Eukaryotic Gene Prediction

Wei Zhu May 2007

"In nature, nothing is perfect ..."

- Alice Walker

Gene Structure



What is Gene Prediction?

<u>Gene prediction</u> is the problem of parsing a sequence into nonoverlapping coding segments (CDSs) consisting of exons separated by introns.

Signal Sensors

A signal sensor evaluates fixed-length features in DNA. Start codons Stop codons Donor sites Acceptor sites Promoters Poly-A signals

Content Sensors

A *content sensor* evaluates variable-length features which extend from one signal to another:

- Exons
- Introns
- Intergenic regions
- UTRs

Gene Prediction Approaches

Intrinsic (*ab initio*)

 GENSCAN, FGENESH, GeneMark.hmm GlimmerM, Genie;

Extrinsic (similarity-based)

- Spliced alignment: GenomeScan, EuGene, FGENESH+, FGENESH_C, GeneId+, etc;
- Genomic comparison: TwinScan, TWAIN, SLAM, SGP, FGENESH_2, etc;

Integrated

 GeneScope, GeneMachine, JIGSAW, RiceGAAS, Ensembl, EVM etc.

ab initio Gene Prediction

 Adopt a rigorous probabilistic model of sequence structure and choose the most probable parse according to that probabilistic model.

Pros

- Fast and efficient
- Remarkable accuracy at the nucleotide level
- Cons
 - Less than 50% accuracy at the gene level

Development of a Gene Finder

Build the model
Train the model to generate the related parameters
Predict/Evaluate

Imperfect Model

GT.....AGA G T....

1-bp intron

Accuracy Evaluation

Nucleotide level
Exon level
Gene level

Nucleotide/Base Level

Prediction accuracy per base coding/non-coding



Exon Level

 Prediction accuracy with respect to exact prediction of exon start and end points



An **exon** is assumed to be **correctly predicted** if the overlap between actual and predicted exon is greater or equal than a given threshold α .

| $Sn = \frac{\text{number of Correct Exons}}{\text{number of Actual Exons}}$ | Sensitivity |
|--|---------------|
| $Sn = \frac{\text{number of Correct Exons}}{\text{number of Predicted Exons}}$ | Specificity |
| $ME = \frac{\text{number of Missing Exons}}{\text{number of Actual Exons}}$ | (Sensitivity) |
| $WE = \frac{\text{number of Wrong Exons}}{\text{number of Predicted Exons}}$ | (Specificity) |
| | |

Gene/Protein Level

Prediction accuracy with respect to the protein product encoded by the predicted gene

A Simple Calculation

Given *x* accuracy at exon level, the accuracy of the prediction at the gene level is:

P = P (all exons correctly predicted) $=x^n$, where *n* is the number of exons in the gene.

Typically, *x*<90% and *n*=5, then *P* = 0.9x0.9x0.9x0.9x0.9 = 59%

Performance

Species-specific setting ■ GC content Gene density Gene/Exon/Intron length distribution Codon usage Benchmark training data set • test data set

Maize Gene Prediction

Plant Molecular Biology (2005) 57:445–460 DOI: 10.1007/s11103-005-0271-1 © Springer 2005

Evaluation of five *ab initio* gene prediction programs for the discovery of maize genes

Hong Yao^{1,4}, Ling Guo^{1,6,†}, Yan Fu^{1,4,†}, Lisa A. Borsuk^{1,6}, Tsui-Jung Wen², David S. Skibbe^{1,5}, Xiangqin Cui^{1,4,9}, Brian E. Scheffler⁸, Jun Cao^{1,4}, Scott J. Emrich⁶, Daniel A. Ashlock^{3,6} and Patrick S. Schnable^{1,2,4–7,*}

¹Department of Genetics, Development, and Cell Biology (*author for correspondence; e-mail schnable@ iastate.edu); ²Department of Agronomy; ³Department of Mathematics; ⁴Inderdepartmental Graduate Programs in Genetics; ⁵Department of Molecular, Cellular and Developmental Biology; ⁶Department of Electrical and Computer Engineering and Department of Bioinformatics and Computational Biology; ⁷Center for Plant Genomics, Iowa State University, Ames, Iowa 50011-3650; ⁸USDA-ARS, Mid South Area Genomics Facility, Stoneville, MS 38776-0038, USA; ⁹Present address: Department of Biostatistics, Birmingham, AL 35294, USA; [†]these authors contributed equally to this work

Received 16 August 2004, accepted in revised form 6 January 2005

Gene Finders

| Programs | Websites | Trained organisms | Type of pro | Type of prediction | | Algorithm models |
|--------------|---|----------------------|-------------|--------------------|------------|---------------------|
| | | - J | Splice site | Exon | Gene model | |
| FGENESH | http://www.softberry.com/ berry.phtml?topic = fgenesh&group = programs&subgroup = gfind | Monocots | Yes | Yes | Yes | GHMM ^a |
| GeneMark.hmm | http://opal.biology.gatech.edu/ GeneMark/eukhmm.cgi?org=H.sapiens | Maize | Yes | Yes | Yes | GHMM |
| GENSCAN | http://genes.mit.edu/GENSCAN.html | Maize | Yes | Yes | Yes | GHMM |
| GlimmerR | http://www.tigr.org/tdb/glimmerm/ glmr_form.html | Rice | Yes | Yes | Yes | IMM ^b |
| Grail | http://compbio.ornl.gov/Grail-1.3/ | Arabidopsis | Yes | Yes | No | neural networks |

^aGHMM, Generalized Hidden Markov Model. ^bIMM, Interpolated Markov Model.

Accuracy

| Programs | Nucleo | tide level | | Exon level | | | | | | | |
|--------------|--------|------------|------|------------|------|-------------|-----|-----|-----|-----|--|
| | SN | SP | CC | SN | SP | (SN + SP)/2 | PE% | OE% | ME% | WE% | |
| FGENESH | 0.97 | 0.94 | 0.93 | 0.86 | 0.88 | 0.87 | 9.4 | 0 | 4.6 | 3.1 | |
| GeneMark.hmm | 0.92 | 0.93 | 0.89 | 0.69 | 0.80 | 0.75 | 14 | 0 | 19 | 5.4 | |
| GENSCAN | 0.81 | 0.95 | 0.82 | 0.54 | 0.81 | 0.68 | 12 | 0 | 39 | 7.0 | |
| GlimmerR | 0.70 | 0.91 | 0.71 | 0.51 | 0.64 | 0.57 | 23 | 5.8 | 23 | 7.7 | |
| Grail | 0.55 | 0.67 | 0.43 | 0.34 | 0.28 | 0.31 | 33 | 7.7 | 17 | 31 | |

Challenges of Intrinsic Approaches

- Alternative splicing
- Nested/overlapped genes
- Extremely long/short genes
- Extremely long introns
- Extremely short exons
- Non-canonical introns
- Frame-shift errors
- Split start codons (that is, the start codon is split by an intron in the genomic sequence)
- UTR introns
- Non-ATG triplet as the start codon
- Polycistronic genes

Gene Prediction Approaches

Intrinsic (*ab initio*)

 GENSCAN, FGENESH, GeneMark.hmm GlimmerM, Genie;

Extrinsic (similarity-based)

- Spliced alignment: GenomeScan, EuGene, FGENESH+, FGENESH_C, GeneId+, etc;
- Genomic comparison: TwinScan, TWAIN, SLAM, SGP, FGENESH_2, etc;

Integrated

 GeneScope, GeneMachine, JIGSAW, RiceGAAS, Ensembl, etc.

Similarity-based Gene Prediction

EST/cDNA spliced alignment
Protein spliced alignment
Genomic comparison

Intra-genomic
Inter-genomic

EST/cDNA Spliced Alignment

Report for cDNA subcluster: 1258

of cluster: 20541 (annotdb_asmbl_id:10197 coords:116784-120621)

Subcluster view.



(+)10197.m00079 [current(v1)]: fgenesh model (+)asmbl_1573-including gene model (a+/s+) asmbl_1573 FL-containing (a+/s+) gi|32987826|dbj|AK102617.1| FL Oryz (a+/s+) gi|32971071|dbj|AK061053.1| FL Oryz (a+/s+) gi|32970061|dbj|AK060043.1| FL Oryz (a+/s+) gi|25996130|gb|CA766875.1|CA766875 (a+/s+) gi|29642352|gb|CB647359.1|CB647359 (a+/s+) gi|32948412|gb|BP184984.1|BP184984 (a+/s+) gi|2312713|gb|C28868.1|C28868 C2886 (a+/s+) gi|44670232|gb|CR283666.1|CR283666 (a+/s+) gi|32947813|gb|BP184385.1|BP184385 (a+/s+) gi|29642353|gb|CB647360.1|CB647360 (a+/s+) gi|25806693|gb|CA762648.1|CA762648 (a+/s+) gi|25806691|gb|CA762657.1|CA762657 (a+/s+) gi|25806694|gb|CA762649.1|CA762649 (a+/s+) gi|25806692|gb|CA762647.1|CA762647 (a+/s+) gi|27920725|gb|CB096533.1|CB096533 (a+/s+) gi|8857146|gb|AU094464.1|AU094464 A (a+/s+) gi|12622130|gb|AU172343.1|AU172343 (a+/s+) gi|27577026|gb|CA999720.1|CA999720 (a+/s?) gi|32947812|gb|BP184384.1|BP184384 (a-/s?) gi|24208723|gb|AU225750.1|AU225750 (a+/s?) gi|1632063|gb|C19792.1|C19792 C1979

Assembly description

| L | sembly | cdnas | annotations linked | status | |
|---|--------|-------------------------------------|--------------------|--------|--|
| | | -1126221201-61ATT172242-11ATT172242 | | | |

Pros and Cons

Pros

- High accuracy
- Cons
 - Unavailability or incompleteness of transcript sequence data
 - Extra computation to generate alignments
 - Diverse sequence quality
 - Incomplete full-length cDNA
 - Contamination
 - Incorrect sequence orientations

Genomic Comparison



Microsynteny between *M. truncatula* and Arabidopsis Hongyan et al, 2003

Gene Structure of Syntenic and non-Syntenic Homologous Genes



Hongyan et al, 2003

Comparative Analysis of Cereal Gene Structures



Comparative Analysis of Cereal Gene Promoters

| Adh1 | DRE- |
|--|--|
| Zm - 401 3b - 392 0s - 273 Hv - 371 | element G-bo) TCCGAGCTAGCGCAGGCGCATCCGACGGCACG t. 2c |
| Zm -201 3b -199 0s -109 Hv -172 | AGGCGGCCAAACCGCACCCTCCTTC tggga .a.tgtc .c.tga 02-Site |
| Zm 763 3b 684 0s 588 Hv 335 | GC G FTTGACTTGC GC CTT CTTGG CG GC TTAT |
| Zm 1088 3b 1012 0s 1176 Hv 792 | A GTG GA CTTTGACA GA TTTAT |
| Zm 1637 3b 1581 0s 1881 Hv 1444 | TTATCTTGAGATGCTGAGTTACA gt. g=gccg cc.g.c |

Pros and Cons

Pros

- Aid to identify low expressed genes
- Identify genes in multiple species simultaneously
- Aid to identify transcription factor binding sites
- Uncover non-protein coding genes

Cons

- Performance will depend on the evolutionary distance between the compared sequences.
- Exon/intron boundaries may not be conserved

Tiling Array



ARTADE

-ARabidopsis Tiling-Array-based Detection of Exons

| Filter by: keyword • C | Arabidopsis thaliana : 1 |
|--|--|
| | |
| Bookmark Swap Loci Zoom In Zoom Out | 8,848,876 bp 16,577 bp 8,865,453 b |
| Add/Config | 🗷 show all 😔 🧉 |
| gene(+) gene(-) Flower Tiling chip(+) | |
| Tiling chip(-) | |
| Tiling chip(+) Tiling chip(-) | er fighet og hang men og en størte det men som som handelse er som en størte er er som en som er som er som er |
| Tiling chip(+) Tiling chip(-) | ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ |
| Light7day | |
| Tiling chip(+) Tiling chip(-) | ulle hal alala han meigen seine seine na an |

Gene Prediction Approaches

Intrinsic (*ab initio*)

 GENSCAN, FGENESH, GeneMark.hmm GlimmerM, Genie;

Extrinsic (similarity-based)

- Spliced alignment: GenomeScan, EuGene, FGENESH+, FGENESH_C, GeneId+, etc;
- Genomic comparison: TwinScan, TWAIN, SLAM, SGP, FGENESH_2, etc;

Integrated

 GeneScope, GeneMachine, JIGSAW (combiner), RiceGAAS, Ensembl, etc.

Gene Discovery via Multiple Gene Finders







Highest Scoring Path Thru Candidate Exons

THE REPORT OF A DECK OF A DECK 120 00 0 ALL MALADALE STREET, MILLER TT - DESCRIPTION T THE R MINE Real Property TABLE REALING MEMORY AND A DESCRIPTION OF A REAL PLACE REAL PROPERTY AND A DESCRIPTION OF A 1.18 IN A REAL PARTY AND A DESCRIPTION OF A D 10.00.00. THE CALL AND A DESCRIPTION OF A REPORT OF A DESCRIPTION OF A DESCRIPTION OF A DESCRIPTION OF A DESCRIPTION OF A 1.0.04 AND IN THE OWNER OF 188 . IN THE REPORT OF THE ACCOUNT OF THE THE REPORT OF BUILDING 11111 1 111 J. I. I.

TIGR Rice Genome Annotation Pipeline



RiceGAAC



Ensembl Gene Prediction Procedure



Summary

- Nothing is perfect
- Each gene identification approach has its own features and limitations;
- Genome annotation is an on-going process, and the accuracy is being improved along with the accumulation of the evidence data;tRNAsnoRNA

Case Study

| _Ds06g02400 | L0C_0s06g02420 | L0C_0s06g02430 | L0C_0s06g02440 | |
|---|---|-----------------------------------|---|----|
| x domain containing protein, expressed LOC_0s06g02410 expressed protein | ATOZI1, putative, expressed | expressed protein | expressed protein | |
| Rice Gene Hodels socg02400.1 LOC_0s06e02410.1 | LOC_0806g02420.1 | L0C_0805g02430.1 | L0C_0506g02440.1 | |
| ank Annotations | <u> </u> | | | ←⊩ |
| esh Predictions | <+⊡ | | | |
| HarkHMM Predictions | BB-BI | | | |
| nerHNN Predictions | <-□ | <u>0-0-0-0-</u> | | |
| tri-cumn | gi 32969371 dbj AK059353.1 gi 32976789 dbj AK066771.1 | gi 32977990 dbj AK067972.1 | gi 32971570 db j AK061552.1 gi 37991788 db j AK122142.1 gi 32984057 db j AK098848.1 | |
| Kice iranscript Hssenblies TA25631_4530 CX100736 | TA5524_4530 B1806488 AU223362 | TA21232_4530 CF196602 BM038818 | CF321424 BI798155 AU222970 | |
| | B9928593 AU058448 CV729971 | | AU082453 TA11189_4530 BU673373 | |

Sorghum-Rice Synteny and EST Read Pair





Create a Gene Model

| < | .9k 839k | 839.1k 839.2k | 839.3k | 839.4k | 839.5k | 839.6k | 839.7k | 839.8k | 839.9k | 840k | 840.1k | 840.2k |
|------------------------------------|------------|---------------|--------|--------|--------|--------|--------|--------|------------|---------|--------|--------|
| Rice BAC Tiling Path | | | | | | | | | | | | |
| AP001389 26827 157519 | | | | | | | | | | | | |
| TIGR Rice Loci | | | | | | | | | | | | |
| | | | | | | | | | LOC_Os06g0 | 2440 | | |
| TTCP Rice Cope Medale | | | | | | | | | expressed | protein | | |
| TTAK KICE DENE NOUEIS | | | | | | | | | LOC_Os06g0 | 2440.1 | | |
| ienBank Annotations | | | | | | | | | | | | |
| Genesh Predictions | | | | | | | | | | | | |
| GeneMarkHMM Predictions | | | | | | | | | | | | |
| WINSCAN Predictions | | | | | | | | | | | | |
| GlinnerHMM Predictions | | | | | | | | | | | | |
| ssembled Sorghum bicolor sequences | : (xyplot) | | 400 | | | | | | | | | |
| | | | LTOO | | | | | | | | | |
| | | | -70 | | | | | | | | | -70 |
| ssembled Zea mays sequences (xyplo | it) | | 40 | | | | | | | | | 40 |
| | | _ | 100 | | | | | | | | | -100 |
| | | | -70 | | | | | | | | | -70 |
| rabidopsis thaliana alignments (xy | plot) | | -40 | | | | | | | | | 40 |
| | | ſ | 100 | | | | | | | | | 100 |
| | | | 70 | | | | | | | | | -70 |
| | | | 40 | | | | | | | | | 40 |
| ice EST Read Pairs | | | | | | | | | | | | |
| strand Pair: CI636560 CI406980 | | | | | | | | | | | | |
| strand Pair: CI634298 CI404700 | | | | | | | | | | | | |
| strand Pair: CI642939 CI414939 | | | | | | | | | | | | |
| strand Pair: CI623474 CI390303 | | | | | | | | | | | | |
| strand Pair: CI650679 CI424051 | | | | | | | | | | | | |
| strand Pair: CI645372 CI417970 | | | | | | | | | | | | |
| VM Predictions (New) | | | | | | | | | | | | |

Expression Data

| Data Type | Data Source |
|--------------|--|
| EST/FL-cDNA | PASA/Manual curation |
| Peptide | Koller et al., PNAS, 2002 (6,296 peptides/2,528 fgenesh models) |
| MPSS | Blake Meyers (http://mpss.udel.edu/rice/) |
| SAGE | 126,663 tags from MGOS (http://www.mgosdb.org/sage/) |
| Microarray | NSF Rice Oligonucleotide Array project |
| | (http://www.ricearrary.org) |
| Tiling array | Deng lab, Yale University |

Expression Data in Gbrowse

| ← ← ← ← ← ← ← ← ← ← ← ← ← ← ← ← ← ← ← | 9.1k 839.2k 839.3k 83 | 39.4k 839.5k 839.6k 839.7 | ************************************** | 840.1k 840.2k 840.3k 840.4k 840.5k |
|--|-----------------------|---------------------------|--|------------------------------------|
| TIGR Rice Loci | | | LDC 0s06g02440 | |
| | | | expressed protein | |
| TIGR Rice Gene Models | | | LDC 0s06s02440 1 | |
| THINSCAN Predictions | | | | |
| Protein Fuidence | | | | |
| | | | | |
| Tale Tilling Hrray Frotile (65669967 Forward | | | | |
| | | | | |
| | | | | |
| | | | | |
| Yale Tiling Array Profile (GSE6996) Reverse | | | | |
| ······ | | | | |
| | | | | |
| | | | | |
| | | | _ | |
| MPSS Tags | | | | |
| | | G | | |
| | | | | |
| SAGE Tags | | | | CATGTAAAAACCTTCAGAATT |
| | | | | CATGAACCGGGCAATGTTG |
| | | | | CATGTAAAATCGAATAT |
| NSF 20k Rice Oligo Microarray | | | | |
| NSF 45k Rice Oligo Microarray | | | | |
| Affymetrix GeneChip Rice Genome Array | | | | |
| | | °ee | | |
| | | | | |
| | | | | |
| Yale Rice Oligo Microarray | | | | — |
| Agilent Rice Oligo Microarray | | | | |
| EVM Predictions (New) | | | | |
| | | | | |

MPSS SEQUENCING TECHNOLOGY

I. Library construction

Brenner et al., PNAS 97:1665-70.



Each bead contains the amplified product derived from the 3' end of a single transcript.



uncover next 4 bases, repeat cycle



©98 Keep

Ovary and mature stigma

| FME | (s) | 0 | 0 | 744,319 | 0 | 1,482,167 | 1.000 | 0 | 2,226,486 | 0.0 |
|-------------|------|----------|--------|---------|-----|-----------|-------|----|-----------|------|
| FRO | (s) | 0 | 0 | 702,046 | | 1,200,706 | 1.000 | 0 | 1,902,752 | 0.0 |
| FRR | (s) | 0 | 0 | 690,560 | | 1,017,481 | 1.000 | 0 | 1,708,041 | 0.0 |
| MC00 | (?) | 0 | 0 | 473,119 | 0 | 639,651 | 1.000 | 0 | 1,112,770 | 0.0 |
| MC24 | (?) | 0 | 0 | 616,970 | 0 | 550,712 | 1.000 | 0 | 1,167,682 | 0.0 |
| MR03 | (?) | 0 | 0 | 549,558 | 0 | 617,519 | 1.000 | 0 | 1,167,077 | 0.0 |
| MR06 | (?) | 0 | 0 | 383,788 | 0 | 450,596 | 1.000 | 0 | 834,384 | 0.0 |
| MR12 | (?) | 0 | 0 | 411,516 | 0 | 583,952 | 1.000 | 0 | 995,468 | 0.0 |
| MR24 | (?) | 0 | 0 | 399,688 | | 663,844 | 1.000 | 0 | 1,063,532 | 0.0 |
| MR48 | (?) | 0 | 0 | 544,207 | | 592,538 | 1.000 | 0 | 1,136,745 | 0.0 |
| MS03 | Ø | 0 | 0 | 583,551 | | 706,204 | 1.000 | 0 | 1,289,755 | 0.0 |
| MS06 | - ČŚ | 0 | 0 | 336,582 | 0 | 489,874 | 1.000 | 0 | 826,456 | 0.0 |
| MS12 | - ČŚ | 0 | 0 | 360,120 | | 452,446 | 1.000 | 0 | 812,566 | 0.0 |
| M\$24 | - ČŚ | 0 | 0 | 366.078 | | 446,143 | 1.000 | 0 | 812.221 | 0.0 |
| MS48 | - ČŚ | 0 | 0 | 500,054 | 0 | 676,083 | 1.000 | 0 | 1.176.137 | 0.0 |
| MS96 | 10 D | 0 | 0 | 333,441 | 0 | 476.251 | 1.000 | 0 | 809.692 | 0.0 |
| NCA | िं | 0 | 0 | 665,983 | 0 | 857.387 | 1.000 | 0 | 1.523.370 | 0.0 |
| NCL | (s) | 0 | 0 | 726,299 | 0 | 993,805 | 1.000 | 0 | 1.720.104 | 0.0 |
| NCR | 6 | ň | Ŭ. | 819,719 | 0 | 1.005.156 | 1.000 | 0 | 1.824.875 | 0.0 |
| NDI | 6 | ň | ů Ú | 868,702 | 0 | 1.057.364 | 1.000 | 0 | 1,926,066 | 0.0 |
| NDR | (0) | ň | 0 | 636,111 | 0 | 944 735 | 1.000 | 0 | 1,580,846 | 0.0 |
| NCD | (0) | ő | 0 | 809.426 | 0 | 1.154.506 | 1.000 | Ő | 1,963,932 | 0.0 |
| NCS | (0) | ő | 0 | 552,552 | 0 | 718 782 | 1.000 | Ő | 1,271,334 | 0.0 |
| NTD | 6 | ŏ | 0 | 000,002 | 0 | 1 152 626 | 1.000 | 0 | 2 042 459 | 0.0 |
| NI 4 | 0 | 0 | 0 | 000,020 | | 0 | 1.000 | 0 | 2,042,407 | 0.0 |
| NLA | 60 | 0 | 0 | 222.521 | 0 | 442 190 | 1.000 | 0 | 774 721 | 0.0 |
| NLD | 6 | <u> </u> | 0 | 444 641 | | 552 004 | 1.000 | 0 | 1.001.545 | 0.0 |
| NLC | 6 | <u> </u> | 0 | 522 002 | | 405.054 | 1.000 | 0 | 950 044 | 0.0 |
| NLD | 6 | <u> </u> | 0 | 442 570 | | 500.000 | 1.000 | 0 | 944 407 | 0.0 |
| NME | 6 | 0 | 0 | 054 500 | - 0 | 1 122 027 | 1.000 | 0 | 2010/55 | 0.0 |
| NOS | 6 | 20 | 22 | 764 207 | - 0 | 1,100,007 | 1,000 | 22 | 2,010,433 | 11.4 |
| NDO | 6 | - 27 | 22 | 011 742 | - 0 | 1,100,070 | 1.000 | 22 | 1,922,000 | 0.0 |
| NPO | | 0 | 0 | 811,742 | | 1,024,436 | 1.000 | 0 | 1,836,178 | 0.0 |
| NRZ NDA | 8 | 0 | 0 | 005.650 | | 1.010.000 | 1.000 | 0 | 0 | 0.0 |
| NRA | | <u> </u> | 0 | 905,659 | | 1,019,990 | 1.000 | 0 | 1,923,649 | 0.0 |
| NRB | | <u> </u> | 0 | 940,221 | | 1,015,594 | 1.000 | 0 | 1,955,815 | 0.0 |
| NSL | | 0 | 0 | 768,410 | | 1,203,268 | 1.000 | 0 | 1,973,678 | 0.0 |
| NSK | (5) | <u> </u> | 0 | 598,714 | | 582,877 | 1.000 | 0 | 1,281,391 | 0.0 |
| NST NUT | (5) | <u> </u> | 0 | 692,117 | | 1.040.005 | 1.000 | 0 | 1,438,617 | 0.0 |
| NTL | (5) | <u> </u> | 0 | 525,418 | | 1,048,895 | 1.000 | 0 | 1,675,313 | 0.0 |
| NTR | (5) | 0 | 0 | 737,011 | | 593,122 | 1.000 | 0 | 1,330,133 | 0.0 |
| PLA | (5) | 0 | 0 | 424,016 | | 503,578 | 1.000 | 0 | 927,594 | 0.0 |
| PLU | (5) | 0 | 0 | 429,898 | | 550,824 | 1.000 | 0 | 980,722 | 0.0 |
| PLW | (5) | 0 | 0 | 363,097 | | 436,803 | 1.000 | 0 | 799,900 | 0.0 |
| PSC | (5) | 0 | 0 | 454,566 | | 591,489 | 1.000 | 0 | 1,046,055 | 0.0 |
| <u>PS1</u> | 5 | 0 | 0 | 435,252 | | 556,580 | 1.000 | 0 | 991,832 | 0.0 |
| PSL | (5) | 0 | 0 | 368,498 | 0 | 516,690 | 1.000 | 0 | 885,188 | 0.0 |
| PSN | (5) | 0 | 0 | 366,140 | 0 | 603,782 | 1.000 | 0 | 969,922 | 0.0 |
| PSY | (s) | 0 | 0 | 411,871 | 0 | 558,414 | 1.000 | 0 | 970,285 | 0.0 |
| XC00 | (5) | 0 | 0 | 448,936 | 0 | 453,985 | 1.000 | 0 | 902,921 | 0.0 |
| XC 06 | (5) | 0 | 0 | 511,865 | 0 | 539,475 | 1.000 | 0 | 1,051,340 | 0.0 |
| <u>XC24</u> | (5) | 0 | 0 | 437,295 | 0 | 476,855 | 1.000 | 0 | 914,150 | 0.0 |
| XR03 | (s) | 0 | 0 | 389,774 | 0 | 484,771 | 1.000 | 0 | 874,545 | 0.0 |
| XR06 | (s) | 0 | 0 | 467,130 | 0 | 495,044 | 1.000 | 0 | 962,174 | 0.0 |

Refine Gene Structure

| 1818.7k 1818.8k 1818.9k 1819k 1 | 1819.1k 1819.2k | 1819.3k 1819. | 4k 1819.5k | 1819.6k |
|---------------------------------|-----------------|-----------------------|------------|-------------------|
| TIGR Rice Loci | | | | |
| hypothetical protein | | | | |
| TIGR Rice Gene Models | | | | |
| FGenesh Predictions | | - | | |
| Twinscan Predictions | | | | |
| GeneMarkHMM Predictions | | J | | |
| GlimmerHMM Predictions | | 1 | | |
| Rice FL-cDNA | | | | |
| TIGR Rice Transcript Assemblies | | | | |
| Haize | | | | |
| AZM5_97785 1458 2300 | | | - | |
| Haize2 | | | | ſ |
| | | | | |
| | | | | |
| Sorghun | | | | |
| Sorghun2 | | | | |
| | | | | |
| | | | | |
| Arabidopsis | | | | |
| | | 2004 17/10/18 17/1154 | 4 ■ | |
| Arabidopsis2 | | 1 | | ſ |
| | | 1 | - | |
| | | | | |
| Plant TA ORF | | J | | |
| Hordeum_vulgare_ORF_30850 | | | | |
| Triticum_aestivum_ORF_43558 | | } | | |
| Zea_mays_ukr_45072 | | J | | |
| n Toroum_acsorvum_ukr_43335 | | | Triticum | aestivum ORE 4355 |
| | | | | |

"Have no fear of perfection you'll never reach it."

- Salvador Dalí