTCTGATGGGTTG TAACAT
>seq_L_lactis_lactis:[368,395]
TTGTTA GTTCAGACATGGATAC TATCCT
>seq_L_lactis_lactis:[399,435]
TTTATA AGTTATTTAGGGTTGCTAAATAGCT TATAAA
>seq_L_lactis_lactis:[470,504]
TTTAAA GCTCAACCGTTTTTAGTAAGAAA TACAAT
>seq_L_lactis_lactis:[530,570]
TTTACT GAATATACTCAAGTCATTGAGACTGTAAG TAAAAA
>seq_L_lactis_lactis:[581,612]
TTGGAA CAGTTACTACTAGCTAATCC TAAACT
>seq_L_lactis_lactis:[675,716]
TTGAAT CTATTTACAAGTATTATAAGAGAAGTTATT TACGAT
>seq_L_lactis_lactis:[769,812]
TTCACA GTACAAGTTAATGGGAAAGACTACAAAGGGTA TAAGAT
>seq_L_lactis_lactis:[927,971]
TTTATA CCATAAATAGTAGTGAGCTTGAAGAAGTAAATA TTAAAT
>seq_L_lactis_lactis:[1054,1096]
TTGCTA TGGAGACGAATATAGAGAAGTATCGGAACAA TATCTT
>seq_L_lactis_lactis:[1148,1185]
TTGTCA GATTTTCTTGGAACACTTTTTTGAC TAAAGT
>seq_L_lactis_lactis:[1323,1356]
TTGAGA ATGATAATTATATTCAAATTGA TTTAAT
>seq_L_lactis_lactis:[1527,1561]
TTGATT CTACATTTGGCATAGGAGCTCCA TATAAT
>seq_L_lactis_lactis:[1578,1613]
ATGACT TTTATGAGTCCGAACCGAGTACTC TATACT
>seq_L_lactis_lactis:[1716,1754]
TGGACT TAGAAAAGAAAAGTGAAGTTCAAGGGT TAGAAT
>seq_L_lactis_lactis:[1766,1812]
TTGGCA AAGGAGTATGAAAAAGATATTTTATTTTAGGGGA TATCGT
>seq_L_lactis_lactis:[2021,2061]
TTTACA AGTACAAGTCACAATGAAGTTCTGTTAAC TAATAT
>seq_L_lactis_lactis:[2123,2152]
TTGAAA TTCTACATTACAAGCATG TACAAT
>seq_L_lactis_lactis:[2211,2257]
TTGGTA ATTTACCTTGGGAAC TTATTACAGAGACTTTGAT TATATT
>seq_L_lactis_lactis:[2300,2334]
TGGAAA ATTTGGGGAAGGGATGTAAATAG TAAGAT
>seq_L_lactis_lactis:[2352,2386]
TTCAAA GCAAAGAAATCCCAAAGAGTTT TATATT
>seq_L_lactis_lactis:[2483,2523]
TTTATA GAGCTACAAGAATATTTTGAAGATGAAAA TATCAT
>seq_L_lactis_lactis:[2648,2678]
TTGCCC TTTAACGAGTGGCTTTATC TAAAGT
>seq_L_lactis_lactis:[2758,2785]
```


TGGAAA TCTATTCTTCCTAAGA TATACT
 >seq_L_lactis_lactis:[3002,3031]
 TTGCTT ACATTGATAAAAGATACT AATAAT
 >seq_L_lactis_lactis:[3162,3205]
 ATGAAA AGATTGAACATTATCTTAAGCTTCTGAAAGTT AATAAT
 >seq_L_lactis_lactis:[3225,3263]
 ATGACA AGAATTTTAAAGAATTAAAGCATGCCA TAAAAA
 >seq_L_lactis_lactis:[3294,3333]
 TTGAAC TTCAGAAAGTTTATTCAATTATTGACAG TATCAT
 >seq_L_lactis_lactis:[3513,3559]
 TTGAAG TTAATTTTTCAGTTGGAAAAACGTTATGCAATTTA TTTAAT
 >seq_L_lactis_lactis:[3591,3622]
 TTGGCT AGTCTTTTATTTATCAAGA TTTAAT
 >seq_L_lactis_lactis:[3712,3758]
 TGGAAA GTTATGTTAGTGGAAAATTTGATATGCGACTTTCT TACAGT
 >seq_L_lactis_lactis:[3793,3830]
 TTGAAT TAAGTGATTATGAGCAGGCTGATATG TATAAT
 >seq_L_lactis_lactis:[3968,3996]
 TTTACT CCTTATTGTTCCAGTAT TATCTT
 >seq_L_lactis_lactis:[4089,4127]
 TTGACT CATGATTTTTCATTTAAAGAAATCAAG TTAAT
 >seq_L_lactis_lactis:[4231,4261]
 TTGATT TCATTTTGAATTTGATAAA TATTCT
 >seq_L_lactis_lactis:[4317,4348]
 TTGGTA AGTCTCATACAAGCTATTTT TAAAAT
 >seq_L_lactis_lactis:[4597,4630]
 TTGTAG GAAAAAATGGTTCAGGGAAAAAG TACACT
 >seq_L_lactis_lactis:[4791,4837]
 TTGAGT GATTGTCTTCTCAATGGGAAGATGAGAAAATTAT TAAAGT
 >seq_L_lactis_lactis:[4954,4997]
 TTGCAT TAGCAAGGACATTCTTTAAGAAAGCTTCAATT TATATT
 >seq_L_lactis_lactis:[5047,5078]
 TTGATT ATTTTGTGCTCTTTCGGAA AATAAT
 >seq_L_lactis_lactis:[5086,5131]
 TTTTCA TTTCTCATAGTTTGAATGCTGCCAGAAAAGCAAA TAAAAT
 >seq_L_lactis_lactis:[5164,5204]
 TTGGAA GTCATGATGTCCTTCTGAGAAGATGTCAA TACTAT
 >seq_L_lactis_lactis:[5338,5369]
 TTGACT CTCTCTACAGGATTGCCTGG TATAAT
 >seq_L_lactis_lactis:[5657,5692]
 TTGAAG GTTTATCTGGGATACTTTCCTATC TATTAT
 >seq_L_lactis_lactis:[5845,5891]
 TTGAAT ATGGGATTAGCACATGGACTTGCTGGAGTGGGCTG TATCTT
 >seq_L_lactis_lactis:[6001,6034]
 TGGAAA GATGGACTTGTAGCAGATGAAT TAAAAA
 >seq_L_lactis_lactis:[6194,6231]
 TTGGTA TTGATTCATATATGATTTGCCATGGC TATTCT
 >seq_L_lactis_lactis:[6239,6273]
 TAGAAA TTTGTTCTTTATTTAAGCGGCTA TTAAT
 >seq_L_lactis_lactis:[6344,6389]
 ATGAAA GTGGCACGGGTTTTCTTGAAGGAATAAGTGGCTG TATACT
 >seq_L_lactis_lactis:[6463,6499]
 TTGAAA GGAGGGAAGAGGAAATGAGAAGATA TTTAAT
 >seq_L_lactis_lactis:[6571,6605]
 TTGACG AAGGAAGTTATACTAATTTTATT TATGAT
 >seq_L_lactis_lactis:[6772,6812]
 CTGACA AGTCGCTATGTATGGGTATTAACGACAGA TACTAT
 >seq_L_lactis_lactis:[6873,6916]
 GTGACA AGCGATATTTCTGATGATAATTTTGTATTGA TAAAAA
 >seq_L_lactis_lactis:[7748,7779]
 TTGATA ATTCTGCTTCTCCTCTATCT TATCGT
 >seq_L_lactis_lactis:[8026,8057]
 TTGTCA AATAGTTTAGGAAATATTT TAAAAA
 >seq_L_lactis_lactis:[8313,8346]
 TGGACA GTATCTTATGATTCAGGATCG TATGAT
 >seq_L_lactis_lactis:[8448,8482]

```

TTTAAA CATACAAGATAACCAAACAATGA TAAACT
>seq_L_lactis_lactis:[8543,8588]
TTGAGG ATGTAATAGCCGTAGGTGGAATAGATGGTTATGG TAATAT
>seq_L_lactis_lactis:[8650,8694]
TTTAAA AAATATGGGCAAGATAAATTTGTCAGTCAGGGT TATTAT
>seq_L_lactis_lactis:[8798,8836]
TTGATA AATATGGAATAAAGAATCCTAACCAAC TAAAAA
>seq_L_lactis_lactis:[9066,9108]
TTTTCA TAACATTTCAAAAGAAGTAATTTTCAGTTGAT TATAAT
>seq_L_lactis_lactis:[9204,9248]
TTTACC GGCTTTAGGTATAGTGTGTATCTCAATCCTTGG TATATT
>seq_L_lactis_lactis:[9307,9353]
TTTTCA ATCCTTGACTTCCGCTTGAAATCACTTGAGCATTG TACAAT
>seq_L_lactis_lactis:[9380,9407]
TTGAGT ACTAAACAATCGGAGG TAAAGT

```

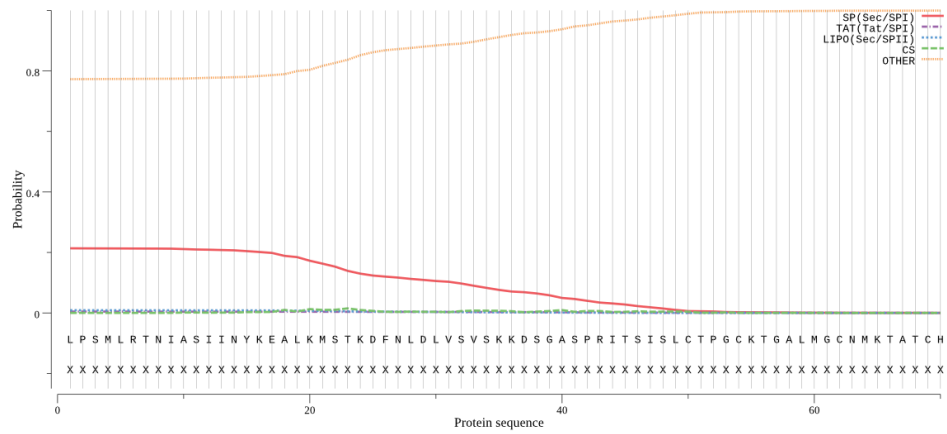
Annexe 7.2 : résultats obtenus avec ScanForMatches et le fragment génomique d'intérêt pour la recherche des sites de fixation du facteur sigma (seuil supérieur à 30)

```

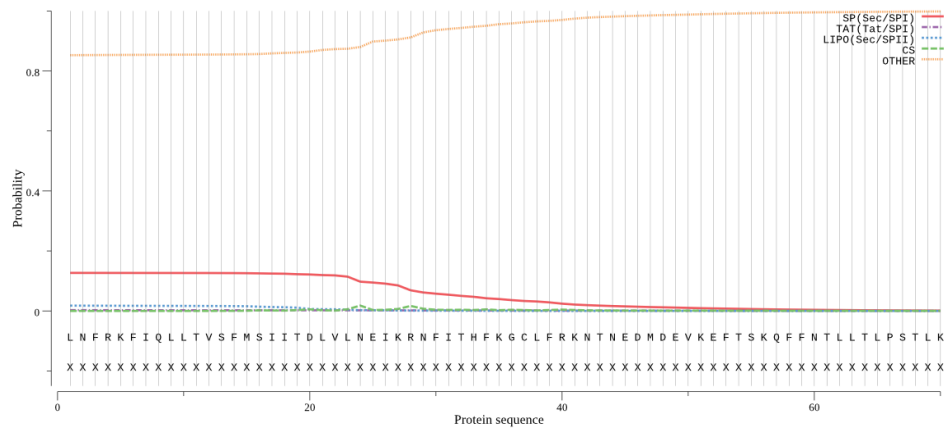
>seq_L_lactis_lactis:[29,59]
AAC CCCGG ATGAGAATT CCGGG GTT TTTTTC
>seq_L_lactis_lactis:[557,582]
ACT GTA AGTAAAAA TAA AGT TTTTTT
>seq_L_lactis_lactis:[1771,1800]
AAA GGAGT ATGAAAAA GATAT TTT TATTTT
>seq_L_lactis_lactis:[3986,4009]
AGT ATT ATCTTT GGT ACT TTTTCT
>seq_L_lactis_lactis:[4404,4424]
TTG TTT ATG GAA CAA CTTTTT
>seq_L_lactis_lactis:[6683,6717]
GTT GATTGTT GACATGAAA AGTGAGA AAC TTTTAT
>seq_L_lactis_lactis:[7224,7248]
GAA AAA AATACTA GGT TTC CTTTTT
>seq_L_lactis_lactis:[8162,8183]
TTA CAG CAAA TGG TAA TATTTT
>seq_L_lactis_lactis:[8670,8696]
TAA ATT TGTCAGTCA GGG TTA TTATTT
>seq_L_lactis_lactis:[9277,9310]
GAA AAAGGGT GGGAATCC ACTCTTT TTC TTTTTT

```

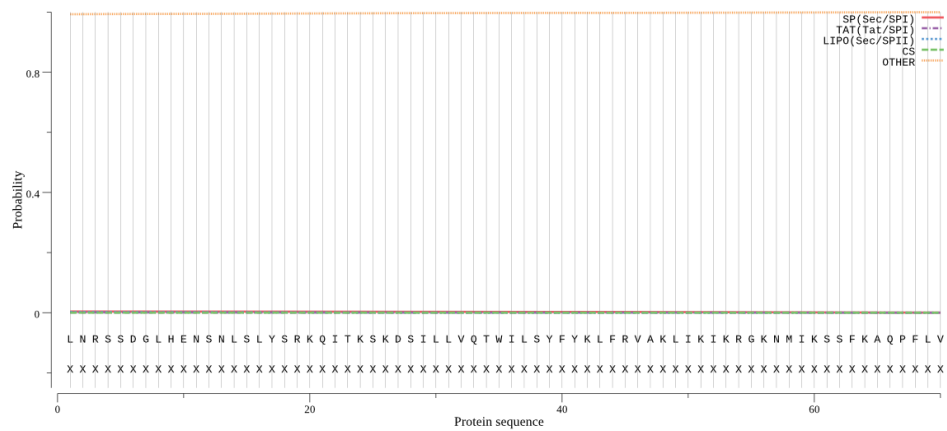
Annexe 7.3 : résultats obtenus avec ScanForMatches et le fragment génomique d'intérêt pour la recherche des sites de terminaison rho-dépendant



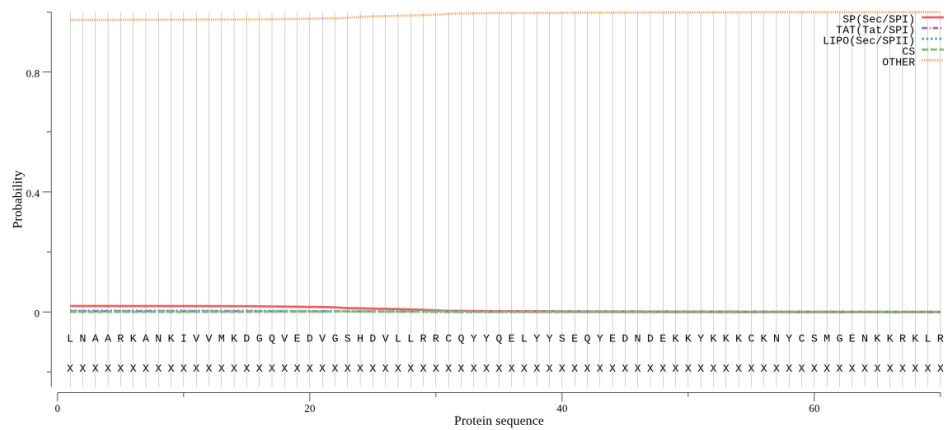
ORF0



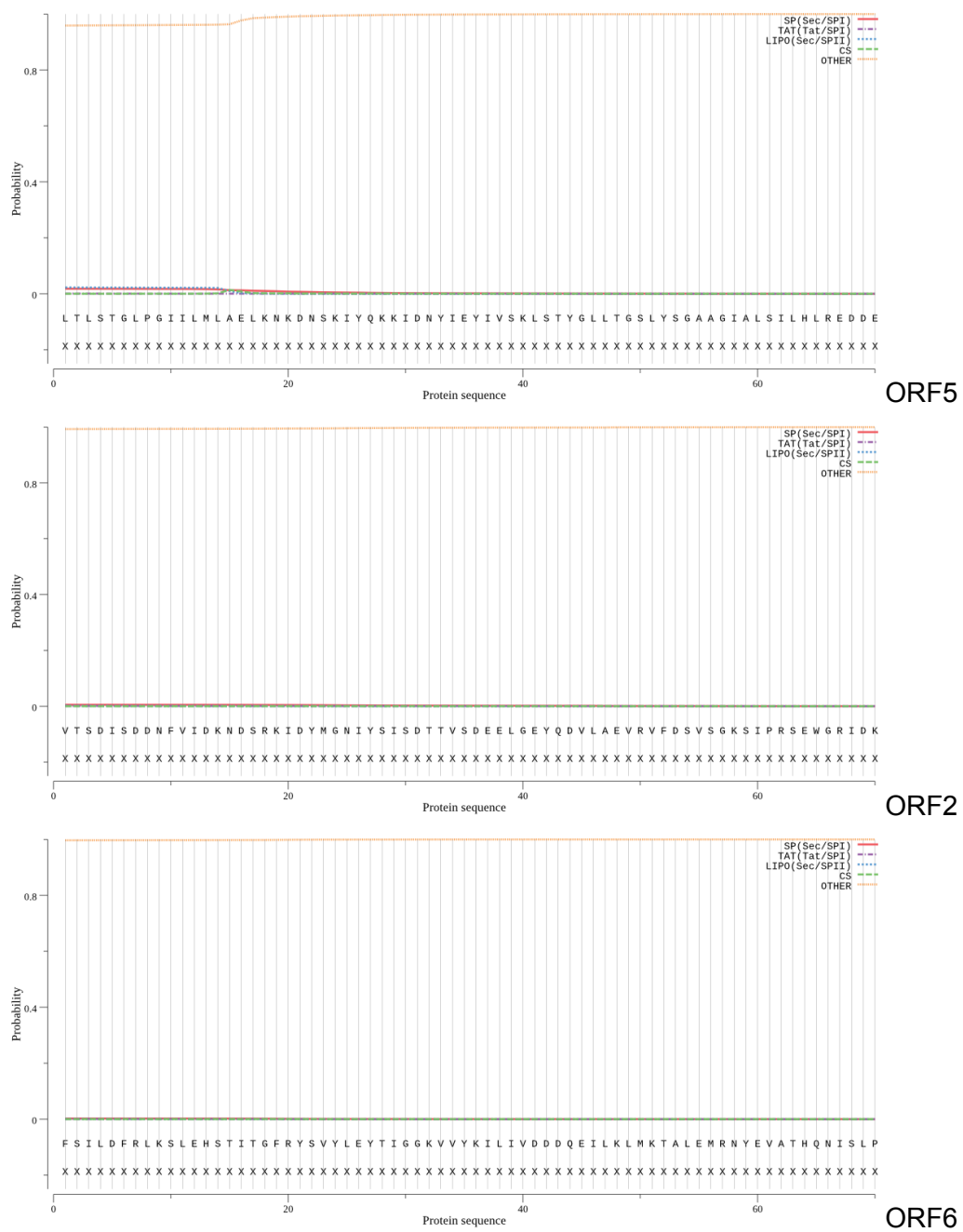
ORF3



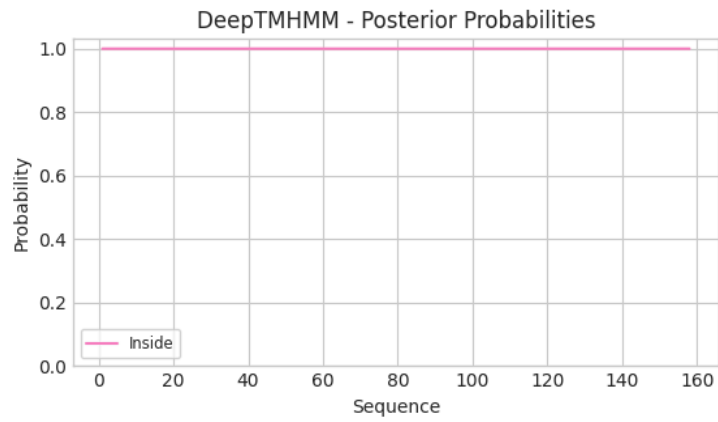
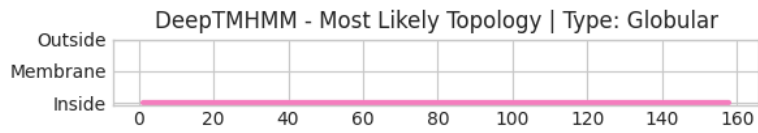
ORF4



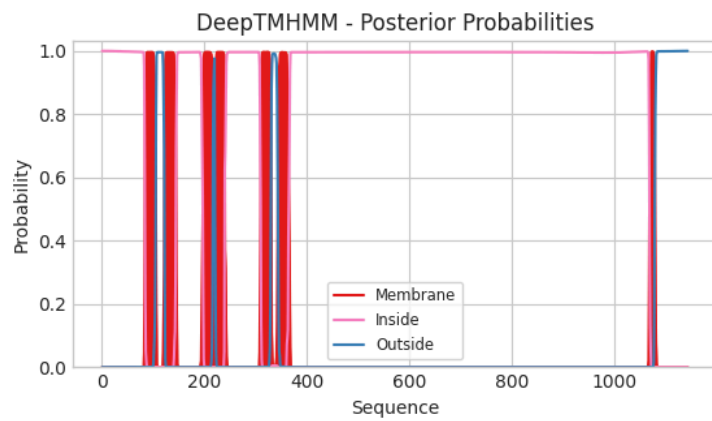
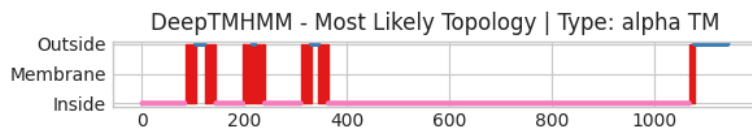
ORF1



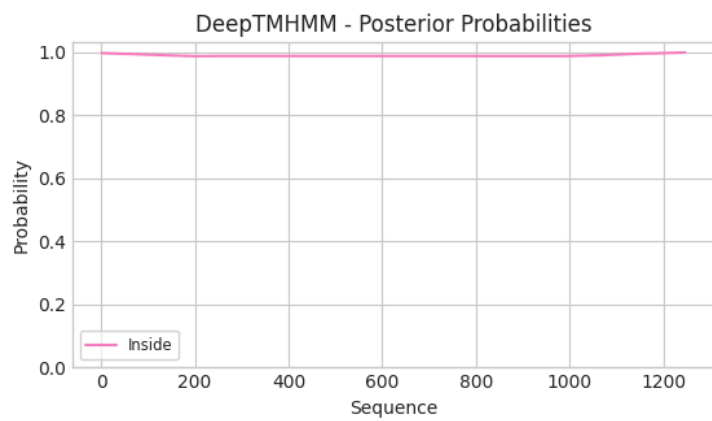
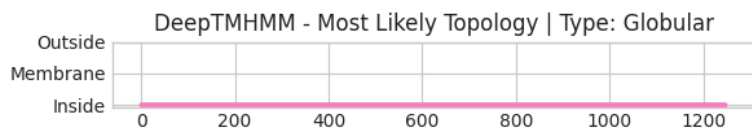
Annexe 8 : résultats obtenus avec SignalIP et les séquences protéiques identifiées



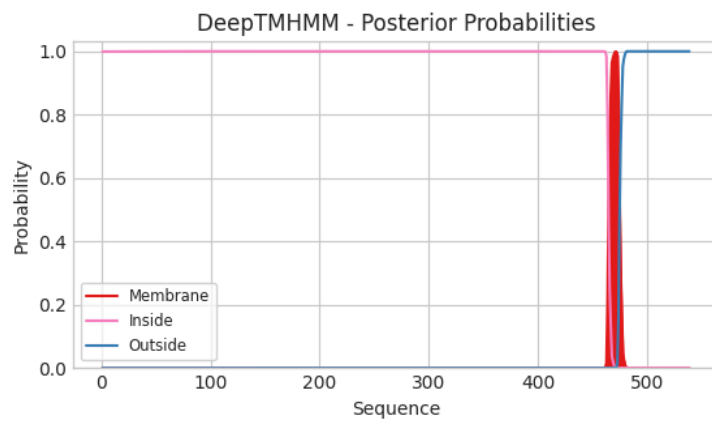
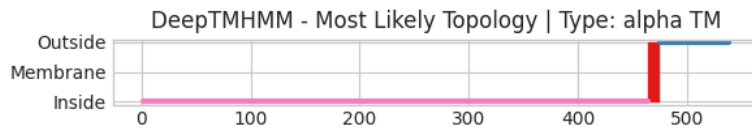
ORF0



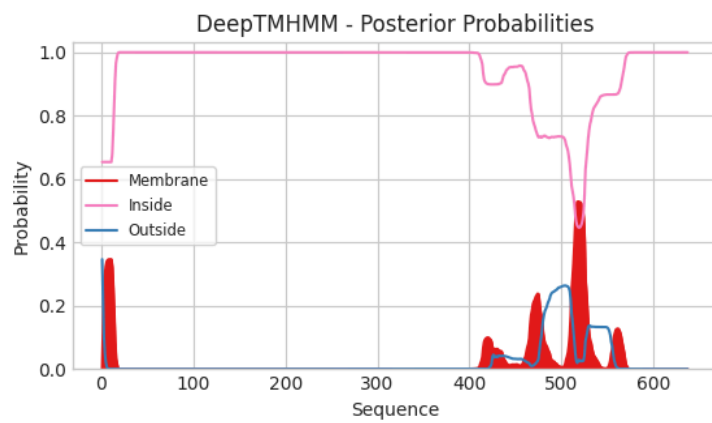
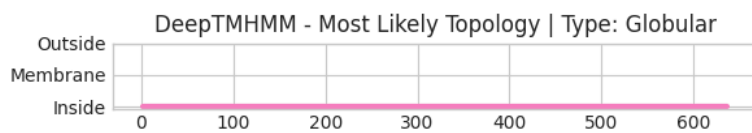
ORF3



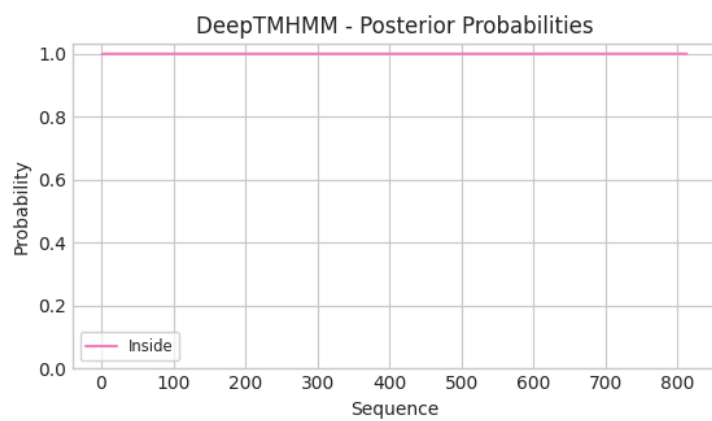
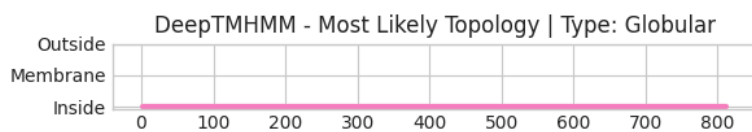
ORF4



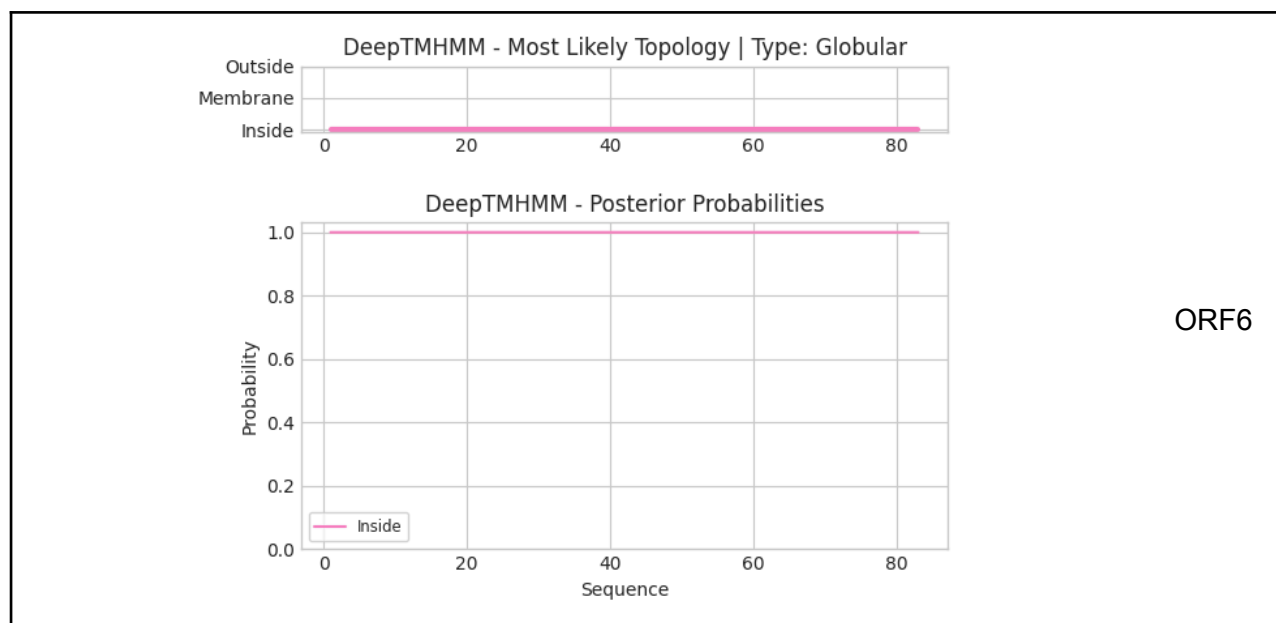
ORF1



ORF5



ORF2



Annexe 9 : résultats obtenus avec DeepTMHMM et les séquences protéiques identifiées

DOCUMENTS COMPLÉMENTAIRES

Document 1 : graphe obtenu avec GeneMark sur la séquence d'intérêt (seuil à 0,5)

Document 2 : graphe obtenu avec GeneMark sur la séquence d'intérêt (seuil à 0,4)

Document 3 : graphe obtenu avec GeneMarkHMM sur la séquence d'intérêt

Document 4 : code python pour la conversion des fichiers en format GFF

Document 5 : résultats obtenus après conversion des fichiers de sortie à l'aide du code fourni (document complémentaire 4)

Document 1 : graphe obtenu avec GeneMark sur la séquence d'intérêt (seuil à 0,5)

GeneMark

Version 2.5p (09.08.06)

Copyright 1993 - M. Borodovsky, J. McIninch

PROGRAM INFORMATION

Sequence : seq_L_lactis_lactis
Analysis Date : 11/29/24 at 4:45:35
Pages : 6
Sequence Length : 9557 bp
GC Content : 31.57%

Window Length : 120 bp
Window Step : 12 bp
Threshold Value : 0.400

PS-Version : 1.2

GeneMark Options : PostScript graph,
Mark ORFs / splice sites,
List ORFs,
List regions and/or splice sites,

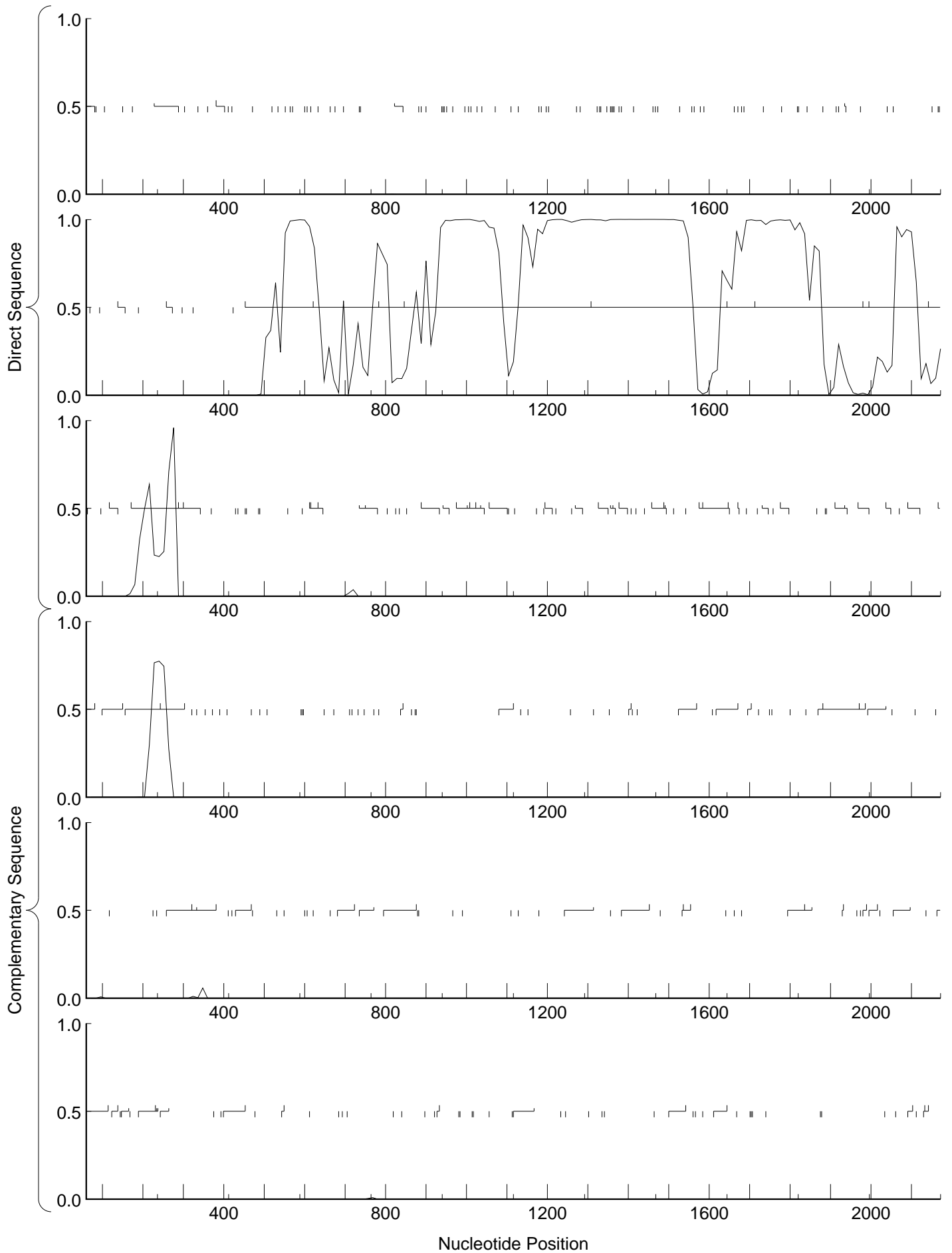
Matrix notes & comments

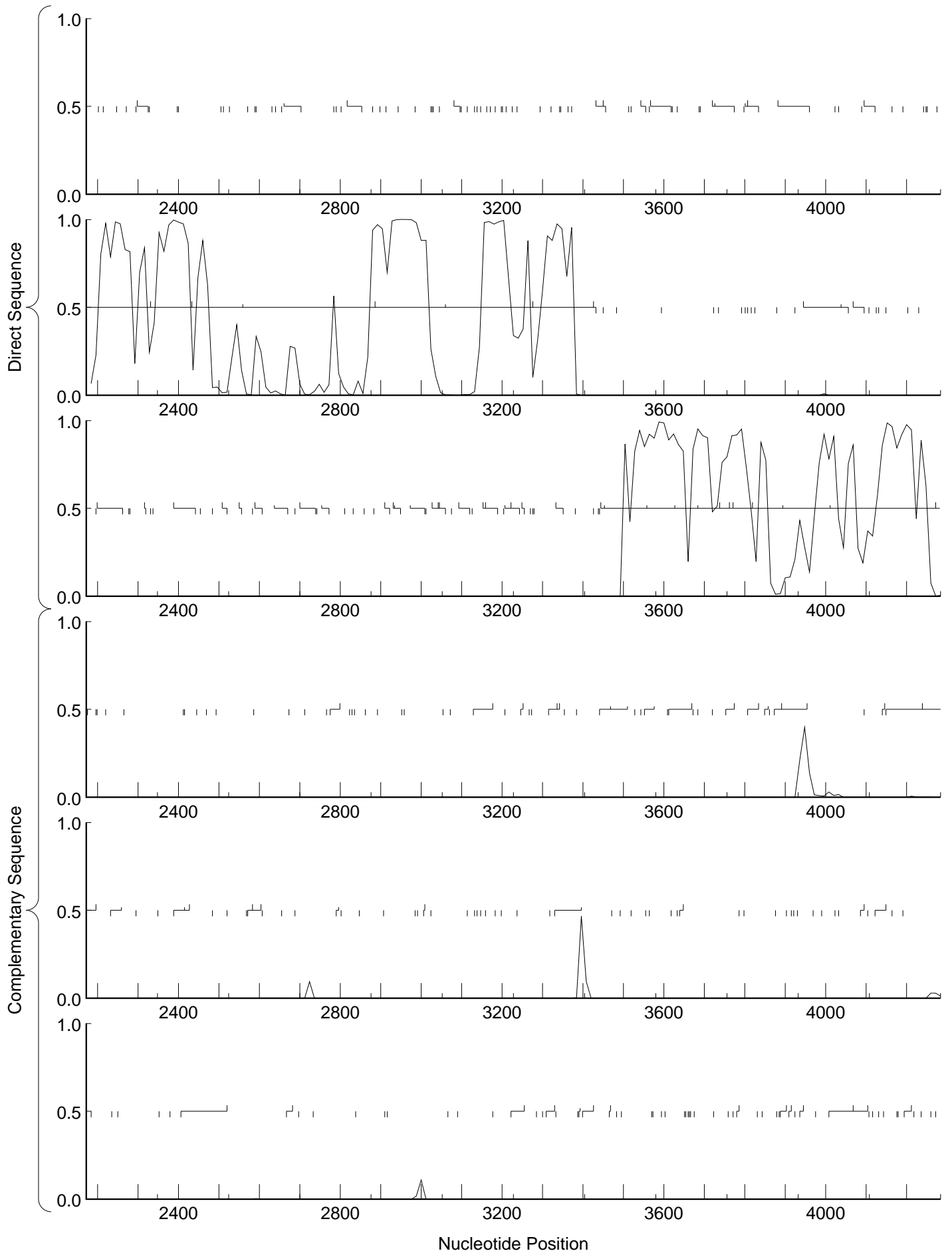
Training set derived by GeneMarkS, 4.27 September 2014
Tue Sep 23 15:01:07 2014

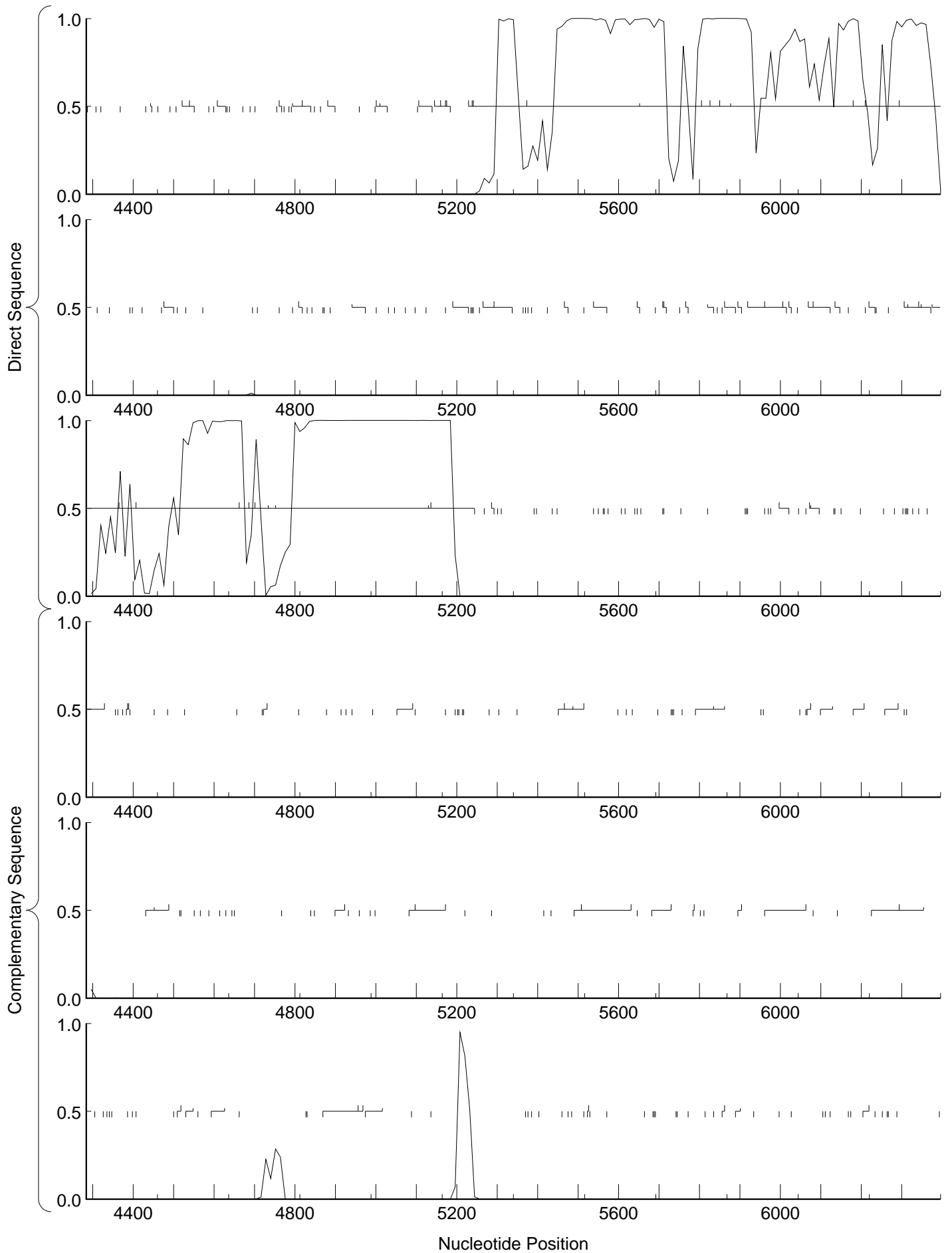
MATRIX INFORMATION

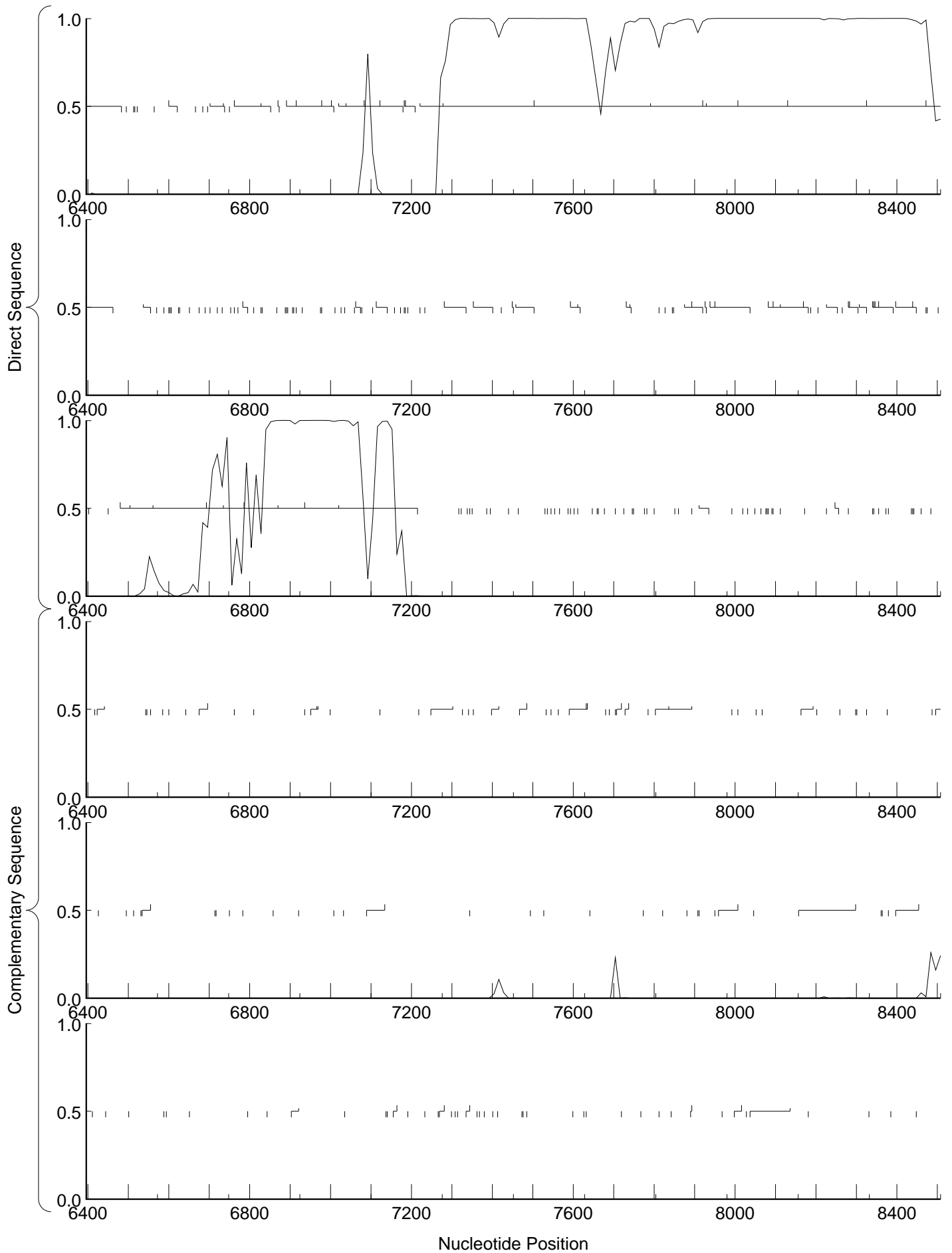
Matrix : Lactococcus_lactis_I11403
Author : -
Order : 4

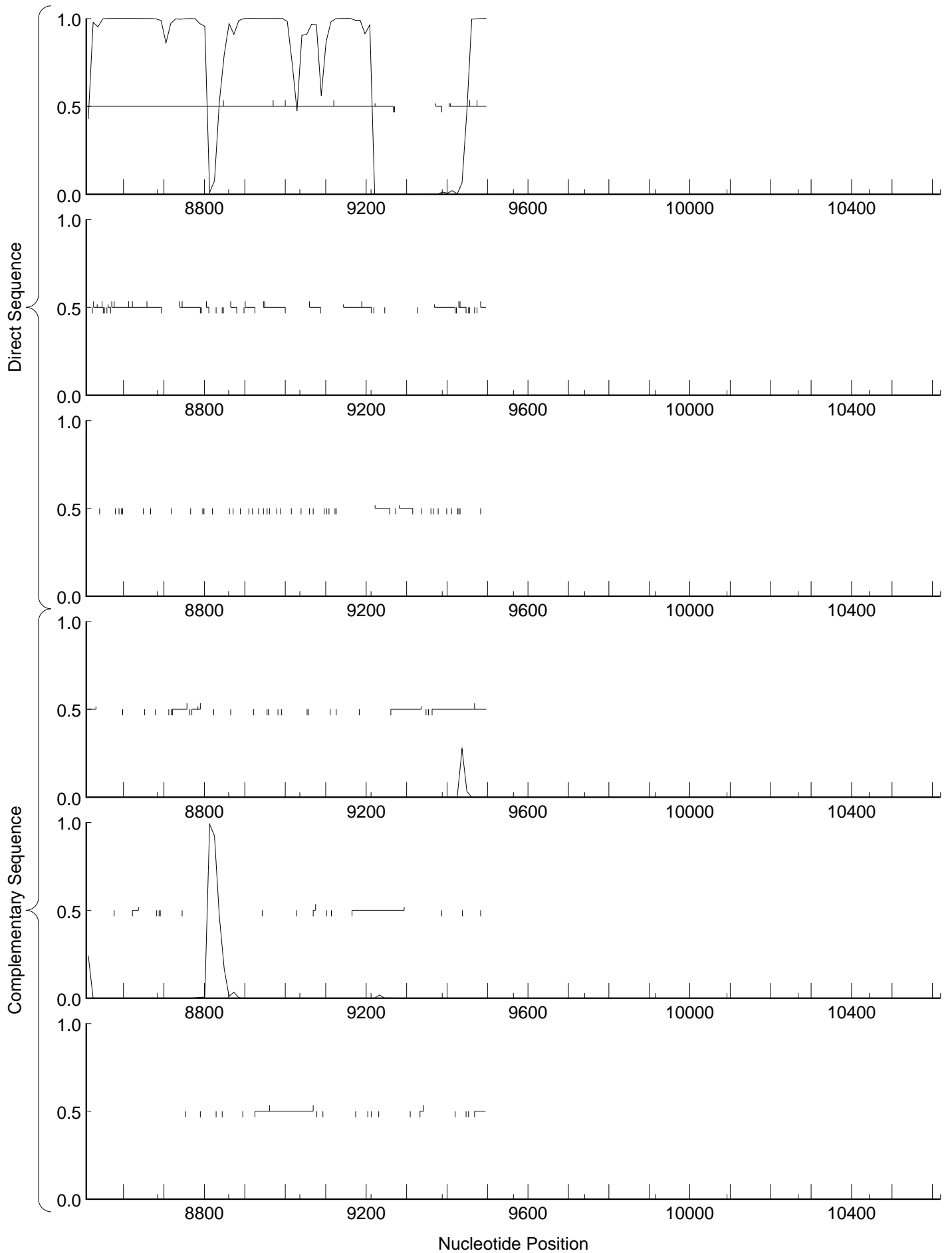
Send questions / comments to:
Dr. M. Borodovsky
Georgia Institute of Technology
School of Biology
Atlanta, GA 30332-0230











Document 2 : graphe obtenu avec GeneMark sur la séquence d'intérêt (seuil à 0,4)

GeneMark

Version 2.5p (09.08.06)

Copyright 1993 - M. Borodovsky, J. McIninch

PROGRAM INFORMATION

Sequence : seq_L_lactis_lactis
Analysis Date : 11/29/24 at 4:40:54
Pages : 6
Sequence Length : 9557 bp
GC Content : 31.57%

Window Length : 120 bp
Window Step : 12 bp
Threshold Value : 0.500

PS-Version : 1.2

GeneMark Options : PostScript graph,
Mark ORFs / splice sites,
List ORFs,
List regions and/or splice sites,

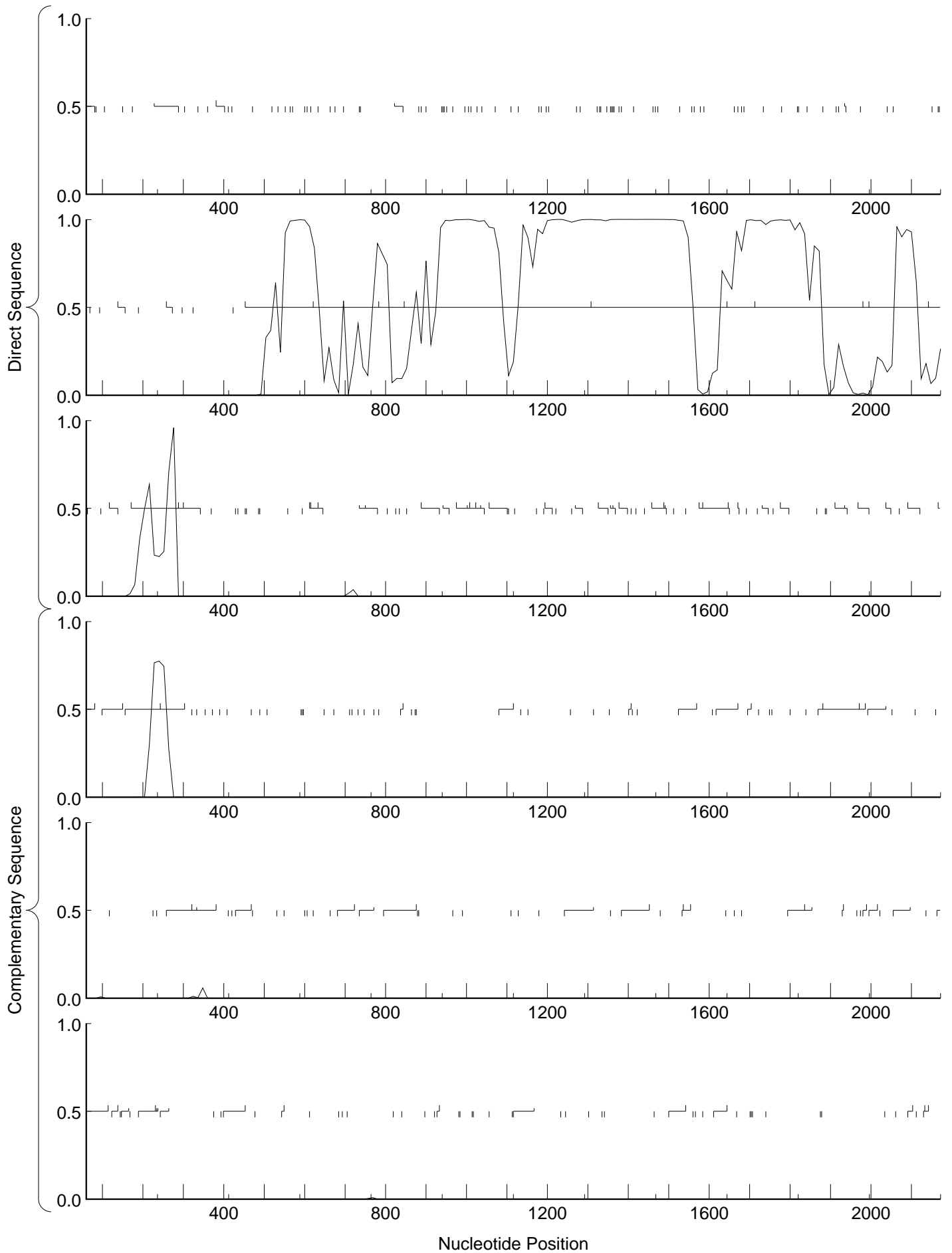
Matrix notes & comments

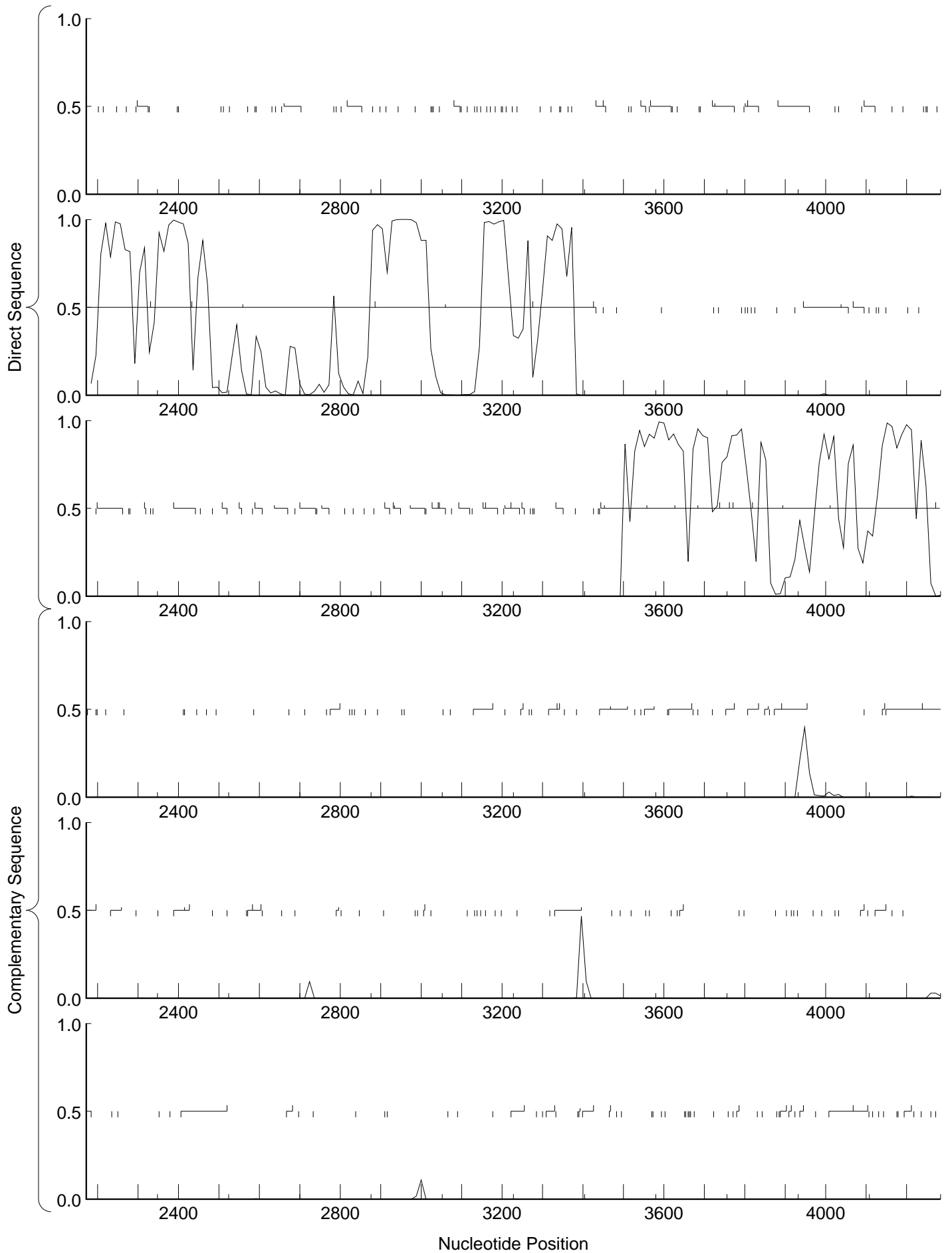
Training set derived by GeneMarkS, 4.27 September 2014
Tue Sep 23 15:01:07 2014

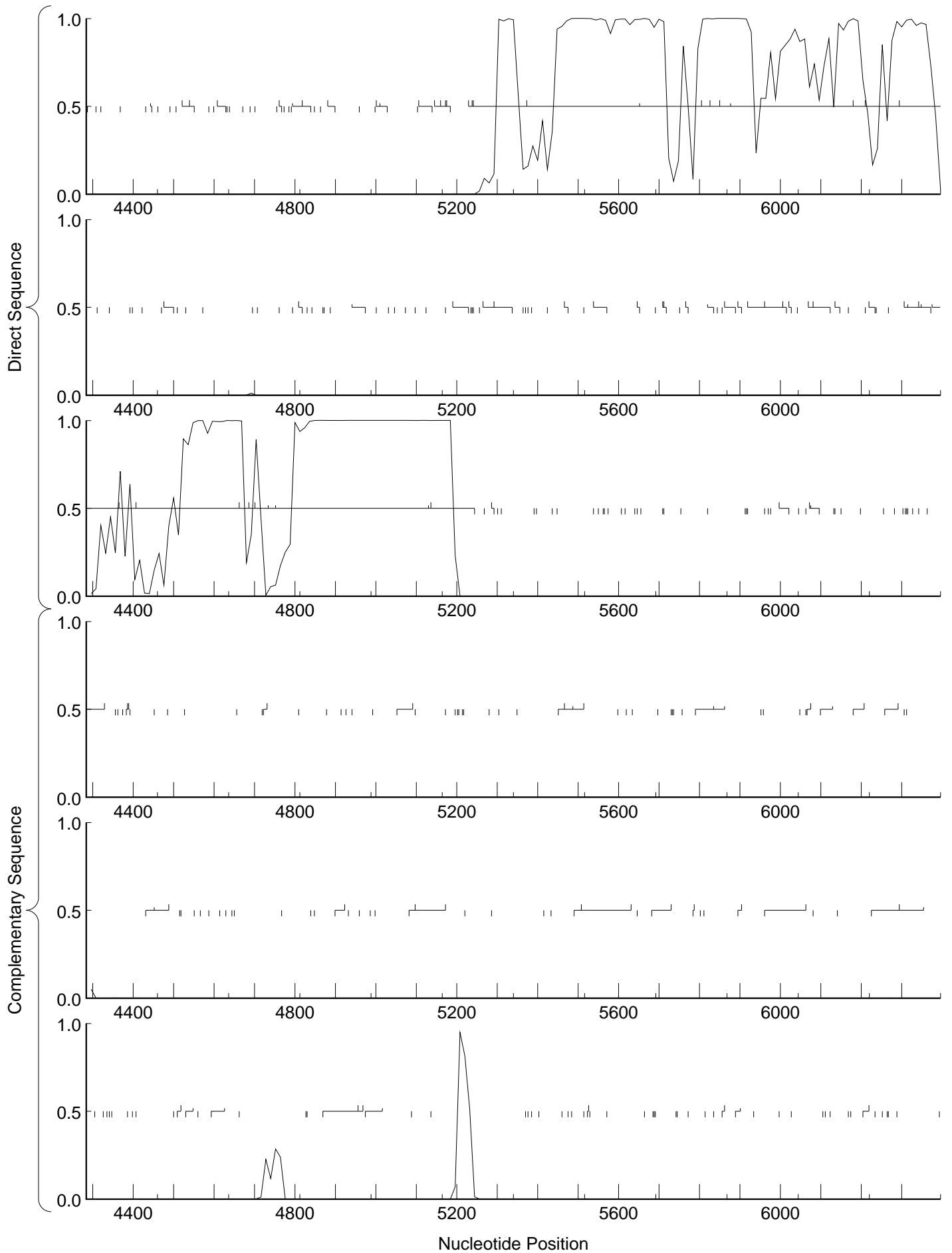
MATRIX INFORMATION

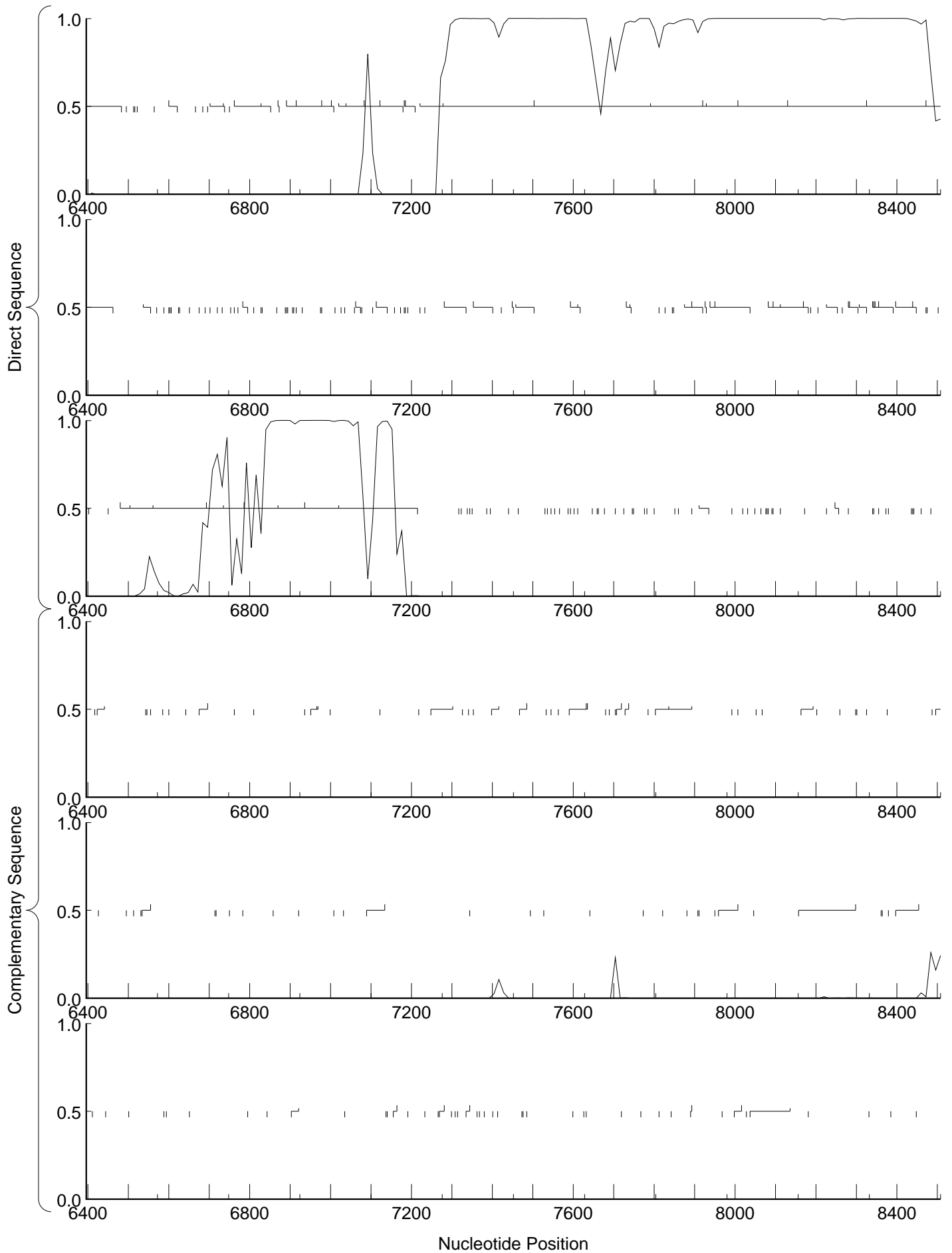
Matrix : Lactococcus_lactis_I11403
Author : -
Order : 4

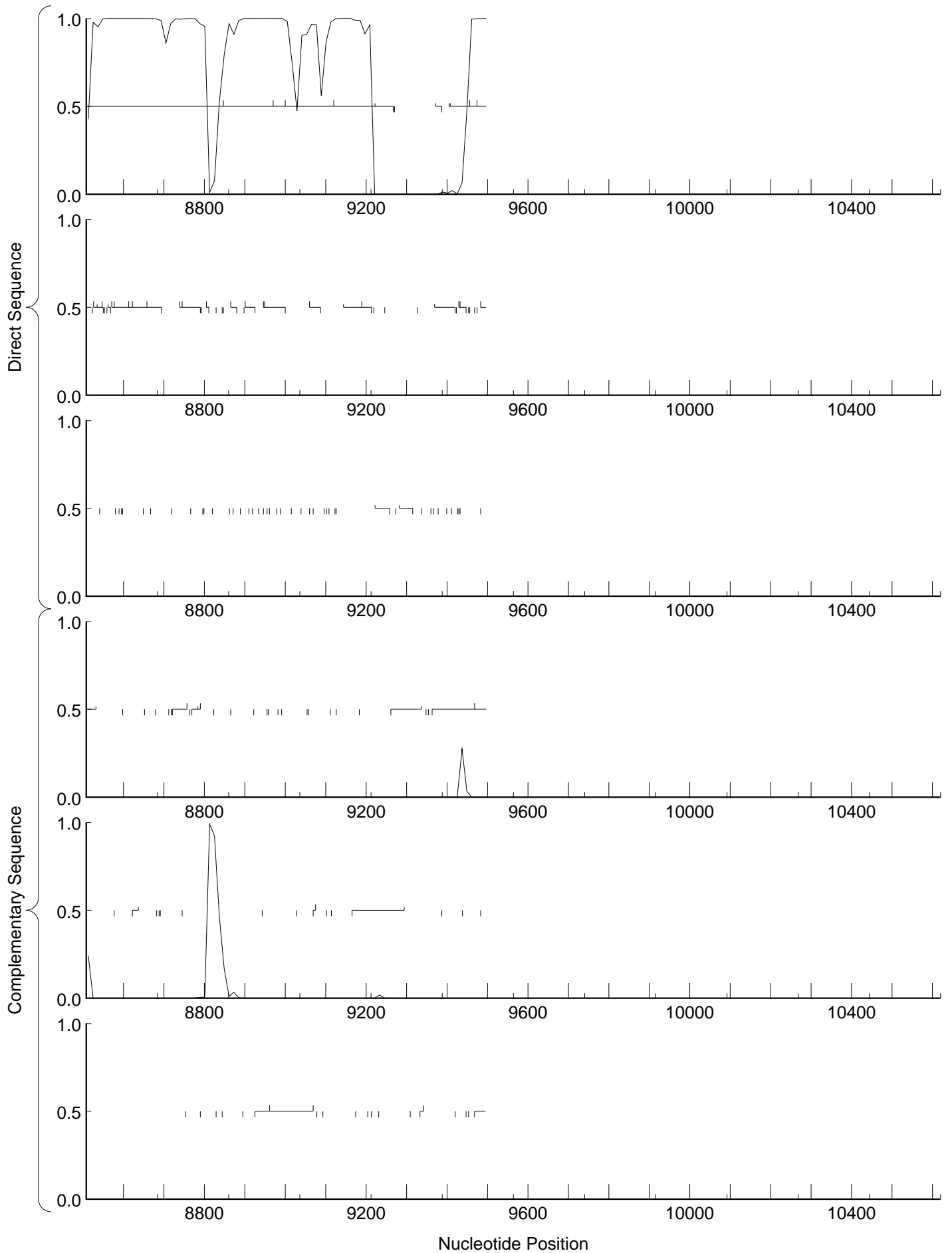
Send questions / comments to:
Dr. M. Borodovsky
Georgia Institute of Technology
School of Biology
Atlanta, GA 30332-0230











Document 3 : graphe obtenu avec GeneMarkHMM sur la séquence d'intérêt

GeneMark

Version 2.5p (09.08.06)

Copyright 1993 - M. Borodovsky, J. McIninch

PROGRAM INFORMATION

Sequence : seq_L_lactis_lactis
Analysis Date : 11/29/24 at 4:52:54
Pages : 6
Sequence Length : 9557 bp
GC Content : 31.57%

Window Length : 96 bp
Window Step : 12 bp
Threshold Value : 0.500

PS-Version : 1.2

GeneMark Options : PostScript graph,
Mark ORFs / splice sites,
List ORFs,
List regions and/or splice sites,

Matrix notes & comments

Training set derived by GeneMarkS, 4.27 September 2014
Tue Sep 23 15:01:07 2014

MATRIX INFORMATION

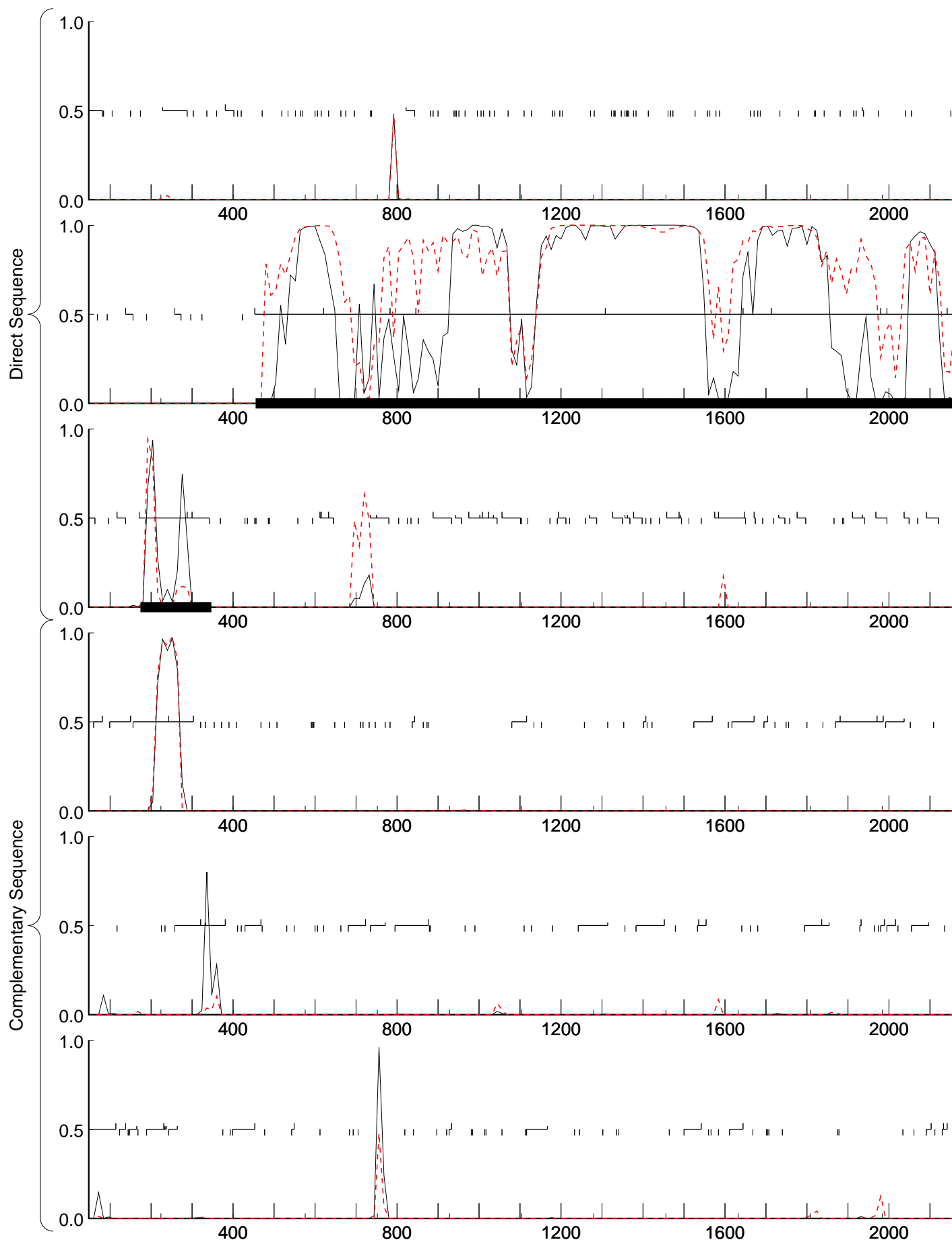
Matrix : Lactococcus_lactis_I11403
Author : -
Order : 4

Send questions / comments to:
Dr. M. Borodovsky
Georgia Institute of Technology
School of Biology
Atlanta, GA 30332-0230

GeneMark.hmm prediction

seq_L_lactis_lactis, Order 4, Window 96, Step 12, 2/6

seq_L_lactis_lactis, Order 2, Window 96, Step 12, 2/6



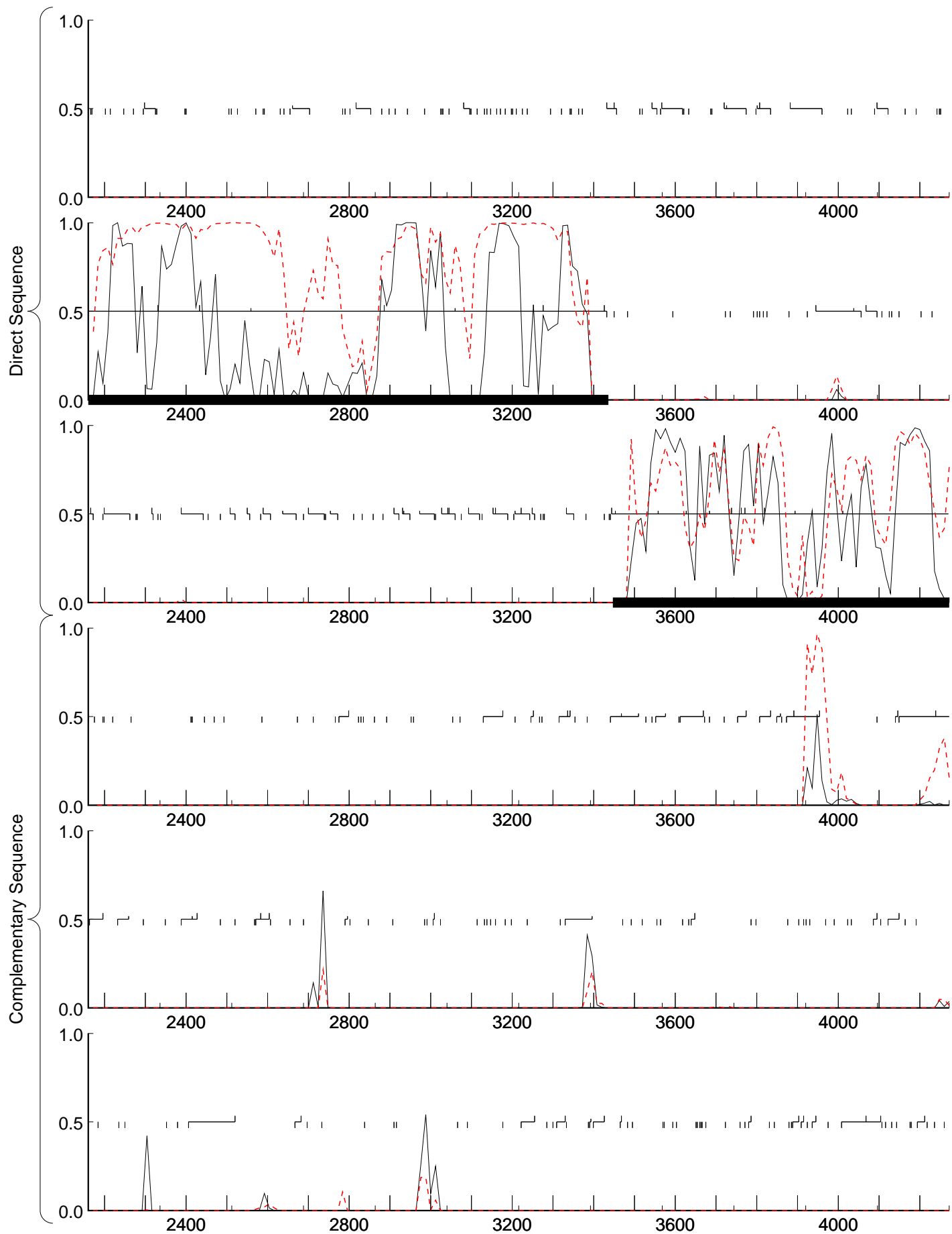
typical.ps

atypical.ps

GeneMark.hmm prediction

seq_L_lactis_lactis, Order 4, Window 96, Step 12, 3/6

seq_L_lactis_lactis, Order 2, Window 96, Step 12, 3/6



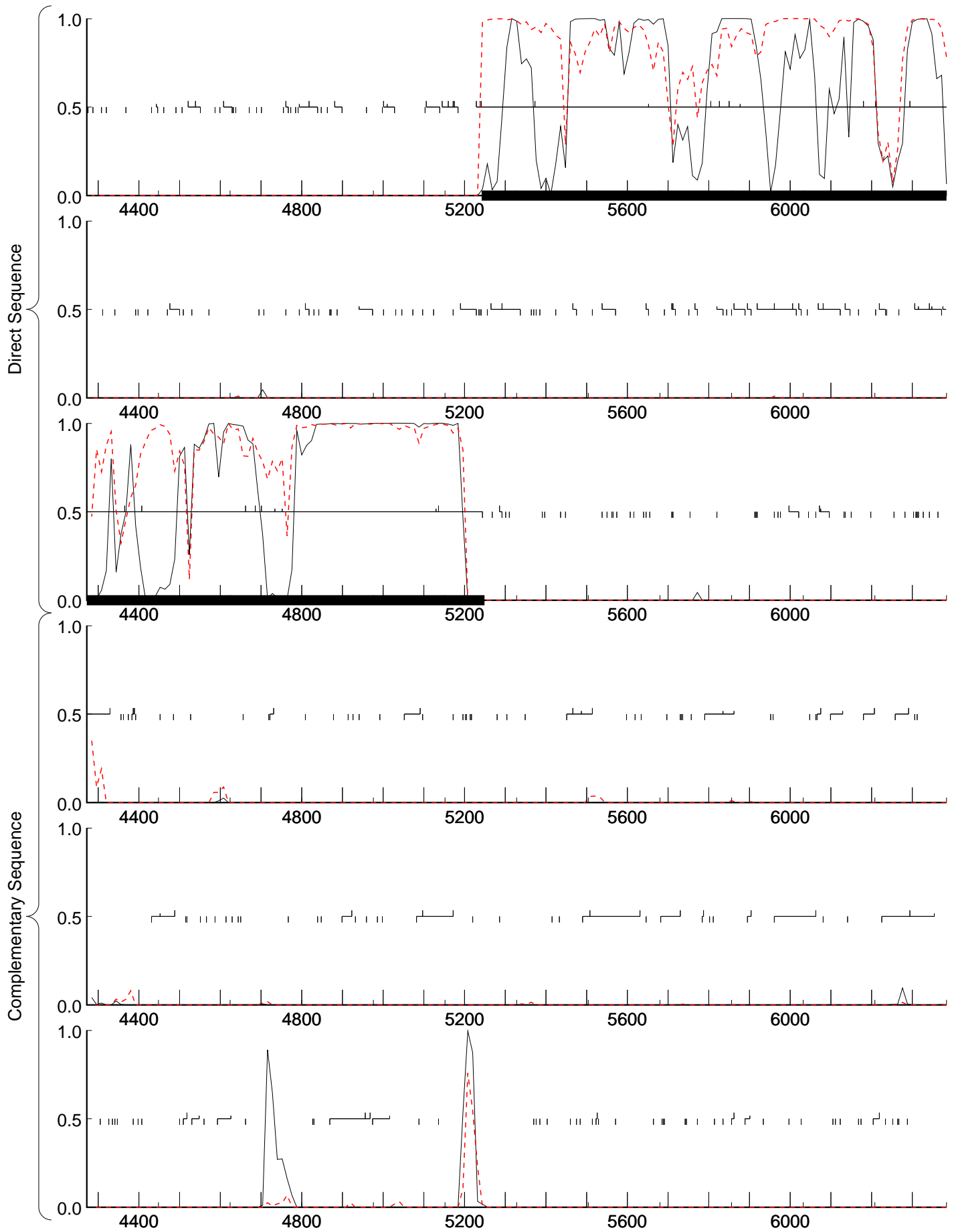
typical.ps

atypical.ps

GeneMark.hmm prediction

seq_L_lactis_lactis, Order 4, Window 96, Step 12, 4/6

seq_L_lactis_lactis, Order 2, Window 96, Step 12, 4/6



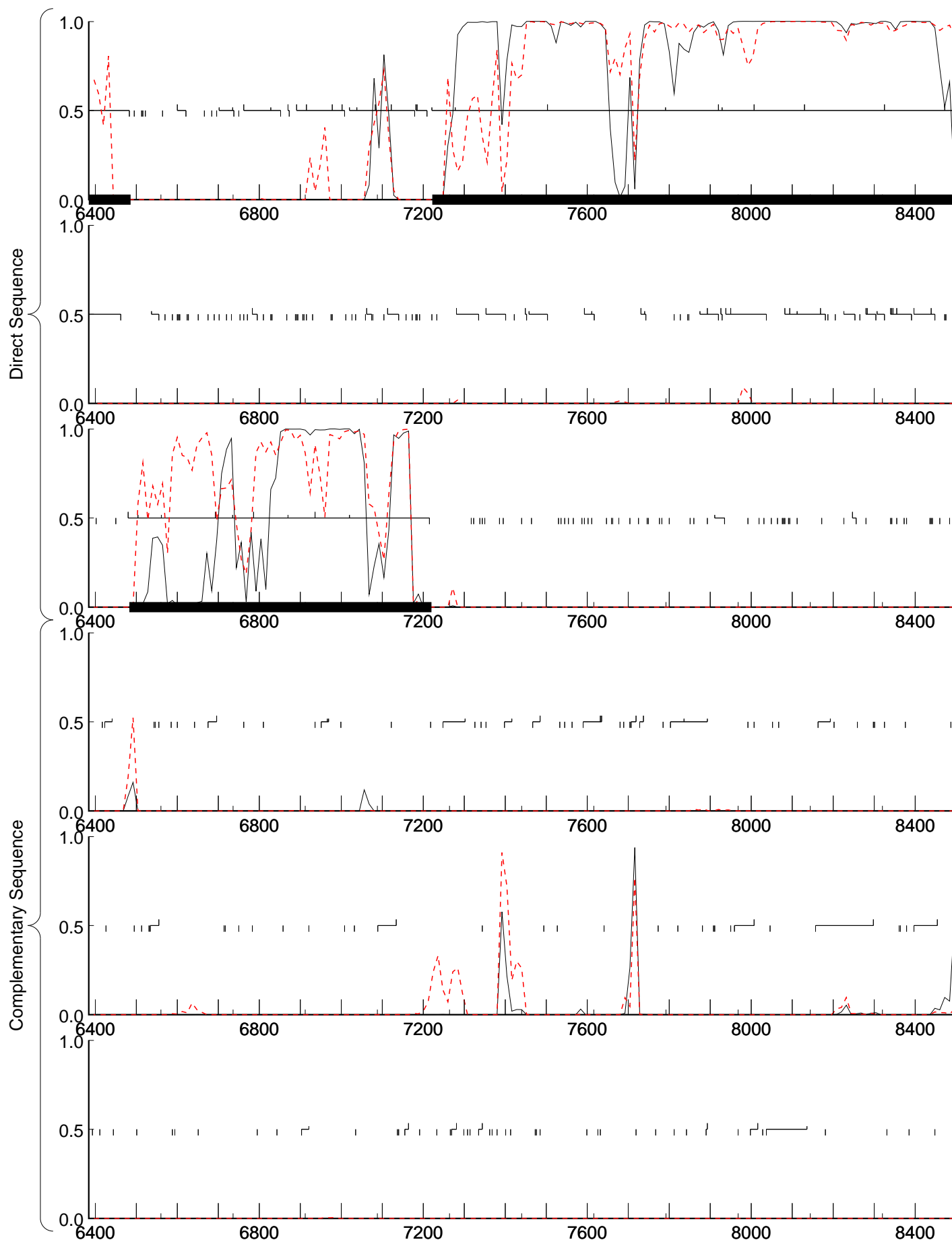
typical.ps

atypical.ps

GeneMark.hmm prediction

seq_L_lactis_lactis, Order 4, Window 96, Step 12, 5/6

seq_L_lactis_lactis, Order 2, Window 96, Step 12, 5/6



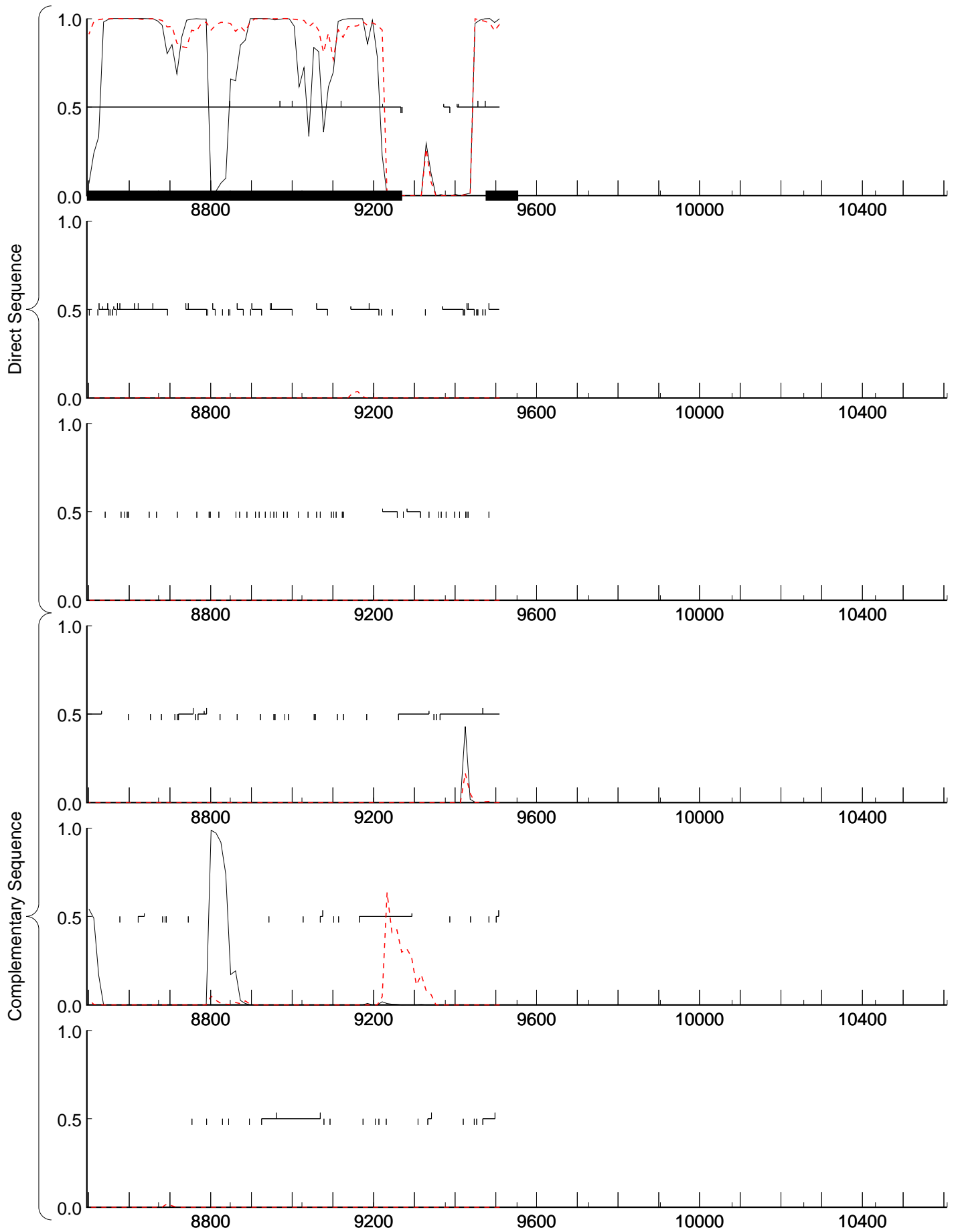
typical.ps

atypical.ps

GeneMark.hmm prediction

seq_L_lactis_lactis, Order 4, Window 96, Step 12, 6/6

seq_L_lactis_lactis, Order 2, Window 96, Step 12, 6/6



typical.ps

Nucleotide Position

atypical.ps

Document 4 : code python pour la conversion des fichiers en format GFF

- `convert_to_GFF.py`

programme principal permettant la lecture du fichier et l'écriture d'un fichier de sortie

- `search_motif.py`

module complémentaire permettant la recherche de motif spécifique afin de créer l'en-tête du fichier de sortie

- `search_data.py`

module complémentaire recherchant les différentes informations du fichier d'entrée et permettant la conversion au format GFF

- `man.py`

module complémentaire recensant les arguments, les erreurs possibles et les bonnes pratiques du programme (manuel d'utilisation)

Nom de la methode : `Error_arg()`

Description : Affiche un message d'erreur pour arguments invalides ou insuffisants et termine le programme.

En-tete : Pas de retour.

Variables locales : Aucune.

Instructions :

1. TERMINER avec le message : "ERR_ARG ! Bad arguments given ... " suivi de la syntaxe correcte d'utilisation.
 2. FIN (`Error_arg`).
-

Nom de la methode : `Error_type()`

Description : Affiche un message d'erreur pour type de fichier invalide et termine le programme.

En-tete : Pas de retour.

Variables locales : Aucune.

Instructions :

1. TERMINER avec le message : "ERR_TYPE ! Wrong type given ... " suivi de la syntaxe correcte d'utilisation.
 2. FIN (`Error_type`).
-

Nom de la methode : `initialise()`

Description : Initialise les parametres en recuperant les arguments de la ligne de commande.

En-tete : Retourne une liste : [filename, fileout, categorie, type].

Variables locales :

- filename (chaîne de caracteres)
- fileout (chaîne de caracteres)
- categorie (chaîne de caracteres)
- type (chaîne de caracteres)

Instructions :

1. SI le nombre d'arguments est > 1 ET le premier argument est "--help" ou "-h" :
 APPELER `man.help()`.
2. Sinon, essayer :
 - AFFECTER type avec le 2e argument (converti en majuscules).

- AFFECTER filename avec le 3e argument.
- SI type == "SFM" :
 - AFFECTER categorie avec le 4e argument (converti en majuscules).
 - SI categorie ≠ "RBS", "PROM", "TERM" :
 - APPELER Error_arg().
 - TENTER d'affecter fileout avec le 5e argument ; sinon, "output.GFF".
 - SINON :
 - AFFECTER categorie = "None".
 - TENTER d'affecter fileout avec le 4e argument ; sinon, "output.GFF".
- 3. Retourner [filename, fileout, categorie, type].
- 4. FIN (initialise).

Nom de la methode : lire(files)

Description : Lit le contenu d'un fichier specifie.

En-tete : Retourne une chaine representant le contenu du fichier.

Variables locales :

- filename (chaine de caracteres)
- input (chaine de caracteres)

Instructions :

1. AFFECTER filename avec le premier element de files.
 2. OUVRIR filename en mode lecture.
 3. LIRE le contenu du fichier et le stocker dans input.
 4. FERMER le fichier.
 5. Retourner input.
 6. FIN (lire).
-

Nom de la methode : ecrire(files, input)

Description : Ecrit les donnees traitees dans un fichier de sortie au format GFF.

En-tete : Pas de retour.

Variables locales :

- type (chaine de caracteres)

Instructions :

1. AFFECTER type avec le 4e element de files.
2. SI type == "GENM" :

```

        APPELER GeneMark(files, input).
3. SINON SI type == "GENMH" :
        APPELER GeneMarkHMM(files, input).
4. SINON :
        APPELER ScanForMatches(files, input).
5. FIN (ecrire).

```

Nom de la methode : ScanForMatches(files, input)

Description : Traite et ecrit les donnees d'un fichier ScanForMatches au format GFF.

En-tete : Pas de retour.

Variables locales :

- fileout (chaîne de caracteres)
- category (chaîne de caracteres)
- RBS (booléen)
- Date (liste de chaînes de caracteres)
- Info (liste contenant taille et DataFrame)
- taille (entier)
- DataFrame (tableau de donnees)

Instructions :

1. AFFECTER fileout avec le deuxieme element de files.
2. AFFECTER category avec le troisieme element de files.
3. SI category == "RBS" :
 - AFFECTER RBS = True.
 SINON :
 - AFFECTER RBS = False.
4. TENTER :
 - EXTRAIRE Date en appelant Date().
 - EXTRAIRE Info en appelant Create_Tab_SFM(input).
 - AFFECTER taille avec Info[0] et DataFrame avec Info[1].
 EN CAS D'ERREUR :
 - APPELER Error_type().
5. OUVRIR fileout en mode ecriture.
6. ECRIRE dans fileout :
 - "##gff-version 2\n".
 - "##source-version ScanForMatches".
 - Date formatee avec :


```

          ##date: " + Date[0] + " " + Date[1] + " " + Date[2] + " " +
          Date[3] + ":" + Date[4] + ":" + Date[5] + " " + Date[6].
          
```
 - "# Sequence file name: -".
 - "# Model file name: -".
 - "# RBS: " + str(RBS).
 - "# Model information: -\n\n".

```

7. POUR chaque index dans [0, taille-1] :
    ECRIRE dans fileout :
        DataFrame[index][0] + "\tScanForMatches\t" + category + "\t" +
        DataFrame[index][1] + "\t" + DataFrame[index][2] + "\t.\t+\t.\tnote \" \" +
        DataFrame[index][3] + "\"\n".
8. FERMER fileout.
9. FIN (ScanForMatches).

```

Nom de la methode : GeneMark(files, input)

Description : Traite et ecrit les donnees d'un fichier GeneMark au format GFF.

En-tete : Pas de retour.

Variables locales :

- fileout (chaîne de caracteres)
- SourceVersion (chaîne de caracteres)
- Date (liste de chaînes de caracteres)
- Seqfilename (chaîne de caracteres)
- ModelInformation (chaîne de caracteres)
- Seq (chaîne de caracteres)
- Info (liste contenant taille et DataFrame)
- taille (entier)
- DataFrame (tableau de donnees)

Instructions :

1. AFFECTER fileout avec le deuxieme element de files.
2. TENTER :
 - EXTRAIRE SourceVersion en appelant SourceVersion_GM(input).
 - EXTRAIRE Date en appelant Date().
 - EXTRAIRE Seqfilename en appelant Seqfilename_GM(input).
 - EXTRAIRE ModelInformation en appelant ModelInformation_GM(input).
 - EXTRAIRE Seq en appelant Seq_GM(input).
 - EXTRAIRE Info en appelant Create_Tab_GM(input).
 - AFFECTER taille avec Info[0] et DataFrame avec Info[1].
- EN CAS D'ERREUR :
 - APPELER Error_type().
3. OUVRIR fileout en mode ecriture.
4. ECRIRE dans fileout :
 - "##gff-version 2\n".
 - "##source-version GeneMark " + SourceVersion.
 - Date formatee avec :
 - "##date: " + Date[0] + " " + Date[1] + " " + Date[2] + " " +
 - Date[3] + ":" + Date[4] + ":" + Date[5] + " " + Date[6].
 - "# " + Seqfilename.
 - "# Model file name: -".

- "# RBS: -".
- "# Model information: " + ModelInformation + "\n\n".

5. POUR chaque index dans [0, taille-1] :

 Ecrire dans fileout :

 Seq + "\tGeneMark\t" + DataFrame[index][6] + "\t" + str(DataFrame[index][1]) + "\t" + str(DataFrame[index][2]) + "\t?\t" + DataFrame[index][0] + "\t.\t" + str(DataFrame[index][7]).

6. FERMER fileout.

7. FIN (GeneMark).

Nom de la methode : GeneMarkHMM(files, input)

Description : Traite et ecrit les donnees d'un fichier GeneMarkHMM au format GFF.

En-tete : Pas de retour.

Variables locales :

- fileout (chaîne de caracteres)
- SourceVersion (liste de deux chaînes de caracteres)
- Date (chaîne de caracteres)
- Seqfilename (chaîne de caracteres)
- Model (chaîne de caracteres)
- RBS (chaîne de caracteres)
- ModelInformation (chaîne de caracteres)
- Seq (chaîne de caracteres)
- Info (liste contenant taille et DataFrame)
- taille (entier)
- DataFrame (tableau de donnees)

Instructions :

1. AFFECTER fileout avec le deuxieme element de files.
2. TENTER :
 - EXTRAIRE SourceVersion en appelant SourceVersion_GMH(input).
 - EXTRAIRE Date en appelant Date_GMH(input).
 - EXTRAIRE Seqfilename en appelant Seqfilename_GMH(input).
 - EXTRAIRE Model en appelant Model_GMH(input).
 - EXTRAIRE RBS en appelant RBS_GMH(input).
 - EXTRAIRE ModelInformation en appelant ModelInformation_GMH(input).
 - EXTRAIRE Seq en appelant Seq_GMH(input).
 - EXTRAIRE Info en appelant Create_Tab_GMH(input).
 - AFFECTER taille avec Info[0] et DataFrame avec Info[1].
- EN CAS D'ERREUR :
 - APPELER Error_type().
3. OUVRIR fileout en mode ecriture.
4. ECRIRE dans fileout :
 - "##gff-version 2\n".

- "##source-version " + SourceVersion[0] + " " + SourceVersion[1].
- "##date: " + Date.
- "# " + Seqfilename.
- "# " + Model.
- "# " + RBS.
- "# " + ModelInformation + "\n\n".

5. POUR chaque index dans [0, taille-1] :

 ECRIRE dans fileout :

 Seq + "_" + str((index+1)) + "\tGeneMark.hmm\tCDS\t" +
 DataFrame[index][1] + "\t" + DataFrame[index][2] + "\t?\t" +
 DataFrame[index][0] + "\t.\tgene GMH_CDS_" + str((index+1)).

6. FERMER fileout.

7. FIN (GeneMarkHMM).

Nom de la methode : SourceVersion_GMH(input)

Description : Extraire le nom de la source et la version depuis une chaîne d'entrée.

En-tete : Retourne un tuple contenant deux chaînes de caractères (source, version).

Variables locales :

- motif_sourceversion (expression reguliere)
- matchSourceversion (objet de correspondance)
- source, version (chaines de caracteres)

Instructions :

1. INITIALISER motif_sourceversion avec l'expression reguliere "(?P<source>.{24})(?P<other>.{10})(?P<version>\\d{1,2}\\.{1,3})".
 2. RECHERCHER une correspondance de motif_sourceversion dans input.
 3. SI une correspondance est trouvee :
 - AFFECTER source = valeur du groupe 'source'.
 - AFFECTER version = valeur du groupe 'version'.
 - SINON :
 - AFFECTER source = "".
 - AFFECTER version = "".
 4. RETOURNER (source, version).
 5. FIN (SourceVersion_GMH).
-

Nom de la methode : Date_GMH(input)

Description : Extraire la date depuis une chaîne d'entrée.

En-tete : Retourne une chaîne représentant la date ou une chaîne vide.

Variables locales :

- motif_date (expression reguliere)
- matchDate (objet de correspondance)
- date (chaîne de caracteres)

Instructions :

1. INITIALISER motif_date avec l'expression reguliere "(Date:)(?P<date>.{24,25})".
2. RECHERCHER une correspondance de motif_date dans input.
3. SI une correspondance est trouvee :
 - AFFECTER date = valeur du groupe 'date'.
- SINON :
 - AFFECTER date = "".
4. RETOURNER date.
5. FIN (Date_GMH).

Nom de la methode : Seqfilename_GMH(input)

Description : Extraire le nom du fichier de séquence depuis une chaîne d'entrée.

En-tete : Retourne une chaîne représentant le nom du fichier ou une chaîne vide.

Variables locales :

- motif_Seqfilename (expression reguliere)
- matchSeqfilename (objet de correspondance)
- Seqfilename (chaîne de caracteres)

Instructions :

1. INITIALISER motif_Seqfilename avec l'expression reguliere "(?P<Seqfilename>Sequence\\sfile\\sname:.{1,}\\s\\.fna)".
2. RECHERCHER une correspondance de motif_Seqfilename dans input.
3. SI une correspondance est trouvee :
 AFFECTER Seqfilename = valeur du groupe 'Seqfilename'.
 SINON :
 AFFECTER Seqfilename = "".
4. RETOURNER Seqfilename.
5. FIN (Seqfilename_GMH).

Nom de la methode : Model_GMH(input)

Description : Extraire le nom du fichier modèle depuis une chaîne d'entrée.

En-tete : Retourne une chaîne représentant le nom du fichier modèle ou une chaîne vide.

Variables locales :

- motif_Model (expression reguliere)
- matchModel (objet de correspondance)
- Model (chaîne de caracteres)

Instructions :

1. INITIALISER motif_Model avec l'expression reguliere "(?P<Model>Model\\sfile\\sname:.{1,})".
2. RECHERCHER une correspondance de motif_Model dans input.
3. SI une correspondance est trouvee :
 AFFECTER Model = valeur du groupe 'Model'.
 SINON :
 AFFECTER Model = "".
4. RETOURNER Model.
5. FIN (Model_GMH).

Nom de la methode : RBS_GMH(input)

Description : Extraire les informations RBS depuis une chaîne d'entrée.
En-tete : Retourne une chaîne représentant les informations RBS ou une chaîne vide.

Variables locales :

- motif_RBS (expression reguliere)
- matchRBS (objet de correspondance)
- RBS (chaîne de caracteres)

Instructions :

1. INITIALISER motif_RBS avec l'expression reguliere "(?P<RBS>RBS:\\s\\w{4,5})".
2. RECHERCHER une correspondance de motif_RBS dans input.
3. SI une correspondance est trouvee :
 AFFECTER RBS = valeur du groupe 'RBS'.
 SINON :
 AFFECTER RBS = "".
4. RETOURNER RBS.
5. FIN (RBS_GMH).

Nom de la methode : ModelInformation_GMH(input)

Description : Extraire les informations du modèle depuis une chaîne d'entrée.

En-tete : Retourne une chaîne représentant les informations du modèle ou une chaîne vide.

Variables locales :

- motif_ModelInformation (expression reguliere)
- matchModelInformation (objet de correspondance)
- ModelInformation (chaîne de caracteres)

Instructions :

1. INITIALISER motif_ModelInformation avec l'expression reguliere "(?P<ModelInformation>Model\\sinformation:.{1,})".
2. RECHERCHER une correspondance de motif_ModelInformation dans input.
3. SI une correspondance est trouvee :
 AFFECTER ModelInformation = valeur du groupe 'ModelInformation'.
 SINON :
 AFFECTER ModelInformation = "".
4. RETOURNER ModelInformation.

5. FIN (ModelInformation_GMH).

Nom de la methode : Seq_GMH(input)

Description : Extrait la ligne de definition FASTA a partir des donnees GeneMarkHMM.

En-tete : Retourne une chaine representant la sequence ou une chaine vide.

Variables locales :

- motif_Seq (expression reguliere)
- matchSeq (objet de correspondance)
- Seq (chaine de caracteres)

Instructions :

1. INITIALISER motif_Seq avec l'expression reguliere correspondante.
 2. RECHERCHER une correspondance de motif_Seq dans input.
 3. SI une correspondance est trouvee :
 AFFECTER Seq = valeur du groupe 'Seq'.
 SINON :
 AFFECTER Seq = "".
 4. RETOURNER Seq.
 5. FIN (Seq_GMH).
-

Nom de la methode : Create_Tab_GMH(input)

Description : Analyse les donnees GeneMarkHMM pour creer une table structuree des caracteristiques de sequence.

En-tete : Retourne une liste contenant le nombre total de caracteristiques et un dictionnaire des details.

Variables locales :

- tab (dictionnaire)
- n (entier initialise a 1)
- motifTab (expression reguliere)
- matchLine (objet de correspondance)

Instructions :

1. INITIALISER motifTab avec l'expression reguliere pour les caracteristiques de sequence.
2. TANT QUE matchLine correspond a motifTab dans input :
 AFFECTER Line avec [Strand, LeftEnd, RightEnd] des groupes de matchLine.
 AJOUTER Line dans tab avec la cle (n-1).
 INCREMENTER n.
 METTRE A JOUR motifTab pour le prochain index n.
3. RETOURNER [n-1, tab].
4. FIN (Create_Tab_GMH).

Nom de la methode : Seq_GM(input)

Description : Extrait la ligne de sequence a partir des donnees GeneMark.

En-tete : Retourne une chaine representant la sequence ou une chaine vide.

Variables locales :

- motif_Seq (expression reguliere)
- matchSeq (objet de correspondance)
- Seq (chaine de caracteres)

Instructions :

1. INITIALISER motif_Seq avec l'expression reguliere correspondante.
 2. RECHERCHER une correspondance de motif_Seq dans input.
 3. SI une correspondance est trouvee :
 - AFFECTER Seq = valeur du groupe 'Seq'.
 - SINON :
 - AFFECTER Seq = "".
 4. RETOURNER Seq.
 5. FIN (Seq_GM).
-

Nom de la methode : Create_Tab_GM(input)

Description : Analyse les donnees GeneMark pour creer une table des informations de genes structurees.

En-tete : Retourne une liste contenant le nombre total de genes valides et un tableau des details des genes.

Variables locales :

- tab (liste)
- taille (entier initialise a 0)
- index, i (entiers)
- Line (liste de chaines de caracteres)
- motifTab (expression reguliere)
- matchLine (objet iterateur de correspondances)

Instructions :

1. INITIALISER motifTab avec l'expression reguliere pour les lignes de donnees GeneMark.
2. POUR chaque element dans matchLine :
 - CALCULER strand en fonction du groupe 'Strand'.
 - AJOUTER Line dans tab avec les informations des groupes de correspondance.
 - INCREMENTER taille si Proba \leq 1.
3. OPTIMISER les probabilites pour marquer les "ATG" et gerer les

annotations des genes.

4. AJOUTER des copies pour les sequences "ATG" et ajuster les annotations des autres.

5. RETOURNER [taille, tab].

6. FIN (Create_Tab_GM).

Nom de la methode : Create_Tab_SFM(input)

Description : Analyse les donnees ScanForMatches pour creer une table des caracteristiques de sequence.

En-tete : Retourne une liste contenant le nombre total de caracteristiques et un dictionnaire des details.

Variables locales :

- tab (dictionnaire)
- n (entier initialise a 0)
- motifTab (expression reguliere)
- matchLine (objet iterateur de correspondances)

Instructions :

1. INITIALISER motifTab avec l'expression reguliere pour les caracteristiques de sequence.

2. POUR chaque element dans matchLine :

AJOUTER une entree dans tab avec les details de la sequence, des extremités et des notes.

INCREMENTER n.

3. RETOURNER [n, tab].

4. FIN (Create_Tab_SFM).

Nom de la methode : pause()

Description : Pause l'execution du programme et attend une action de l'utilisateur.

En-tete : Pas de retour.

Variables locales :

- response (chaîne de caracteres)

Instructions :

1. LIRE response avec le message : "Appuyez sur une touche pour continuer ou 'q' pour quitter".
 2. SUPPRIMER la ligne de saisie du terminal avec les commandes d'effacement du curseur.
 3. SI response en minuscules est egal a 'q' :
 TERMINER le programme.
 4. FIN (pause).
-

Nom de la methode : help()

Description : Affiche un manuel utilisateur en plusieurs sections avec des pauses pour la lecture.

En-tete : Pas de retour.

Variables locales : Aucune.

Instructions :

1. AFFICHER le titre : "User Manual for 'convert_to_GFF.py'".
2. AFFICHER les sections suivantes :
 - Purpose : Explique le but du script et les formats pris en charge (GENM, GENMH, SFM).
 - Features : Liste les fonctionnalites comme la validation des fichiers et le format de sortie GFF.
 - Usage : Explique la syntaxe pour l'execution du script.
 - Arguments : Detaille les arguments obligatoires et optionnels.
 - Exemples : Fournit des exemples d'utilisation pour les differents formats.
 - Output : Decrit le contenu du fichier de sortie (annotations, metadonnees).
 - Error Messages : Donne les messages d'erreur pour des arguments ou formats invalides.
 - How It Works : Resume les etapes de traitement du script (initialisation, lecture, traitement, ecriture).
 - Additional Modules : Liste les dependances et leurs fonctionnalites ('search_data.py', 'search_motif.py').
 - Contact : Donne les informations pour contacter le developpeur.

3. APPELER `pause()` apres chaque section.
4. TERMINER le programme apres l'affichage complet.
5. FIN (`help`).

Document 5 : résultats obtenus après conversion des fichiers de sortie à l'aide du code fourni (document complémentaire 4)

- output_GENMH.GFF

fichier de sortie obtenu après conversion du fichier input_GENMH.LSF (résultats de l'étude menée avec GeneMarkHMM) avec le programme convert_to_GFF.py

- output_GENM.GFF

fichier de sortie obtenu après conversion du fichier input_GENM.txt (résultats de l'étude menée avec GeneMark) avec le programme convert_to_GFF.py

- output_SFM.GFF

fichier de sortie obtenu après conversion du fichier input_SFM.txt (résultats de l'étude menée avec ScanForMatches) avec le programme convert_to_GFF.py

```
##gff-version 2
##source-version GeneMark.hmm PROKARYOTIC 3.42
##date: Sun Dec 1 13:49:12 2024
# Sequence file name: seq.fna
# Model file name:
/home/genemark/parameters/prokaryotic/Lactococcus_lactis_Il1403/GeneMark_hmm_combined.mod
# RBS: true
# Model information: Lactococcus_lactis_Il1403
```

seq_L_lactis_lactis_1	GeneMark.hmm	CDS	174	347	?	+
. gene GMH_CDS_1						
seq_L_lactis_lactis_2	GeneMark.hmm	CDS	455	3436	?	+
. gene GMH_CDS_2						
seq_L_lactis_lactis_3	GeneMark.hmm	CDS	3447	5249	?	+
. gene GMH_CDS_3						
seq_L_lactis_lactis_4	GeneMark.hmm	CDS	5242	6486	?	+
. gene GMH_CDS_4						
seq_L_lactis_lactis_5	GeneMark.hmm	CDS	6483	7220	?	+
. gene GMH_CDS_5						
seq_L_lactis_lactis_6	GeneMark.hmm	CDS	7222	9270	?	+
. gene GMH_CDS_6						
seq_L_lactis_lactis_7	GeneMark.hmm	CDS	9475	>9555	?	+
. gene GMH_CDS_7						

```
##gff-version 2
##source-version GeneMark 4.27
##date: Wed Jan 01 13:44:31 2025
# Sequence file: seq.fna
# Model file name: -
# RBS: -
# Model information: Lactococcus_lactis_Il1403
```

seq_L_lactis_lactis	GeneMark	CDS	455	3436	?	+
. gene GM_CDS_1.1						
seq_L_lactis_lactis	GeneMark	CDS	623	3436	?	+
. gene GM_CDS_1.2						
seq_L_lactis_lactis	GeneMark	CDS	785	3436	?	+
. gene GM_CDS_1.3						
seq_L_lactis_lactis	GeneMark	ATG	848	850	?	+
. .						
seq_L_lactis_lactis	GeneMark	CDS	848	3436	?	+
. gene GM_CDS_1.4						
seq_L_lactis_lactis	GeneMark	CDS	1310	3436	?	+
. gene GM_CDS_1.5						
seq_L_lactis_lactis	GeneMark	ATG	3447	3449	?	+
. .						
seq_L_lactis_lactis	GeneMark	CDS	3447	5249	?	+
. gene GM_CDS_2.1						
seq_L_lactis_lactis	GeneMark	CDS	3456	5249	?	+
. gene GM_CDS_2.2						
seq_L_lactis_lactis	GeneMark	CDS	3561	5249	?	+
. gene GM_CDS_2.3						
seq_L_lactis_lactis	GeneMark	CDS	3630	5249	?	+
. gene GM_CDS_2.4						
seq_L_lactis_lactis	GeneMark	CDS	3687	5249	?	+
. gene GM_CDS_2.5						
seq_L_lactis_lactis	GeneMark	CDS	5230	6486	?	+
. gene GM_CDS_3.1						
seq_L_lactis_lactis	GeneMark	CDS	5239	6486	?	+
. gene GM_CDS_3.2						
seq_L_lactis_lactis	GeneMark	CDS	5242	6486	?	+
. gene GM_CDS_3.3						
seq_L_lactis_lactis	GeneMark	CDS	5374	6486	?	+
. gene GM_CDS_3.4						
seq_L_lactis_lactis	GeneMark	CDS	5653	6486	?	+
. gene GM_CDS_3.5						
seq_L_lactis_lactis	GeneMark	ATG	5806	5808	?	+
. .						
seq_L_lactis_lactis	GeneMark	CDS	5806	6486	?	+
. gene GM_CDS_3.6						
seq_L_lactis_lactis	GeneMark	CDS	5827	6486	?	+
. gene GM_CDS_3.7						

seq_L_lactis_lactis	GeneMark	CDS	6483	7220	?	+
. gene GM_CDS_4.1						
seq_L_lactis_lactis	GeneMark	ATG	6507	6509	?	+
. .						
seq_L_lactis_lactis	GeneMark	CDS	6507	7220	?	+
. gene GM_CDS_4.2						
seq_L_lactis_lactis	GeneMark	CDS	6564	7220	?	+
. gene GM_CDS_4.3						
seq_L_lactis_lactis	GeneMark	CDS	6696	7220	?	+
. gene GM_CDS_4.4						
seq_L_lactis_lactis	GeneMark	CDS	6738	7220	?	+
. gene GM_CDS_4.5						
seq_L_lactis_lactis	GeneMark	CDS	6789	7220	?	+
. gene GM_CDS_4.6						
seq_L_lactis_lactis	GeneMark	CDS	6873	7220	?	+
. gene GM_CDS_4.7						
seq_L_lactis_lactis	GeneMark	CDS	6939	7220	?	+
. gene GM_CDS_4.8						
seq_L_lactis_lactis	GeneMark	CDS	7222	9270	?	+
. gene GM_CDS_5.1						
seq_L_lactis_lactis	GeneMark	ATG	7279	7281	?	+
. .						
seq_L_lactis_lactis	GeneMark	CDS	7279	9270	?	+
. gene GM_CDS_5.2						
seq_L_lactis_lactis	GeneMark	CDS	7504	9270	?	+
. gene GM_CDS_5.3						
seq_L_lactis_lactis	GeneMark	CDS	7792	9270	?	+
. gene GM_CDS_5.4						
seq_L_lactis_lactis	GeneMark	CDS	9406	9558	?	+
. gene GM_CDS_6.1						
seq_L_lactis_lactis	GeneMark	ATG	9409	9411	?	+
. .						
seq_L_lactis_lactis	GeneMark	CDS	9409	9558	?	+
. gene GM_CDS_6.2						

##gff-version 2
##source-version ScanForMatches
##date: Tue Dec 10 22:30:59 2024
Sequence file name: -
Model file name: -
RBS: True
Model information: -

seq_L_lactis_lactis	ScanForMatches	RBS	161	176	.	+
.	note " GGAGG CACTCAAA ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	282	298	.	+
.	note " GGAGC TCTGATGGG TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	1802	1814	.	+
.	note " GGGGA TATCG TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	2317	2335	.	+
.	note " GGATG TAAATAGTAAG ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	2531	2549	.	+
.	note " GGAGA AAAGGGGAGAG TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	4581	4599	.	+
.	note " GGAGA ATTAAGTCTA TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	5497	5512	.	+
.	note " GGAGC AGCTGGCA TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	5611	5625	.	+
.	note " GAAGG ATTTAAT TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	5796	5808	.	+
.	note " GGAGA ATCAG ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	6340	6352	.	+
.	note " GGAGA TGAAA GTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	6367	6379	.	+
.	note " GAAGG AATAA GTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	6469	6485	.	+
.	note " GGAGG GAAGAGGAA ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	6836	6854	.	+
.	note " GGGGG CGGTCAAAGCT TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	6990	7005	.	+
.	note " GGAGA ATATCAGG ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	7067	7086	.	+
.	note " GGGGA GAATTGATAAGG ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	7205	7224	.	+
.	note " GAAGG TAGGAACTAGA GTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	7449	7462	.	+
.	note " GGATG TAACTG GTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	7593	7606	.	+
.	note " GGATG AGCTAC TTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	7942	7954	.	+
.	note " GGAGA AAGCT ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	8071	8086	.	+
.	note " GGAGG GTTTGATA ATG "					

seq_L_lactis_lactis	ScanForMatches	RBS	8547	8566	.	+
.	note " GGATG TAATAGCCGTAG GTG "					
seq_L_lactis_lactis	ScanForMatches	RBS	8608	8626	.	+
.	note " GGAGC GGATGCAATTT ATG "					
seq_L_lactis_lactis	ScanForMatches	RBS	9397	9411	.	+
.	note " GGAGG TAAAGTG GTG "					