

1st transPLANT user training workshop

November 13th, 2012

INRA Versailles

Annotating *Triticeae* sequences: the TriAnnot pipeline



Leroy P, Guilhot N, Sakai H, Bernard A, Theil S, Choulet F, Reboux S, Viseux C, Amano N, Giacomoni F, Alaux M, Legeai F, Cerutti L,

Itoh T, Quesneville H, Feuillet C

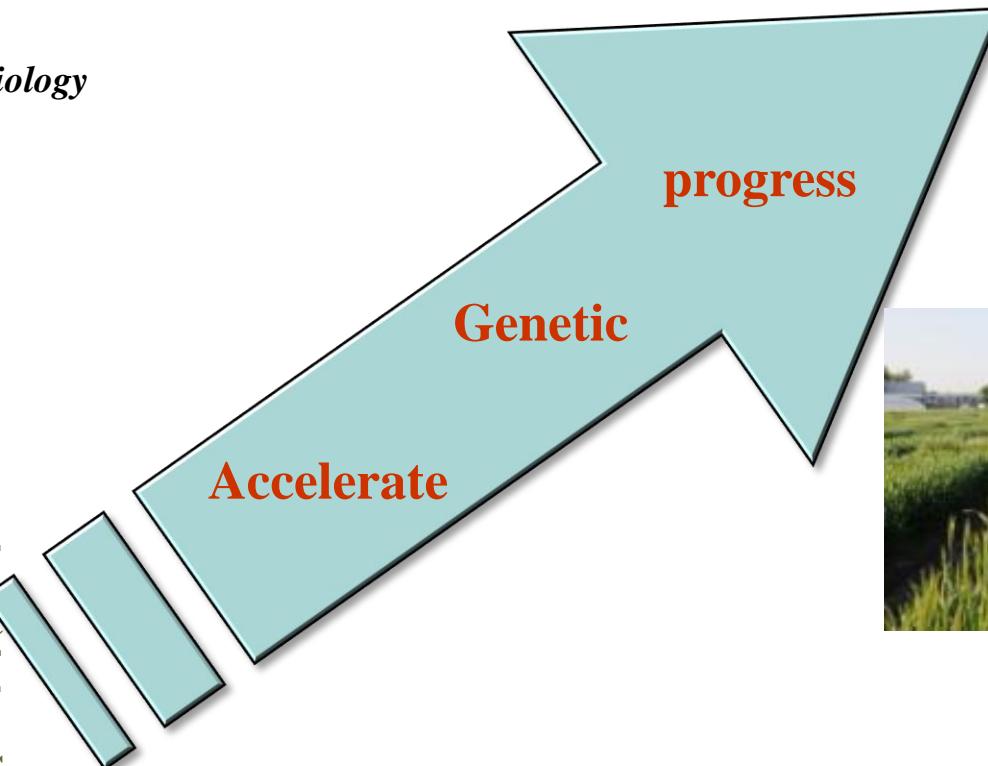
Plant Genomics

“Today genomics is the primary driver for unification in biological science”

Cook and Varshney (2010)

Current Opinion in Plant Biology
13:115-118.

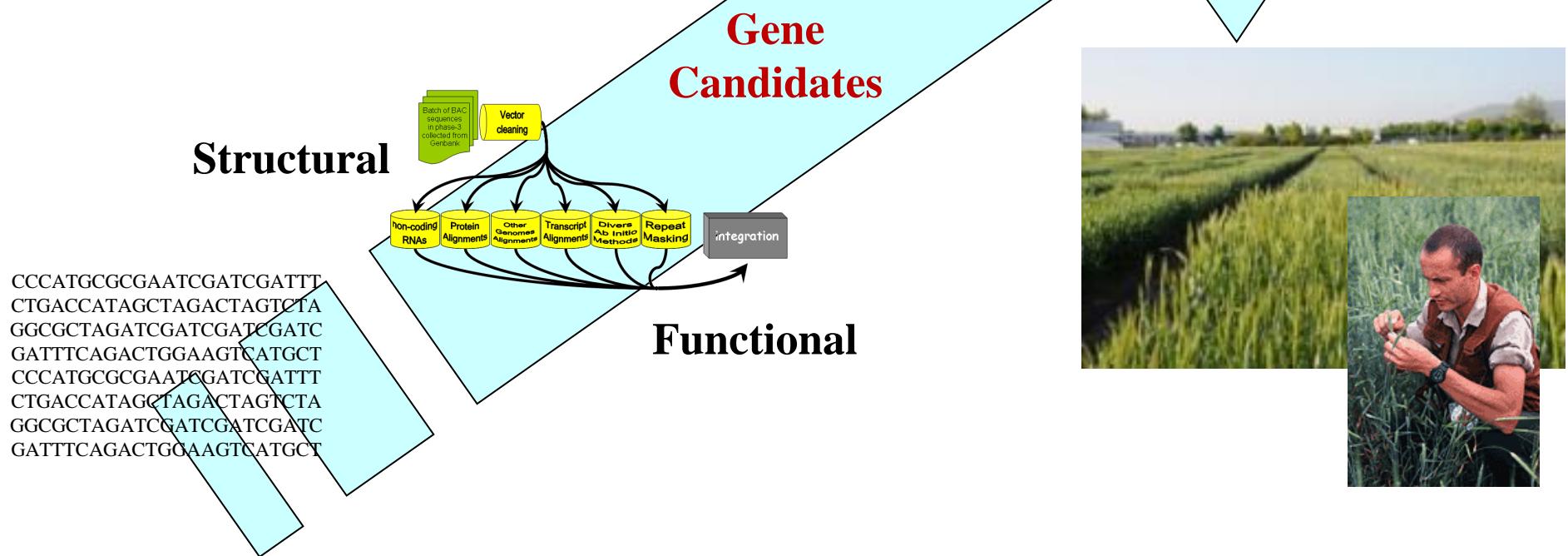
CCCATGCGCGAATCGATCGATT
CTGACCAGATAGCTAGACTAGTCTA
GGCGCTAGATCGATCGATCGATC
GATTTCACTGGAAAGTCATGCT
CCCATGCGCGAATCGATCGATT
CTGACCAGATAGCTAGACTAGTCTA
GGCGCTAGATCGATCGATCGATC
GATTTCACTGGAAAGTCATGCT



Genome Annotation

“Annotation is one of the most difficult tasks in genome sequencing projects, yet it is essential for connecting genome sequence to biology”

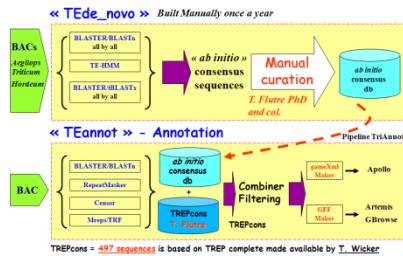
Elsik *et al.* 2006 Genome Research



Pipelines

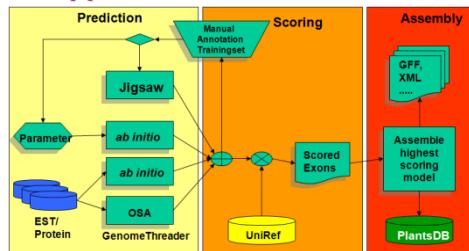
REPET

REPET Quesneville et al. (2005) PLoS Computational Biology

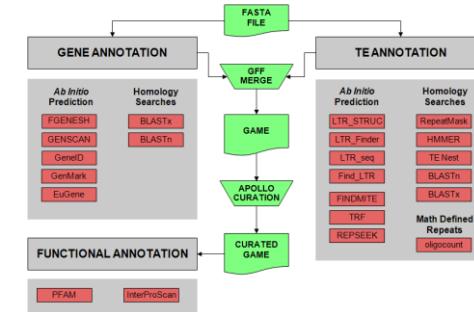


MIPS pipeline

MIPS pipeline

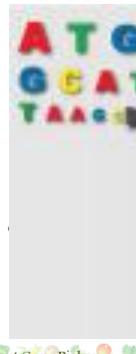


DAWG-PAWs



PASA

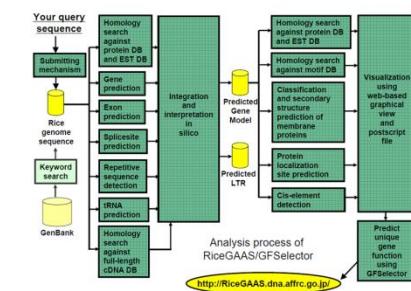
PASA and



Gnomon

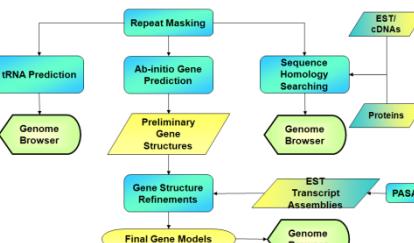


RiceGAAs



TIGR

Gene Finding: Flowchart (TIGR)





FranceAgriMer



Funding from the
European Community's
FP7 Programme - GA
FP7-212019



The International Wheat Genome Sequencing Consortium

Launched in 2005 on the initiative of Kansas Growers



www.wheatgenome.org

Executive director

K. Eversole

Cochairs

R. Appels
J. Dvorak
C. Feuillet
B. Gill
B. Keller
Y. Ogihara

Sponsors (23)



Coordinating Committee (64)

General members (>500)



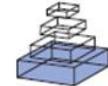
Leroy *et al.* (2012) Frontiers in Plant Science 3:1-14

frontiers in
PLANT SCIENCE

METHODS ARTICLE

published: 31 January 2012

doi: 10.3389/fpls.2012.00005



TriAnnot: a versatile and high performance pipeline for the automated annotation of plant genomes

Philippe Leroy^{1*}, Nicolas Guilhot¹, Hiroaki Sakai², Aurélien Bernard^{1,3}, Frédéric Choulet¹, Sébastien Theil¹, Sébastien Reboux⁴, Naoki Amano^{2,5}, Timothée Flutre⁴, Céline Pelegrin¹, Hajime Ohyanagi^{6,7}, Michael Seidel⁸, Franck Giacomoni⁹, Mathieu Reichstadt¹⁰, Michael Alaux⁴, Emmanuelle Gicquel¹, Fabrice Legeai¹¹, Lorenzo Cerutti¹², Hisataka Numa², Tsuyoshi Tanaka², Klaus Mayer⁸, Takeshi Itoh², Hadi Quesneville⁴ and Catherine Feuillet^{1*}

TriAnnot Architecture

Modular

N. Guilhot
GDEC



Panels

Transposable Elements

I

TEannot BLASTx
(TREPcons) (TREPprot)

Initial sequence

k-mer frequency
(cs 3B 2X)

Protein Coding Genes

II

ab initio
FgeneSH, Augustus
GeneMarkHMM, geneid

Similarity

Exonerate / Gmap
Plant FL-cDNA, mRNA
& ESTs

Exonerate
Plant proteomes

RepeatMasker
Masking

Gene Modeling
Transcripts
SIMsearch
ab initio & Transcripts
EuGene
ab initio
Augustus

Merge

Functional Annotation

known function
putative function
domain containing protein
expressed sequence
conserved unknown function
hypothetical protein

IWGSC
guideline

InterProscan (Pfam, Prosite, Smart, GO)

Best Hit (plant proteomes)

III

Non Coding Sequences

Conserved NCS

BLASTn

A. thaliana, O. sativa, Z. mays, S. bicolor,
B. distachyon, plastids & mitochondria,
EMBL plant wgs / tsa

BLASTx

SwissProt, TrEMBL, Uniref, nr,
NCBI RefSeq proteins, Plant
proteomes

Panel I & II
masking

ncRNAs
tRNAscan

Update
112 databases
21 programs

Wheat
Barley
Rice
Oak

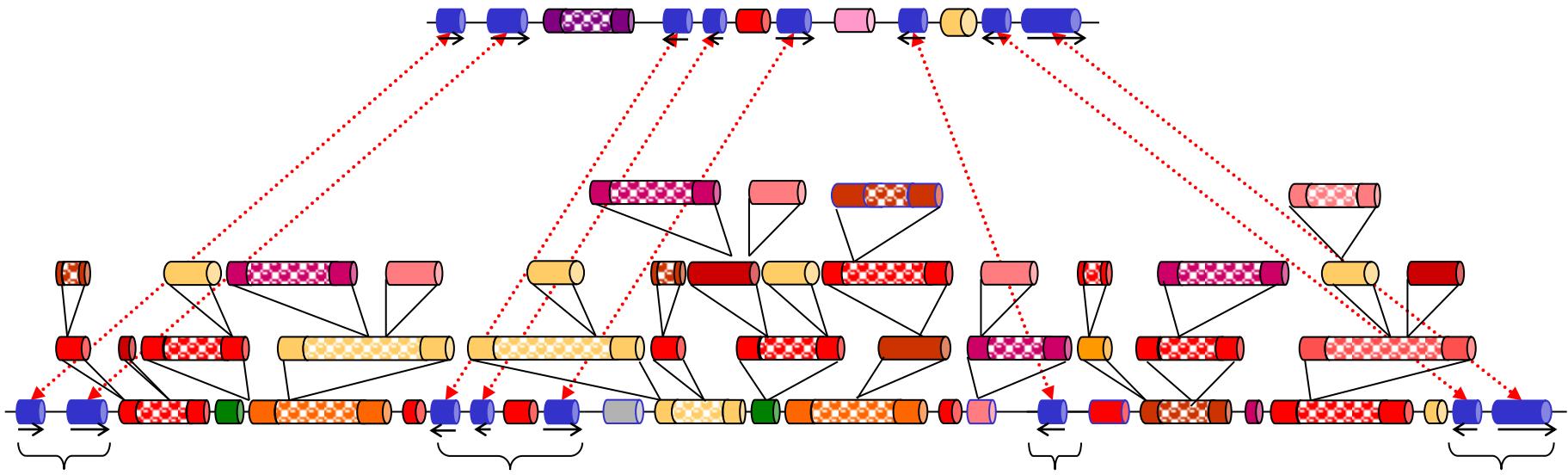
IV

Molecular Markers

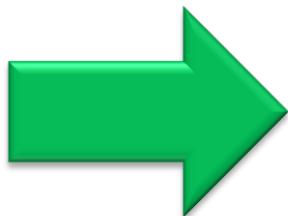
Microsatellites (SSRs)

Tandem Repeat Finder (TRF)

Transposable Elements



Erayman et al., 2004



To be masked

Transposable Elements “code”

Wicker et al. (2007)

A unified classification system for eukaryotic transposable elements.

Nature Reviews Genetics 8:973-982

Classification		Structure	TSD	Code	Occurrence
Order	Superfamily				
<i>Class I (retrotransposons)</i>					
LTR	Copia	→ GAG AP INT RT RH →	4–6	RLC	P,M,F,O
	Gypsy	→ GAG AP RT RH INT →	4–6	RLG	P,M,F,O
	Bel-Pao	→ GAG AP RT RH INT →	4–6	RLB	M
	Retrovirus	→ GAG AP RT RH INT ENV →	4–6	RLR	M
	ERV	→ GAG AP RT RH INT ENV →	4–6	RLE	M
DIRS	DIRS	→ GAG AP RT RH YR →	0	RYD	P,M,F,O
	Ngaro	→ GAG AP RT RH YR → →	0	RYN	M,F
	VIPER	→ GAG AP RT RH YR → →	0	RYV	O
PLE	Penelope	← → RT EN →	Variable	RPP	P,M,F,O
LINE	R2	— RT EN —	Variable	RIR	M
	RTE	— APE RT —	Variable	RIT	M
	Jockey	— ORF1 — APE RT —	Variable	RIJ	M
	L1	— ORF1 — APE RT —	Variable	RIL	P,M,F,O
	I	— ORF1 — APE RT RH —	Variable	RII	P,M,F
SINE	tRNA	—	Variable	RST	P,M,F
	7SL	—	Variable	RSL	P,M,F
	5S	—	Variable	RSS	M,O
<i>Class II (DNA transposons) - Subclass 1</i>					
TIR	Tc1-Mariner	→ Tase* ←	TA	DTT	P,M,F,O
	hAT	→ Tase* ←	8	DTA	P,M,F,O
	Mutator	→ Tase* ←	9–11	DTM	P,M,F,O
	Merlin	→ Tase* ←	8–9	DTE	M,O
	Transib	→ Tase* ←	5	DTR	M,F
	P	→ Tase ←	8	DTP	P,M
	PiggyBac	→ Tase ←	TTAA	DTB	M,O
	PIF-Harbinger	→ Tase* → ORF2 ←	3	DTH	P,M,F,O
	CACTA	← → Tase → ORF2 ← →	2–3	DTC	P,M,F
Crypton	Crypton	— YR —	0	DYC	F
<i>Class II (DNA transposons) - Subclass 2</i>					
Helitron	Helitron	— RPA — / — Y2 HEL —	0	DHH	P,M,F
Maverick	Maverick	→ C-INT → ATP → / — CYP — POL B —	6	DMM	M,F,O

TriAnnot Architecture

Panels

Transposable Elements

I

TEannot
(TREPcons) BLASTx
(TREPprot)

Initial sequence

k-mer frequency
(cs 3B 2X)

Protein Coding Genes

II

ab initio

FgeneSH, Augustus
GeneMarkHMM, geneid

Similarity

Exonerate / Gmap
Plant FL-cDNA, mRNA
& ESTs

Exonerate

Plant proteomes

RepeatMasker
Masking

Gene Modeling
Transcripts
SIMsearch
ab initio & Transcripts
EuGene
ab initio
Augustus

Merge

Functional Annotation

known function
putative function
domain containing protein
expressed sequence
conserved unknown function
hypothetical protein

IWGSC
guideline

InterProscan (Pfam, Prosite, Smart, GO)

Best Hit (plant proteomes)

Non Coding Sequences

Conserved NCS

III

BLASTn
A. thaliana, O. sativa, Z. mays, S. bicolor,
B. distachyon, plastids & mitochondria,
EMBL plant wgs / tsa

BLASTx

SwissProt, TrEMBL, Uniref, nr,
NCBI RefSeq proteins, Plant
proteomes

Panel I & II
masking

ncRNAs
tRNAscan

Update
112 databanks
21 programs

Wheat
Barley
Rice
Oak

Molecular Markers

IV

Microsatellites (SSRs)

Tandem Repeat Finder (TRF)

Protein-coding gene models

Similarity based

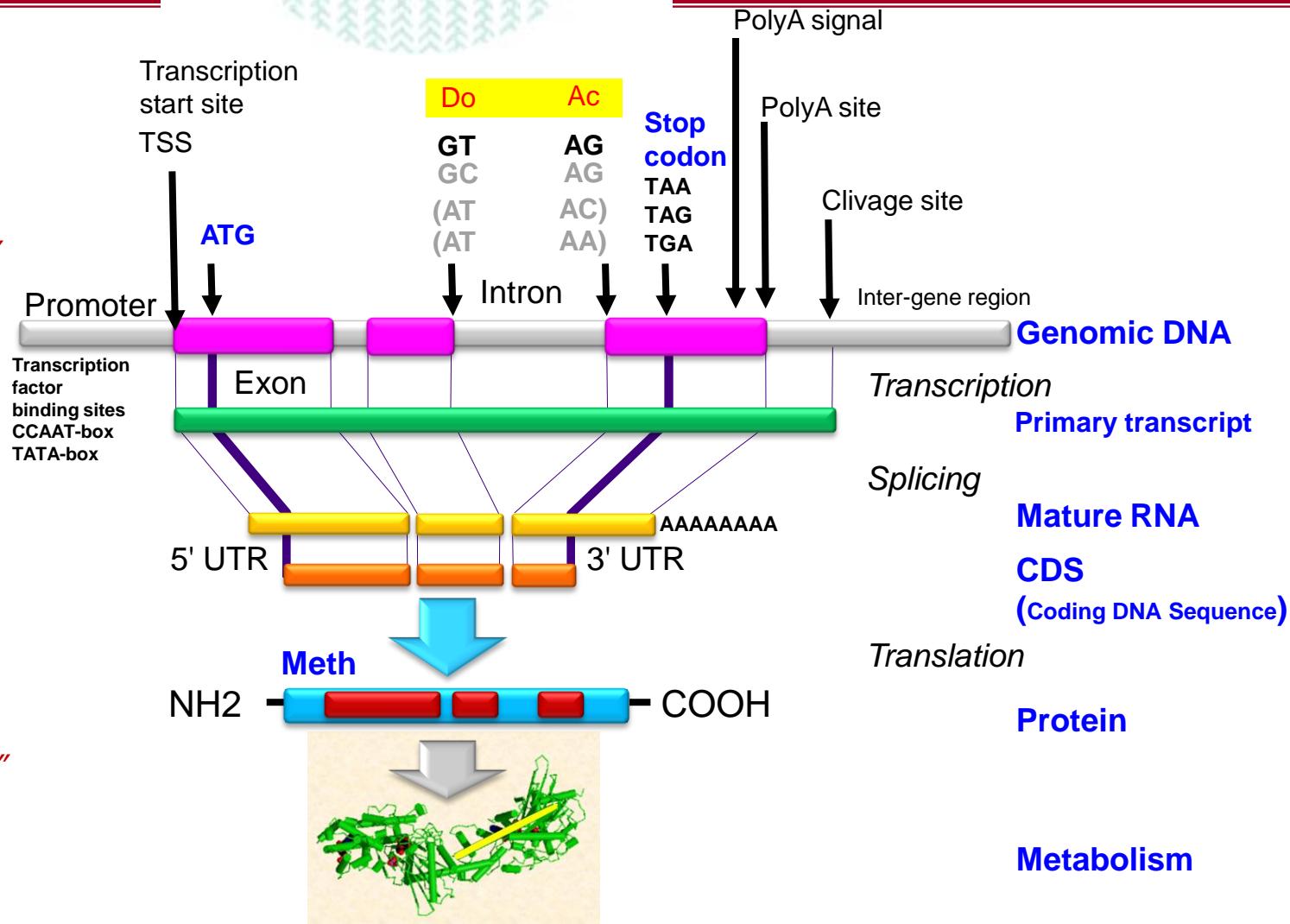
Biological evidences
"Extrinsic approach"

Combiner

Biological evidences &
training data set

ab initio

Training data set
"Intrinsic approach"



Use a color code system to assess the confidence of the structural annotation

Panels

Transposable Elements

I TEannot
(TREPcons) BLASTx
(TREPprot)

Initial sequence

RepeatMasker
Masking

k-mer frequency
(CS 3B 2x)

Protein Coding Genes

ab initio

FgeneSH, Augustus
GeneMarkHMM, geneid

Similarity

Exonerate / Gmap
Plant FL-cDNA, mRNA
& ESTs
Exonerate
Plant proteomes

Gene Modeling
Transcripts
SIMsearch
ab initio & Transcripts
EuGene
ab initio
Augustus

Merge

Functional Annotation

known function
putative function
domain containing protein
expressed sequence
conserved unknown function
hypothetical protein

IWGSC
guideline

InterProscan (Pfam, Prosite, Smart, GO)

Best Hit (plant proteomes)

II

III

IV

Non Coding Sequences

Conserved NCS

BLASTn

A. thaliana, O. sativa, Z. mays, S. bicolor,
B. distachyon, plastids & mitochondria,
EMBL plant wgs / tsa

BLASTx

SwissProt, TrEMBL, Uniref, nr,
NCBI RefSeq proteins, Plant
proteomes

ncRNAs
tRNAscan

Update

112 databanks
21 programs

Wheat
Barley
Rice
Oak

Molecular Markers

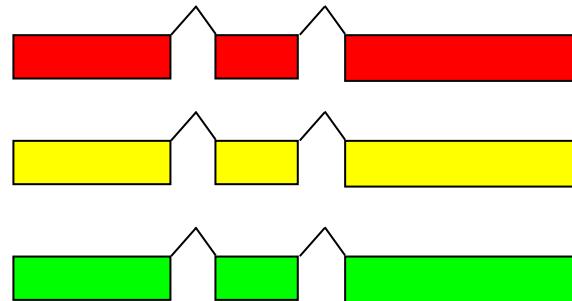
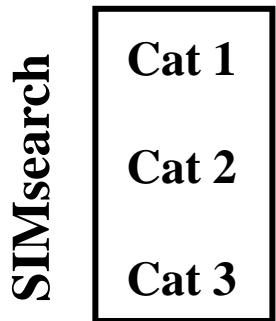
Microsatellites (SSRs)

Tandem Repeat Finder (TRF)

Structural Annotation

5 Categories

Similarity

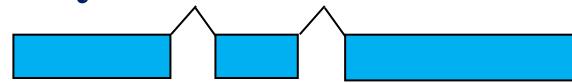


Wheat FL-cDNAs + wheat RNAseq

CDS-genes derived from reference genomes annotation (*Os-irgsp*, *Bd*)
Poaceae FL-cDNAs

ab initio & Similarity

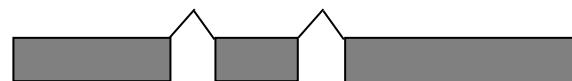
Cat 4



EuGene Combiner :
Augustus + wheat-ESTs +
SIMnuc + SIMprot

ab initio

Cat 5



Augustus
ab initio gene prediction only

Use a color code system to assess the confidence of the structural annotation



Panels

I

Transposable Elements

TEannot
(TREPcons) BLASTx
(TREPprot)

II

Protein Coding Genes

ab initio
FgeneSH, Augustus
GeneMarkHMM, geneid

Similarity

Exonerate / Gmap
Plant FL-cDNA, mRNA
& ESTs

Exonerate

Plant proteomes

III

Non Coding Sequences

Conserved NCS

BLASTn BLASTx

*A. thaliana, O. sativa, Z. mays, S. bicolor,
B. distachyon*, plastids & mitochondria,
EMBL plant wgs / tsa



k-mer frequency
(CS 3B 2X)

Functional Annotation

known function
putative function
domain containing protein
expressed sequence
conserved unknown function
hypothetical protein

IWGSC
guideline

Merge

InterProscan (Pfam, Prosite, Smart, GO)

Best Hit (plant proteomes)

IV

Molecular Markers

Microsatellites (SSRs)

Tandem Repeat Finder (TRF)

Panel I & II
masking

ncRNAs
tRNAscan

Update

112 databanks
21 programs

Wheat
Barley
Rice
Oak

Functional Annotation

6 Classes

known-function

**>80% identity
>80% coverage**

[UniProtKB/Swiss-Prot](#)

putative-function

**>45% identity
>50% coverage**

[UniProtKB/Swiss-Prot](#)

[UniProtKB/TrEMBL](#)

Not annotated as putative or hypothetical

domain-containing-protein

Pfam

expressed-sequence

**>45% identity
>50% coverage**

[ESTs databanks](#)

conserved-unknown-function

**>45% identity
>50% coverage**
annotated as putative or hypothetical

[UniProtKB/Swiss-Prot](#)

[UniProtKB/TrEMBL](#)

hypothetical-protein

no known function



TriAnnot Architecture

Panels

Transposable Elements

I

TEannot
(TREPcons) BLASTx
(TREPprot)

Initial sequence

RepeatMasker
Masking

k-mer frequency
(cs 3B 2X)

Protein Coding Genes

II

ab initio
FgeneSH, Augustus
GeneMarkHMM, geneid

Similarity

Exonerate / Gmap
Plant FL-cDNA, mRNA
& ESTs

Exonerate
Plant proteomes

Gene Modeling
Transcripts
SIMsearch
ab initio & Transcripts
EuGene
ab initio
Augustus

Merge

Functional Annotation

known function
putative function
domain containing protein
expressed sequence
conserved unknown function
hypothetical protein

IWGSC
guideline

InterProscan (Pfam, Prosite, Smart, GO)

Best Hit (plant proteomes)

III

Non Coding Sequences

IV

Conserved NCS
BLASTn BLASTx

A. thaliana, O. sativa, Z. mays, S. bicolor,
B. distachyon, plastids & mitochondria,
EMBL plant wgs / tsa

Panel I & II
masking

ncRNAs
tRNAscan

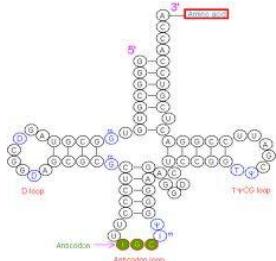
Update
112 databanks
21 programs
Wheat
Barley
Rice
Oak

Molecular Markers

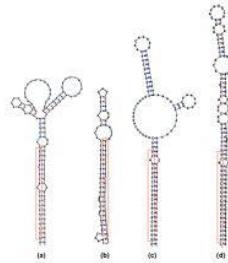
Microsatellites (SSRs)
Tandem Repeat Finder (TRF)

Non-coding RNAs

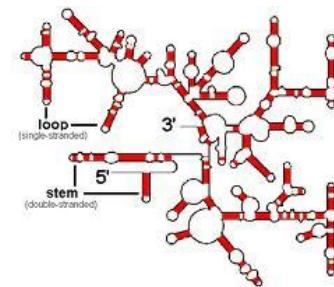
* tRNAs



*miRNA



* rRNA



Rfam

Rfam 11.0 (August 2012, 2208 families)

*snORNA



→ many other families
piRNA, siRNA, paRNA, uaRNA

TriAnnot Architecture

Panels

Transposable Elements

I

TEannot
(TREPcons) BLASTx
(TREPprot)

Initial sequence

k-mer frequency
(cs 3B 2X)

Protein Coding Genes

II

ab initio
FgeneSH, Augustus
GeneMarkHMM, geneid

Similarity
Exonerate / Gmap
Plant FL-cDNA, mRNA
& ESTs

Exonerate
Plant proteomes

RepeatMasker
Masking

Gene Modeling
Transcripts
SIMsearch
ab initio & Transcripts
EuGene
ab initio
Augustus

Merge

Functional Annotation

known function
putative function
domain containing protein
expressed sequence
conserved unknown function
hypothetical protein

IWGSC
guideline

InterProscan (Pfam, Prosite, Smart, GO)

Best Hit (plant proteomes)

Non Coding Sequences

III

Conserved NCs
BLASTn
A. thaliana, O. sativa, Z. mays, S. bicolor, B. distachyon, plastids & mitochondria, EMBL plant wgs / tsa

BLASTx

SwissProt, TrEMBL, Uniref, nr, NCBI RefSeq proteins, Plant proteomes

Panel I & II
masking

ncRNAs
tRNAscan

RNAspace

Update
112 databanks
21 programs

Wheat
Barley
Rice
Oak

Molecular Markers

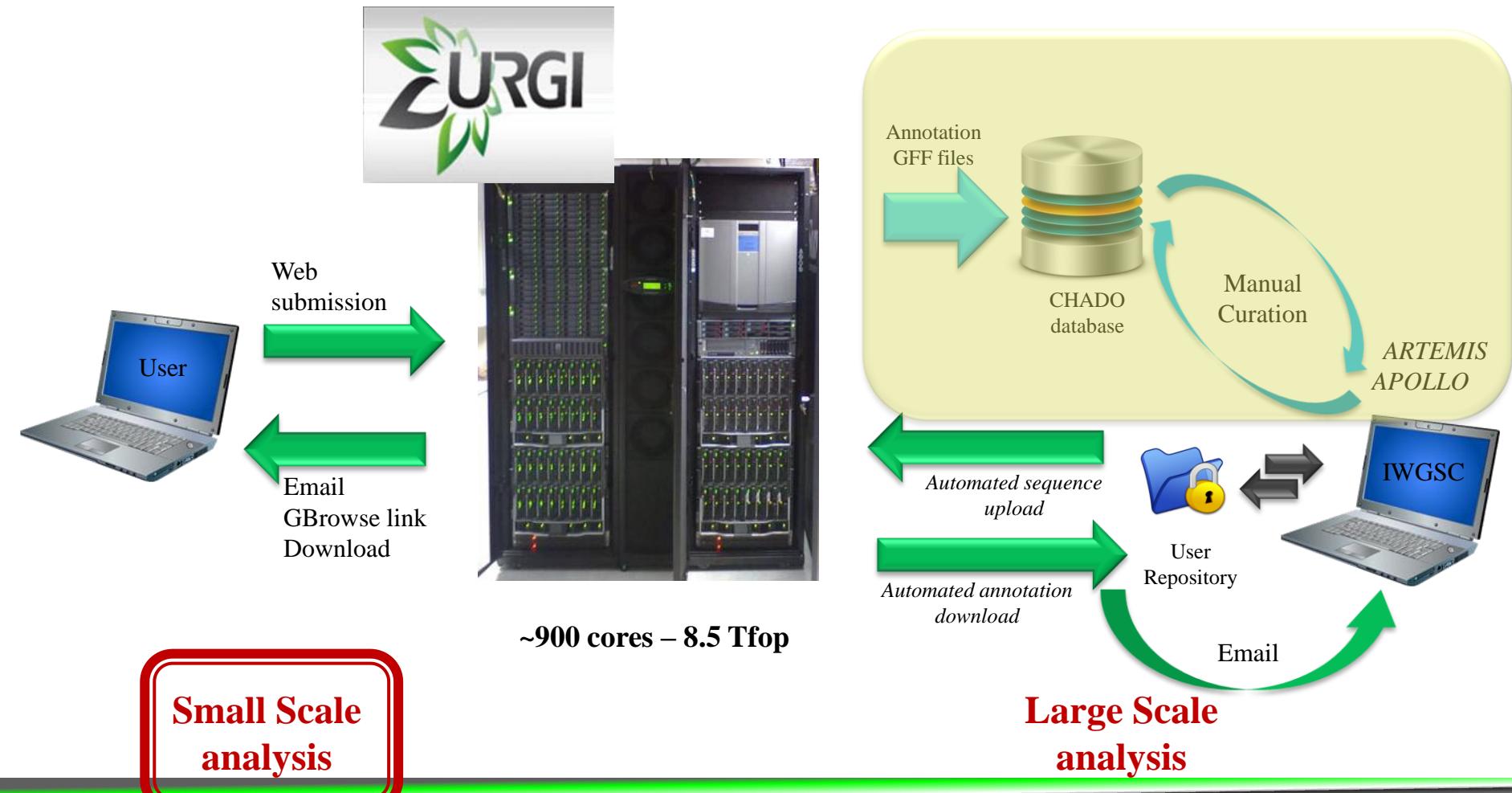
IV

Microsatellites (SSRs)
Tandem Repeat Finder (TRF)

ISBP

SNPs

How to use the TriAnnot pipeline ?



TriAnnot v3.6

Friendly web interface

Login/Password

- ~ 125 users registered since the beginning
- ~ 59 users at present
- ~ 33 actively using the pipeline
→ 726 analysis since January 2012

Submission

Automatic annotation

Graphical Viewing (Gbrowse)

Download data

EMBL, GFF, align

<http://www.clermont.inra.fr/triannot>

The screenshot shows the URGI PLANT AND FUNGI DATA INTEGRATION website. The top navigation bar includes links for FEEDBACK, CONTACT, SITE MAP, FUNDING, and a registration link. Below the navigation is a green header bar with links for About us, Projects, Research, Data, Tools, Species, and a SEARCH bar. A sidebar on the left is titled 'Tools' and contains links for GnPLS, Triannot Pipeline, Help, News, Architecture, Usage, Softwares, Databases, Defaults, and Links. A red box highlights the 'Run Pipeline' button in the 'Tools' sidebar. The main content area is titled 'Triannot Pipeline'. It features a sub-header 'You are here : Home / Tools / Triannot Pipeline', a large image of a wheat ear, and a detailed text block about wheat's history and genome. At the bottom of the content area, another red box highlights the text 'In the last 50 years, extensive genetic and cytogenetic (aneuploide and deletion lines) studies have led to the...'.



triannot

Environ 2 060 résultats (0,17 secondes)

[Triannot Pipeline - URGI](#) - [Traduire cette page]

19 Nov 2010 ... A first version of **TriAnnot** has been deposited on January 31th, 2006 to the French APP ("Agence pour la Protection des Programmes") ...
urgi.versailles.inra.fr/.../Triannot-Pipeline - En cache - Pages similaires

“TriAnnot Pipeline” TAB

The screenshot shows the "Analysis page" of the TriAnnot Pipeline. At the top, there is a navigation bar with tabs: "TriAnnot Pipeline", "My Analyses", and "New". Below the navigation bar, the title "TRIANNOT PIPELINE" is displayed. The main form area has a header "Analysis page". It contains two sections: "Analysis title" and "Pipeline template". The "Analysis title" section includes a red-bordered input field labeled "Give a descriptive title to your analysis (maximum 100 characters)". The "Pipeline template" section includes a dropdown menu labeled "Please choose a template...". Below these sections is a large input area for "Enter Query Sequence" with a placeholder "Enter FASTA sequence (> 1000 bp and < 3 Mbp)". There is also a "Browse..." button and a "Clear" button for file uploads. At the bottom of the form is a "Submit analysis" button.

The template is the “menu” : the **step.xml** which defines the type of programs to be used with which parameters and which databanks. A lot of combination can be written. A full step is available for wheat and a default analysis is usually used for wheat and other species.

Several possibilities for analysis

TriAnnot Pipeline My Analyses My profile Admin About TriAnnot

TRIANNOT PIPELINE ANALYSIS SUBMISSION

Analysis parameters

Analysis title 

Give a descriptive title to your analysis (maximum 255 characters)

Pipeline template 

Please choose a template... 

Please choose a template...
Wheat default IWGSC Annotation
Wheat full analysis
Rice default analysis
Barley default analysis
Oak default analysis

Query Sequence
(length < 3 Mbp) 

Or, upload file 

Browse... Clear

Submit analysis

Status



TriAnnot Pipeline My Analyses My profile Admin About TriAnnot

MY TRIANNOT PIPELINE ANALYSES

Pending (0) Running (0) Finished (4) Failed (2) All (6)

Search: Show 15 entries

	Status	Submitted	Progress	Title	Sequence	Sequence Length	Started	Finished	Pipeline Template	TriAnnot version
<input type="checkbox"/>		12-01-12 10:33	<div style="width: 0%;">0%</div>	Just a test on a short sequence	ANF4514	4367	12-01-12 10:35	12-01-12 10:45	Wheat IWGSC Annotation	3.5
<input type="checkbox"/>		12-01-12 10:33	<div style="width: 0%;">0%</div>	Another test	ANF4514	1257	12-01-12 10:35	12-01-12 10:45	Wheat IWGSC Annotation	3.5
<input type="checkbox"/>		12-01-03 11:18	<div style="width: 100%; background-color: #2e6b2e;">100%</div>	ctg0464b_00000001_00020000	ctg0464b_00000001_00020000	200000	12-01-03 11:20	12-01-03 12:02	Wheat IWGSC Annotation	3.5
<input type="checkbox"/>		11-12-31 12:14	<div style="width: 100%; background-color: #2e6b2e;">100%</div>	BACsynth12	Synth12	127860	11-12-31 12:15	11-12-31 13:33	Wheat IWGSC Annotation	3.5
<input type="checkbox"/>		11-12-20 12:18	<div style="width: 100%; background-color: #2e6b2e;">100%</div>	Scaffolds A20	A20	117015	11-12-20 12:20	11-12-20 13:09	Wheat IWGSC Annotation	3.5
<input type="checkbox"/>		11-12-07 11:51	<div style="width: 100%; background-color: #2e6b2e;">100%</div>	Test_3Mb	ctg954_3Mb	2999940	11-12-23 18:15	11-12-24 11:39	Wheat IWGSC Annotation	3.5

Showing 1 to 6 of 6 entries

With selected rows: [Download Results](#) [Delete](#)

- Analyses is waiting in the queue for a free slot
- Analyses is in progress
- Analyses has been processed successfully
- Analyses has failed

- View results with GBrowse
- Download results file

Email



De Moi <triannot-support@clermont.inra.fr>

Répondre

Répondre à tous

Transférer

Archiver

Indésirable

Supprimer

Sujet [TriAnnot] Analyses results notification: finished (1)

31/10/2012 18:56

Pour leroy@clermont.inra.fr

Autres actions

Dear Philippe Leroy,

Results for the following TriAnnot analysis are available for download:

739 - "case study transPLANT" submitted on "2012-10-31 17:46:51"
(5 genes predicted)

You can download results files from "My Analyses" section on TriAnnot pipeline web interface.

<http://urgi.versailles.inra.fr/triannot/?results>

GBrowse



phleroy Version 3.6

TriAnnot Pipeline My Analyses My profile Admin About TriAnnot

MY TRIANNOT PIPELINE ANALYSES

Pending (0) Running (0) Finished (18) Failed (0) All (18)

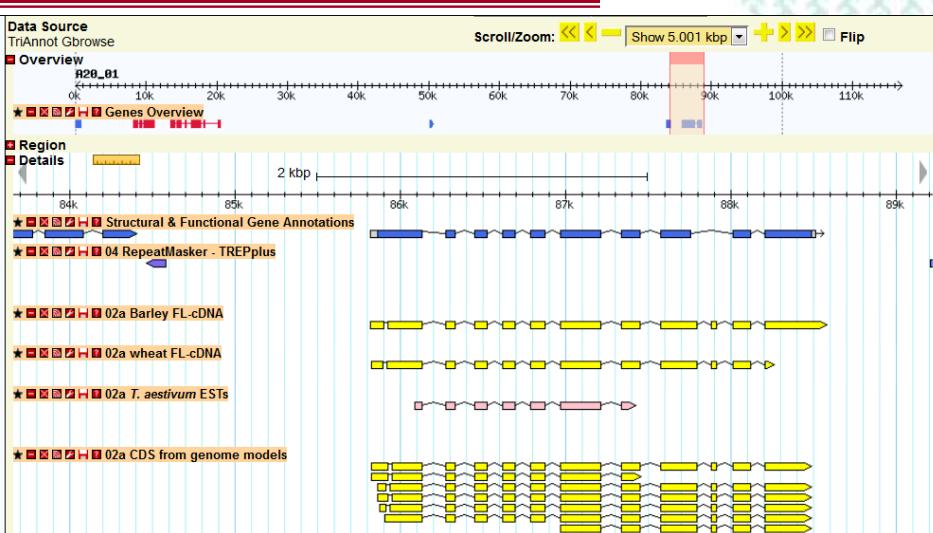
Search: Show 15 entries

<input type="checkbox"/> Status	Submitted	Progress	# predicted genes	Title	Sequence	Sequence Length	Started	Finished	Pipeline Template	TriAnnot version
<input type="checkbox"/>	12-10-24 14:41		102	Test Rice Default	Rice_test	998220	12-10-24 14:44	12-10-25 08:01	Rice default analysis	3.6
<input type="checkbox"/>	12-10-24 14:40		22	Test Oak Default	ASN_A_5E10_P2_137232bp	137232	12-10-24 14:44	12-10-24 22:16	Oak default analysis	3.6
<input type="checkbox"/>	12-10-24 14:40		8	Test Barley Default	seq01	191655	12-10-24 14:44	12-10-25 10:36	Barley default analysis	3.6
<input type="checkbox"/>	12-10-24 14:39		18	Test bacSynth12 Default	Synth12	127860	12-10-24 14:40	12-10-24 16:27	Wheat default IWGSC Annotation	3.6

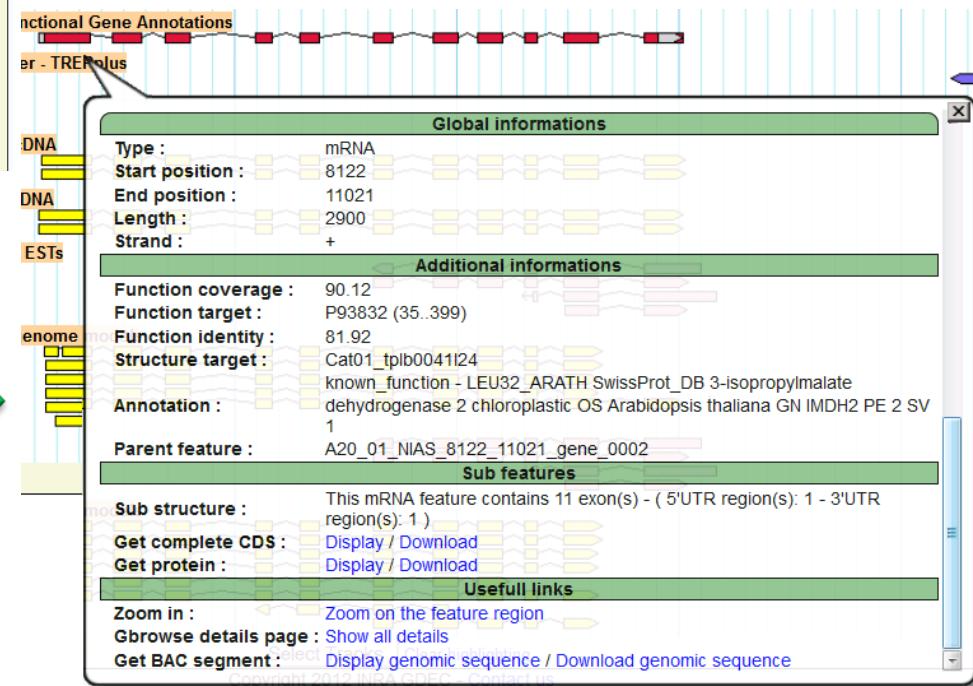
- Analyses is waiting in the queue for a free slot
- Analyses is in progress
- Analyses has been processed successfully
- Analyses has failed

View results with GBrowse
 Download results file

Gbrowse - v2.33



A quick graphical overview of gene models with a color code system to assess the confidence of structural annotation and biological evidences.



Pop up window to retrieve gene structural & functional annotation information, and download sequences (mRNA, CDS, protein).



Download data



phleroy Version 3.6

TriAnnot Pipeline My Analyses My profile Admin About TriAnnot

MY TRIANNOT PIPELINE ANALYSES

Pending (0) Running (0) Finished (18) Failed (0) All (18)

Search: Show 15 entries

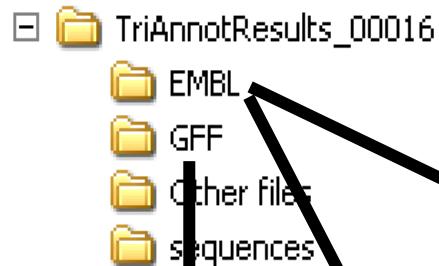
Status	Submitted	Progress	# predicted genes	Title	Sequence	Sequence Length	Started	Finished	Pipeline Template	TriAnnot version
	12-10-24 14:41		102	Test Rice Default	Rice_test	998220	12-10-24 14:44	12-10-25 08:01	Rice default analysis	3.6
	12-10-24 14:40		22	Test Oak Default	ASN_A_5E10_P2_137232bp	137232	12-10-24 14:44	12-10-24 22:16	Oak default analysis	3.6
	12-10-24 14:40		8	Test Barley Default	seq01	191655	12-10-24 14:44	12-10-25 10:36	Barley default analysis	3.6
	12-10-24 14:39		18	Test bacSynth12 Default	Synth12	127860	12-10-24 14:40	12-10-24 16:27	Wheat default IWGSC Annotation	3.6

- Analyses is waiting in the queue for a free slot
- Analyses is in progress
- Analyses has been processed successfully
- Analyses has failed

View results with GBrowse

Download results file

Manual curation



GenomeView (1838)

<http://genomeview.org/>

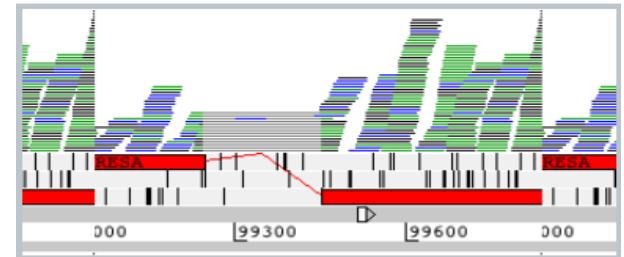
Abeel *et al.* (2011) Nucleic Acids Res. doi: 10.1093/nar/gkr995



Artemis (13.2.0)

<http://www.sanger.ac.uk/resources/software/artemis/>

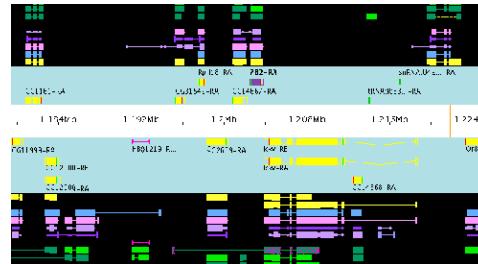
Carver *et al.* (2008) Bioinformatics 24, 2672-2676



Apollo (1.11.6)

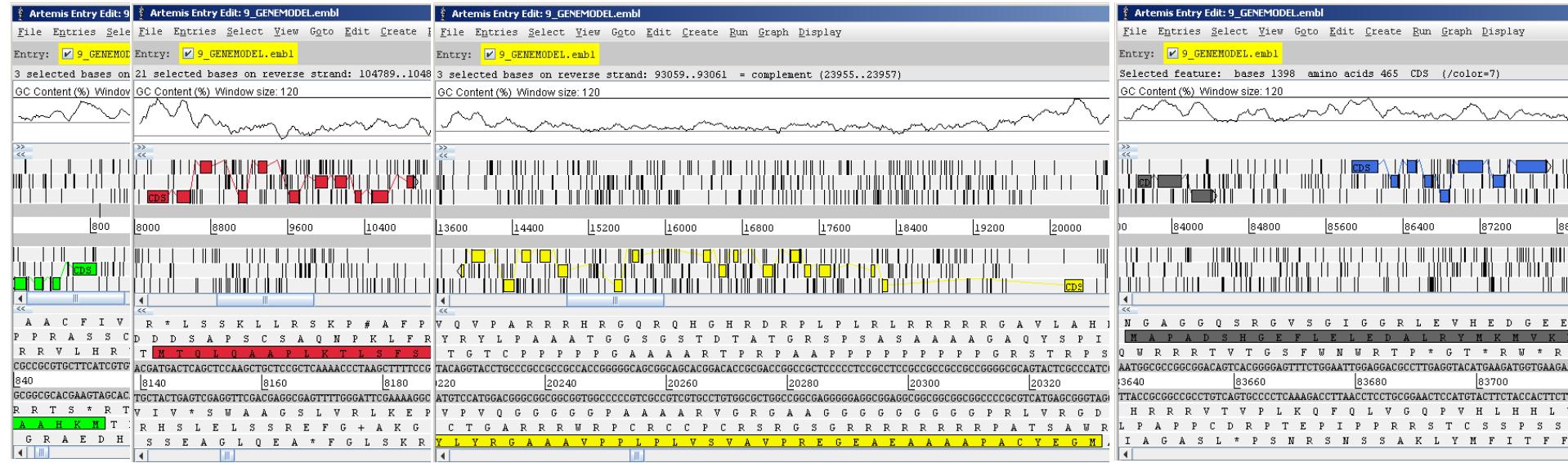
<http://apollo.berkeleybop.org/current/index.html>

Lewis *et al.* (2002) Genome Biology 3, research0082



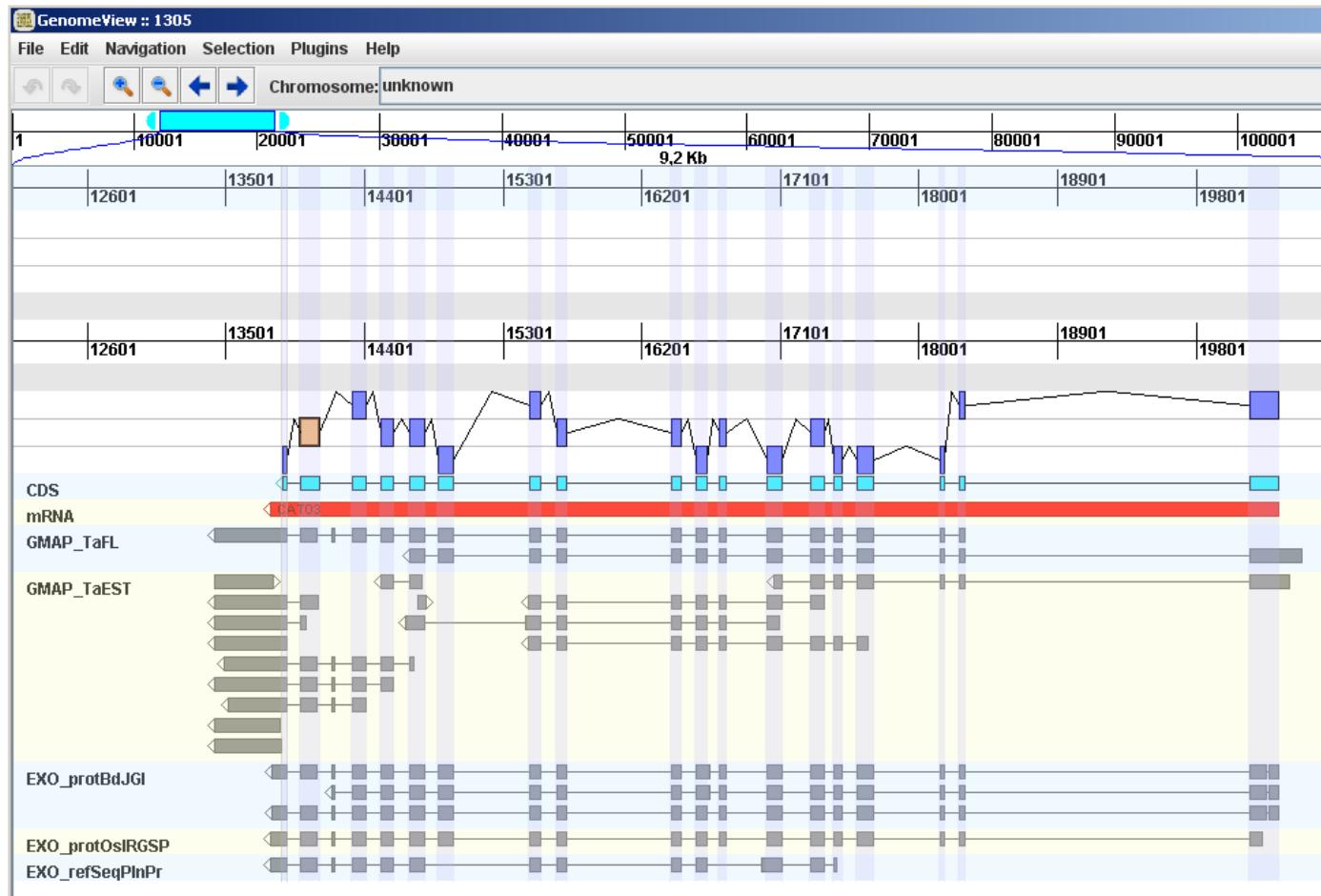
Artemis

Rutherford et al. (2000)
Bioinformatics 16:944-945



GenomeView

Thomas Abeel
Broad Institute





A case study A corse ასული

Sequences



Name	Size	Type
initial.seq	116 KB	SEQ File
proteins.seq	3 KB	SEQ File
TE_and_Genes_lower_case_masked.seq	116 KB	SEQ File
TE_and_Genes_N_masked.seq	116 KB	SEQ File
TE_lower_case_masked.seq	116 KB	SEQ File
TE_N_masked.seq	116 KB	SEQ File



>A20_SIMsearch_85867_87288_gene_0001_mRNA_0001
MKHAAALLLILAAAALVLLILAPAAHVLLPSAQYHWLNPEIHDAGGLSRQ
GFPEGFVFGTATSAVQVEGMAEQGGRGPSIWDAFIKIPGTIPGNGTANVA
VDEYHRYKEDVDIMKNMGFDAYRFSISWSRIFPDGTGVNWEGVGYYNRL
IDYMLQQGITPYANLYQYDLPLALHQYLGWLSPRIVGAFADYADFCFKV
FGDRVKNWFTNEPRCVAALGYDIGLHAPGRCSQCSAGGDSTTEPYLAAH
HLILSHAAAVERRYRDKYQVRLTEDYLQVISQFKNHIFNGSA*

>A20_SIMsearch_13417_20382_gene_0002_mRNA_0001
MGEYCAPAAAAAEGERPVAVSVLPLPPVAAAGRYLYGDYDRCSTKQVFD
NLHGNISLDPLAREFVDTEEFQQLRDLKQLGLTYLVFPGAHVTRFEHSLG
VYRLAGEAMNNLQKYQGNELGIDRIDVQTVKLAGLLHDIGHGPFSHLFEH
EFLPRVHPGSTWSHEHMSALLDSIVDKHAIDIEPDYLKVIKEMIVASSD
VSTAEGVKEKFLYDIVANGRNGIDVDKFDYIDRDCRACGIGSNFQHWRL
LEGMRVMGDEICYPAKDYLISHLKLTTRADLHRTVYTHAKVKAELMLVD
ALVEANEYLGISLHADDPEDFWKLDDTIVKSIEATPNDLKAKEIIQRI
RRRELYKFCNQYSVPDKLDFKNITAQDIVCSQITSKVLKKEEDVAVSN
VKIDLTRGKDNPLESSVVMPGFRMRREVPDNGRAREPPAARVQPGQDREG
VRQEAGAGGGVG*
>A20_SIMsearch_8122_11021_gene_0003_mRNA_0001
MTQLQAAPLTKLSFGRRSAVAAARPRHATATFRCSDAARSYIITLLPGD
GIGPEVVAVAKDVLSAAGAKEGVELRFREMLMGGALDAVGVP LPDET LA
AAQASDAILLGAIGGYKWDSNEKHLKPETGLLNIRAGLGVFANLRPATVL
PQLVDASTLKEVAEGVDIMVVREL TGGIYFGQPRGF GTNDMGE EIGFNT
EIYSVSEIDRIARVAFEVAR KRGKGKLCV DKANVLEASMLWRKRVTAIAS
EFPDVELSHMYVDNASMQLVRNPQKFDTIVTNNIFGDILSDEASMITGSI
GMLPSASVGESGPGLFEP IHSAPDIAGQDKANPLATILSAAMLLKYGLG
AETAAKRIETAVTETLDNGFRTGDIYSPGTTLVGCKRMGE EVLKALE SQK
*
>A20_EUGENE_16_854_gene_0001_mRNA_0001
MKHAAALLLILAAAHVLLPWAQCHRLNPEIHDAGGLSRQGFPEGFVFGT
AASAYQVEGMAEQGGRGPSIWDAFIKIPGTIAGNGTADVAVDEYHRYKVC
IEDVDIMKNMGFDAYRFSISWSRIFPDGTGVNWEGV DYYNR LIDYMLQQ
GNWNFYTRQEA*

>A20_EUGENE_50225_50590_gene_0004_mRNA_0001
METQAPPASLDLSLALATMPQLPPPAAAAPPLSLQAAGDAVSSAVAGAG
WKVF SCLFYEKKF LKS QALGGHQN AHRKDRGAAGWN ASLYLPAADRPWPP
TTATSHPEIGDENQLDSLKL*

Rules of GFF & EMBL files naming

<http://urgi.versailles.inra.fr/Species/Wheat/Triannot-Pipeline/Usage>

- EMBL folder
 - In each folder files are tagged and follows the following rules:
 - A number related to the step number
 - The type of programs *i.e.* REPEATMASKER; AUGUSTUS; BLASTX; BLASTN;BLASTP; EXONERATE; EUGENE; GENEMODEL; BESTHIT; TRNASCAN-SE; TRF.
 - The databank used (see [databanks](#)). When no databank is used, just the type of program is displayed *i.e.* TRF. For *ab initio* gene prediction programs the matrix used is displayed *i.e.* AUGUSTUS_wheat.
 - Extension .embl
 - Few examples
 - Step1 - *Transposable Elements annotation & masking*
 - 1_REPEATMASKER_TREP_plus.embl
 - Step2 - *BLASTx against TREPprot*
 - 2_BLASTX_TREP_prot.embl
 - Step5 - *ab initio gene prediction*
 - 5_AUGUSTUS_wheat.embl
 - Step6 - *BLASTn / Exonerate*
 - 6_BLASTN_cdsBdJGI.embl
 - 6_BLASTN_cdsOsIRGSP.embl
 - etc.

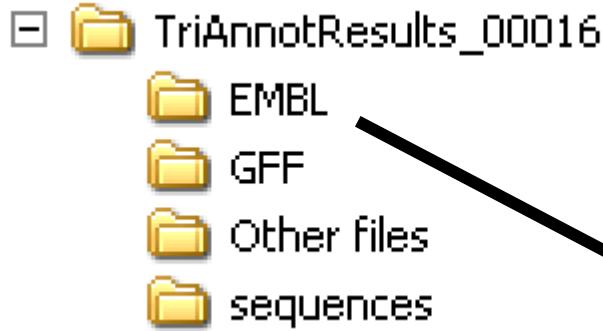
2_BLASTX_TREP_prot.embl

Protein alignments



- 10_BESTHIT_BLASTP_nr.align
- 10_BESTHIT_BLASTP_protBdJGI.align
- 10_BESTHIT_BLASTP_protHOR.align
- 10_BESTHIT_BLASTP_protOsIRGSP.align
- 10_BESTHIT_BLASTP_protPlant.align
- 10_BESTHIT_BLASTP_protZmMSO.align

EMBL files



- 1_REPEATMASKER_TREP_plus.embl
- 2_BLASTX_TREP_prot.embl
- 5_AUGUSTUS_wheat.embl
- 6_BLASTN_cdsBdJGI.embl
- 6_BLASTN_cdsOsIRGSP.embl
- 6_BLASTN_cdsOsMSU.embl
- 6_BLASTN_cdsSbDOE.embl
- 6_BLASTN_cdsZmMSO.embl
- 6_BLASTN_HvFL.embl
- 6_BLASTN_SIMnuc_wheat.embl
- 6_BLASTN_TaEST.embl
- 6_BLASTN_TaFL.embl

- 1_REPEATMASKER_TREP_plus.embl
- 2_BLASTX_TREP_prot.embl
- 5_AUGUSTUS_wheat.embl 
- 6_BLASTN_cdsBdJGI.embl
- 6_BLASTN_cdsOsIRGSP.embl
- 6_BLASTN_cdsOsMSU.embl
- 6_BLASTN_cdsSbDOE.embl
- 6_BLASTN_cdsZmMSO.embl
- 6_BLASTN_HvFL.embl
- 6_BLASTN_SIMnuc_wheat.embl
- 6_BLASTN_TaEST.embl
- 6_BLASTN_TaFL.embl
- 6_EXONERATE_cdsBdJGI.embl
- 6_EXONERATE_cdsOsIRGSP.embl
- 6_EXONERATE_cdsOsMSU.embl
- 6_EXONERATE_cdsSbDOE.embl
- 6_EXONERATE_cdsZmMSO.embl
- 6_EXONERATE_HvFL.embl
- 6_EXONERATE_SIMnuc_wheat.embl
- 6_EXONERATE_TaEST.embl
- 6_EXONERATE_TaFL.embl
- 7_BLASTX_protBdJGI.embl
- 7_BLASTX_protHOR.embl
- 7_BLASTX_protOsIRGSP.embl
- 7_BLASTX_protOsMSU.embl
- 7_BLASTX_protSbDOE.embl
- 7_BLASTX_protTRI.embl
- 7_BLASTX_protZmMSO.embl
- 7_BLASTX_refSeqPInProt.embl
- 7_BLASTX_SIMprot_wheat.embl
- 7_BLASTX_uniprot_sprot.embl
- 7_EXONERATE_protBdJGI.embl
- 7_EXONERATE_protHOR.embl
- 7_EXONERATE_protOsIRGSP.embl
- 7_EXONERATE_protOsMSU.embl
- 7_EXONERATE_protSbDOE.embl
- 7_EXONERATE_protTRI.embl
- 7_EXONERATE_refSeqPInProt.embl
- 7_EXONERATE_SIMprot_wheat.embl
- 7_EXONERATE_uniprot_sprot.embl
- 8_EUGENE.embl 
- 9_GENEMODEL.embl 
- 10_BESTHIT_BLASTP_nr.embl
- 10_BESTHIT_BLASTP_protAtTAIR.embl
- 10_BESTHIT_BLASTP_protBdJGI.embl
- 10_BESTHIT_BLASTP_protHOR.embl
- 10_BESTHIT_BLASTP_protOsIRGSP.embl
- 10_BESTHIT_BLASTP_protOsMSU.embl
- 10_BESTHIT_BLASTP_protPlant.embl
- 10_BESTHIT_BLASTP_protSAC.embl
- 10_BESTHIT_BLASTP_protSbDOE.embl
- 10_BESTHIT_BLASTP_protTRI.embl
- 10_BESTHIT_BLASTP_protZmMSO.embl
- 10_EXONERATE_nr.embl
- 10_EXONERATE_protAtTAIR.embl
- 10_EXONERATE_protBdJGI.embl
- 10_EXONERATE_protHOR.embl
- 10_EXONERATE_protOsIRGSP.embl
- 10_EXONERATE_protOsMSU.embl
- 10_EXONERATE_protPlant.embl
- 10_EXONERATE_protSAC.embl
- 10_EXONERATE_protSbDOE.embl
- 10_EXONERATE_protTRI.embl
- 10_EXONERATE_protZmMSO.embl
- 11_INTERPROSCAN.embl 
- 12_BLASTN_genoAtTAIR.embl
- 12_BLASTN_genoBdJGI.embl
- 12_BLASTN_genoOsIRGSP.embl
- 12_BLASTN_genoOsMSU.embl
- 12_BLASTN_genoSbDOE.embl
- 12_BLASTN_genoZmMSO.embl
- 12_BLASTN_refSeq_chloro.embl
- 12_BLASTN_refSeq_mito.embl
- 13_BLASTX_protAtTAIR.embl
- 13_BLASTX_protBdJGI.embl
- 13_BLASTX_protOsIRGSP.embl
- 13_BLASTX_protOsMSU.embl
- 13_BLASTX_protSbDOE.embl
- 13_BLASTX_protZmMSO.embl
- 14_TRNASCAN-SE.embl
- 15_TRF.embl

Gene Annotation

Avoid to open empty files with GenomeView!

9_GENEMODEL.embl

ID	unknown; SV 1; linear; unassigned DNA; STD; UNC; 117015 BP.
XX	
AC	unknown;
XX	
XX	
FH	Key
FT	gene
FT	mRNA
FT	exon
FT	CDS
FT	gene

Location/Qualifiers

```
complement(13417..20382)
/locus_tag="A20_CAT01_1"
/color="223 38 56"
/label=CAT01
/id="A20_NIAS_13417_20382_gene_0001"
complement(13417..20382)
/locus_tag="A20_CAT01_2"
/color="223 38 56"
/label=CAT01
/id="A20_NIAS_13417_20382_gene_0001_mRNA_0001"
complement(13417..13642)
/locus_tag="A20_CAT01_3"
/color="223 38 56"
/label=CAT01
/id="A20_NIAS_13417_13642_gene_0001_mRNA_0001_exon_0001"
/note="Ontology_term: SO:0000202"
complement(13660..13894)
/locus_tag="A20_CAT01_4"
/color="223 38 56"
/label=CAT01
/id="A20_NIAS_13660_13894_gene_0001_mRNA_0001_exon_0002"
/note="Ontology_term: SO:0000004"
complement(13984..14111)
complement(20152..20382)
/locus_tag="A20_CAT01_21"
/color="223 38 56"
/label=CAT01
/id="A20_NIAS_20152_20382_gene_0001_mRNA_0001_exon_0019"
/note="Ontology_term: SO:0000200"
complement(join(13869..13894, 13984..14111, 14318..14409,
14507..14587, 14693..14794, 14883..14975, 15469..15542,
15652..15707, 16392..16454, 16555..16622, 16706..16744,
17016..17106, 17298..17382, 17450..17501, 17598..17704,
18140..18162, 18260..18297, 20152..20331))
/locus_tag="A20_CAT01_22"
/color="223 38 56"
/label=CAT01
/id="A20_NIAS_13417_20382_gene_0001_mRNA_0001_joinedCDS"
/note="Similar_to: HMMScan - domain-containing_protein -->
PF01966.15 - HD domain (1.6e-13)"
/note="Structure_target: Cat01_tplb0006o12"
8122..11021
```

Functional annotation





V3.6

This is not a tutorial about how
to annotate a sequence 😊

GenomeView (1838)

<http://genomeview.org/>



Abeel *et al.* (2011) Nucleic Acids Res. doi: 10.1093/nar/gkr995



Check list 1

- ✓ Build a curation repertory without empty files / BLAST files / SIMsearch files
- ✓ File / Load data... / local file/**9_GENEMODEL.embl**
 - Zoom and move (4 ways - the best = up/down/right/left keys) => **2nd gene**
 - Track list



With GenomeView the user needs the correspondence between the embl file name and the feature name of the track keys (wheat full analysis- TriAnnot 3.6).

1_REPEATMASKER_TREP_plus → repeat_region
1_REPEATMASKER_TREP_nr → repeat_region
1_REPEATMASKER_TREP_total → repeat_region
1_REPEATMASKER_ALL_Repbase → repeat_region
1_REPEATMASKER_MIPS_repeat → repeat_region
1_REPEATMASKER_TIGR_Fam_Repeats → repeat_region
1_REPEATMASKER_TIGR_GSS_Repeats → repeat_region

1_REPEATMASKER_univec → repeat_region
1_REPEATMASKER_Ecoli → repeat_region

2_BLASTX_TREP_prot → BLASTX_TREP_pro

5_AUGUSTUS_wheat → AUGUSTUS_wheat
5_FGENESH → FGeneSH
5_GENEID → GeneID
5_GMHMM_wheat → GMHMM_wheat

8_EUGENE → Eugene

9_GENEMODEL → gene / mRNA / exon / CDS



6_EXONERATE_AtFL → EXO_N_AtFL
6_EXONERATE_HvFL → EXO_N_HvFL
6_EXONERATE_OsFL → EXO_N_OsFL
6_EXONERATE_PlFL → EXO_N_PlFL
6_EXONERATE_PoFL → EXO_N_PoFL
6_EXONERATE_PpFL → EXO_N_PpFL
6_EXONERATE_PtFL → EXO_N_PtFL
6_EXONERATE_RosiFL → EXO_N_RosiFL
6_EXONERATE_TaFL → EXO_N_TaFL
6_EXONERATE_ZmFL → EXO_N_ZmFL

6_EXONERATE_HvEST → EXO_N_HvEST
6_EXONERATE_OsEST → EXO_N_OsEST
6_EXONERATE_QuerEST → EXO_N_QuerEST
6_EXONERATE_SoEST → EXO_N_SoEST
6_EXONERATE_TaEST → EXO_N_TaEST
6_EXONERATE_ZmEST → EXO_N_ZmEST

6_EXONERATE_rnaSeq_wheat → EXO_N_rnaSeq_wh

6_EXONERATE_AT_unigene → EXO_N_AT_unigen
6_EXONERATE_BD_unigene → EXO_N_BD_unigen
6_EXONERATE_HV_unigene → EXO_N_HV_unigen
6_EXONERATE_OS_unigene → EXO_N_OS_unigen
6_EXONERATE_PP_unigene → EXO_N_PP_unigen
6_EXONERATE_PT_unigene → EXO_N_PT_unigen
6_EXONERATE_QR_unigene → EXO_N_QR_unigen
6_EXONERATE_SB_unigene → EXO_N_SB_unigen
6_EXONERATE_SO_unigene → EXO_N_SO_unigen
6_EXONERATE_TA_unigene → EXO_N_TA_unigen
6_EXONERATE_VV_unigene → EXO_N_VV_unigen
6_EXONERATE_ZM_unigene → EXO_N_ZM_unigen

6_EXONERATE_cdsAtTAIR → EXO_N_cdsAtTAIR
6_EXONERATE_cdsBdJGI → EXO_N_cdsBdJGI
6_EXONERATE_cdsOsIRGSP → EXO_N_cdsOsIRGSP
6_EXONERATE_cdsOsMSU → EXO_N_cdsOsMSU
6_EXONERATE_cdsPp → EXO_N_cdsPp
6_EXONERATE_cdsPt → EXO_N_cdsPt
6_EXONERATE_cdsSbDOE → EXO_N_cdsSbDOE
6_EXONERATE_cdsVv → EXO_N_cdsVv
6_EXONERATE_cdsZmMSO → EXO_N_cdsZmMSO

6_EXONERATE_PlCpltCDS → EXO_N_PlCpltCD
6_EXONERATE_PlMrnaSTD → EXO_N_PlMrnaST
6_EXONERATE_PoaMrnaSTD → EXO_N_PoaMrnaST



7_EXONERATE_pepATH → EXO_X_pepATH
 7_EXONERATE_pepBDI → EXO_X_pepBDI
 7_EXONERATE_pepOSA_IRGSP → EXO_X_pepOSA_IR
 7_EXONERATE_pepOSA_IMSU → EXO_X_pepOSA_MS
 7_EXONERATE_pepPPE → EXO_X_pepPPE
 7_EXONERATE_pepPTR → EXO_X_pepPTR
 7_EXONERATE_pepSBI → EXO_X_pepSBI
 7_EXONERATE_pepVVI → EXO_X_pepVVI
 7_EXONERATE_pepZMA → EXO_X_pepZMA

 7_EXONERATE_protATH → EXO_X_protATH
 7_EXONERATE_protBDI → EXO_X_protBDI
 7_EXONERATE_protHOR → EXO_X_protHOR
 7_EXONERATE_protOSA → EXO_X_protOSA
 7_EXONERATE_protPPE → EXO_X_protPPE
 7_EXONERATE_protPTR → EXO_X_protPTR
 7_EXONERATE_protSAC → EXO_X_protSAC
 7_EXONERATE_protSBI → EXO_X_protSBI
 7_EXONERATE_protTRI → EXO_X_protTRI
 7_EXONERATE_protVVI → EXO_X_protVVI
 7_EXONERATE_protZMA → EXO_X_protZMA

 7_EXONERATE_protPlant → EXO_X_protPlant
 7_EXONERATE_protPoa → EXO_X_protPoa
 7_EXONERATE_protRos → EXO_X_protRos

 7_EXONERATE_refSeqPlnProt → EXO_X_refSeqPln
 7_EXONERATE_uniprot_sprot → EXO_X_uniprot_s
 7_EXONERATE_uniprot_trembl → EXO_X_uniprot_t

10_EXONERATE_pepATH → EXO_P_pepATH
 10_EXONERATE_pepBDI → EXO_P_pepBDI
 10_EXONERATE_pepOSA_IRGSP → EXO_P_pepOSA_IR
 10_EXONERATE_pepOSA_IMSU → EXO_P_pepOSA_MS
 10_EXONERATE_pepPPE → EXO_P_pepPPE
 10_EXONERATE_pepPTR → EXO_P_pepPTR
 10_EXONERATE_pepSBI → EXO_P_pepSBI
 10_EXONERATE_pepVVI → EXO_P_pepVVI
 10_EXONERATE_pepZMA → EXO_P_pepZMA

 10_EXONERATE_protHOR → EXO_P_protHOR
 10_EXONERATE_protSAC → EXO_P_protSAC
 10_EXONERATE_protTRI → EXO_P_protTRI

 10_EXONERATE_protPlant → EXO_P_protPlant
 10_EXONERATE_protPoa → EXO_P_protPoa
 10_EXONERATE_protRos → EXO_P_protRos

 10_EXONERATE_nr → EXO_P_nr
 10_EXONERATE_refSeqPlnProt → EXO_P_refSeqPln
 10_EXONERATE_uniprot_sprot → EXO_P_uniprot_s
 10_EXONERATE_uniprot_trembl → EXO_P_uniprot_t



12_BLASTN_embl_pln → **BLASTN_embl_pln**
12_BLASTN_embl_est_pln → **BLASTN_embl_est**
12_BLASTN_embl_gss_pln → **BLASTN_embl_gss**
12_BLASTN_embl_htc_pln → **BLASTN_embl_htc**
12_BLASTN_embl_htg_pln → **BLASTN_embl_htg**
12_BLASTN_embl_pat_pln → **BLASTN_embl_pat**
12_BLASTN_embl_sts_pln → **BLASTN_embl_sts**
12_BLASTN_embl_tsa_pln → **BLASTN_embl_tsa**
12_BLASTN_embl_wgs_pln → **BLASTN_embl_wgs**

12_BLASTN_genoAtTAIR → **BLASTN_genoAtTA**
12_BLASTN_genoBdJGI → **BLASTN_genoBdJG**
12_BLASTN_genoOsIRGSP → **BLASTN_genoOsIR**
12_BLASTN_genoOsMSU → **BLASTN_genoOsMS**
12_BLASTN_genoPpe → **BLASTN_genoPpe**
12_BLASTN_genoPtr → **BLASTN_genoPtr**
12_BLASTN_genoSbDOE → **BLASTN_genoSbDO**
12_BLASTN_genoVvi → **BLASTN_genoVvi**
12_BLASTN_genoZmMSO → **BLASTN_genoZmMS**

12_BLASTN_refSeq_chloro → **BLASTN_refSeq_c**
12_BLASTN_refSeq_mito → **BLASTN_refSeq_m**

13_EXONERATE_pepATH → **BLASTX_pepATH**
13_EXONERATE_pepBDI → **BLASTX_pepBDI**
13_EXONERATE_pepOSA_IRGSP → **BLASTX_pepOSA_IR**
13_EXONERATE_pepOSA_IMSU → **BLASTX_pepOSA_MS**
13_EXONERATE_pepPPE → **BLASTX_pepPPE**
13_EXONERATE_pepPTR → **BLASTX_pepPTR**
13_EXONERATE_pepSBI → **BLASTX_pepSBI**
13_EXONERATE_pepVVI → **BLASTX_pepVVI**
13_EXONERATE_pepZMA → **EblastX_pepZMA**

11_INTERPROSCAN → **misc_feature**

14_TRNASCAN-SE → **tRNA**

15_TRF → **tandem_repeat**



Check list 1

- ✓ Build a curation repertory without empty files / BLAST files / SIMsearch files
- ✓ File / Load data... / local file/**9_GENEMODEL.embl**
 - Zoom and move (4 ways - the best = up/down/right/left keys) => **2nd gene**
 - Track list
 - Exon/CDS/Gene Structure (hidden) – keep exon
 - Clic on exon / CDS - shadow bands (last exon)
 - Details on selected items (windows on the right)
 - CTRL E on exon (label=Category / functional Annotation)
 - Change label=3-isopropylmalate dehydrogenase 2 chloroplastic1
 - **problem with save**
 - See Query: **Google**, NCBI (Ensembl, Plaza)
 - Add Track exon, first exon ➔ Start / Stop / Intron-Exon junctions (problem first Do/Ac)
- ✓ Biological Evidences (BLASTn; BLASTx; BLASTp – best hits)
 - EXO_N / EXO_X / EXO_P ➔ move tracks / compact tracks
 - Correction of Do/Ac first 2 junctions – expand last exon
 - PoaFL / HvFL / CDS



Check list 2

- ✓ Biological Evidences (BLASTn; BLASTx; BLASTp – best hits)
 - EXO_X
 - EXO_P (case of empty file → see files on repertory)
 - See “Other files” → multiple alignment pepBDI & protHor
- ✓ Case of the last Gene (wrong gene model as seen on the Gbrowse)
 - Open all the following evidences
 - Wheat RNA-seq / FL-cDNA (keep FL-cDNA)
 - Ta_unigene (delete an item)
 - Hv_FL (keep FL-cDNA)
 - CDS / EXO_X (compact) / EXO_P
 - InterProScan

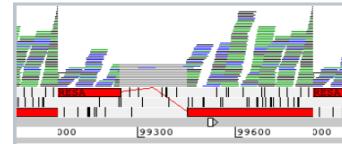


Check list 3

- An example with **Augustus Gene Model**
 - **Clone** feature – change type to CDS
 - Select 2 exons → **split** (CTRL/Caps)
 - Take the first Met
 - Clone wheat FL-cDNA and keep the missing exon => change type to CDS
 - **Merge** the two CDS
 - Check Expasy (*Viridiplantaea*, no filter, 50, 50, 0,0001)
 - Navigation/**Search** aa sequence
 - Keep this gene (label=pseudo)
- ✓ Conserved Non-coding Sequences – CNSs : 12_ ; 13_ ; tRNAs
- ✓ SSRs
- ✓ TEs (Wicker's code)

Check list 4

- ✓ Save annotation with an other name (gene/CDS/exon/mRNA)
 - See the file with wordpad
- ✓ Save Session
- ✓ Re open a previous analysis
- ✓ Recover previous annotation file → trick
- ✓ See the saved annotation with ARTEMIS
 - ✓ Sub sequence function





10 last minutes

10 last minutes

How to use the TriAnnot pipeline ?



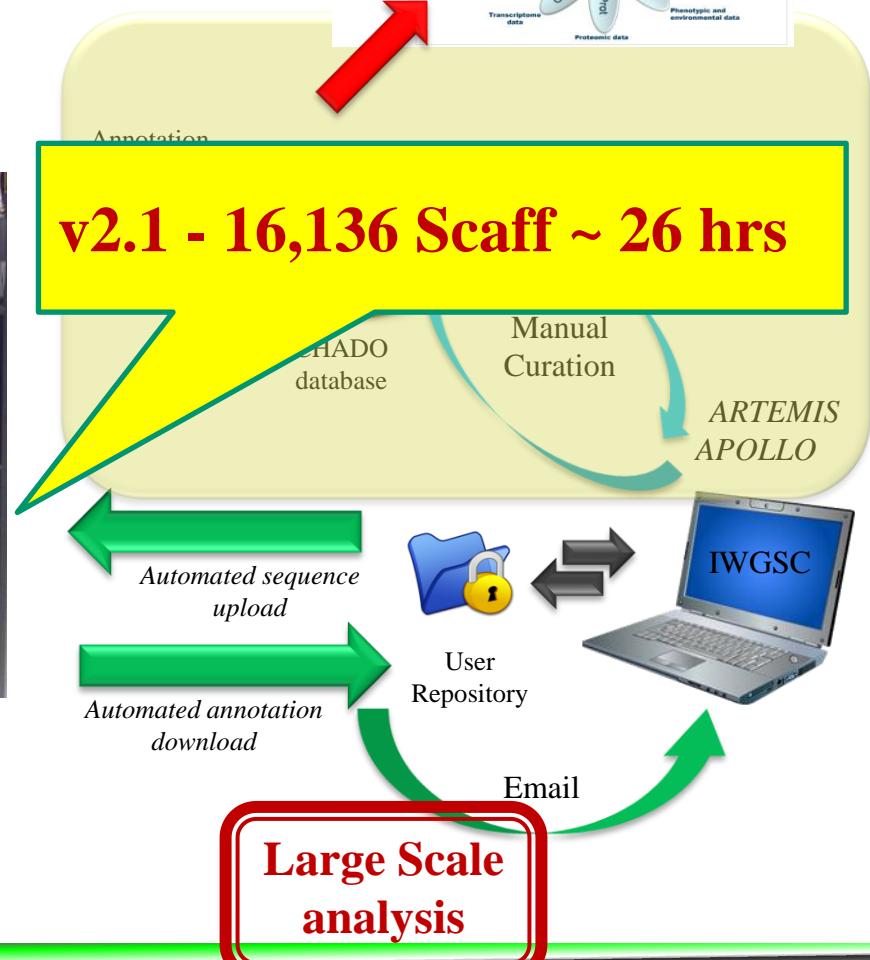
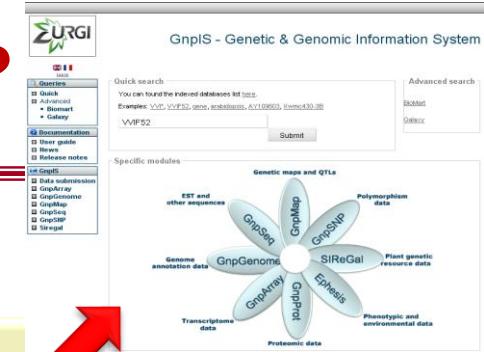
Web submission



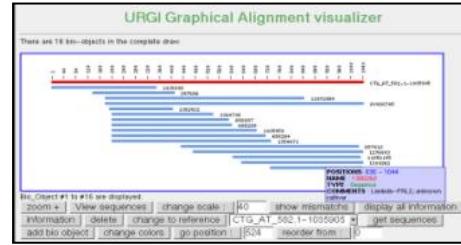
Email
GBrowse link
Download

~900 cores – 8.5 Tfop

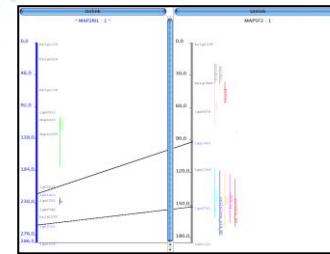
**Small Scale
analysis**



Integrative Databases for interoperability



Genetic maps and QTLs



EST and other sequences

Genome annotation data

Transcriptome data

GnpSeq

GnpGenome

GnpArray

GnpMap

GnpSNP

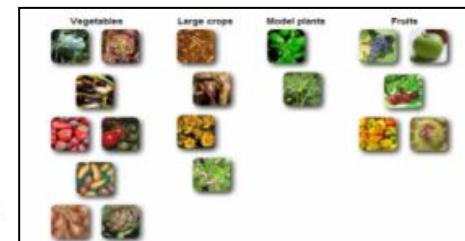
GnpProt

Polymorphism data

Plant genetic resource data

2	35	37	44	67	69	73	76	84	85	110	112	120	121	129	152	190	192	211	223	236	234	235	250	291	297
C	T	T	A	D	D	T	C	T	T	G	G	T	C	G	C	G	T	C	T	A	C	C	T		
C	T	T	A	T	-	T	C	T	G	G	T	C	G	C	B	T	C	C	T	A	C	C	T		
C	T	T	A	T	-	T	C	T	G	G	T	C	G	C	G	T	C	C	T	A	C	C	T		
C	T	T	A	T	-	T	C	T	G	G	T	C	G	C	G	A	G	T	C	T	A	C	C	T	

Phenotypic and environmental data



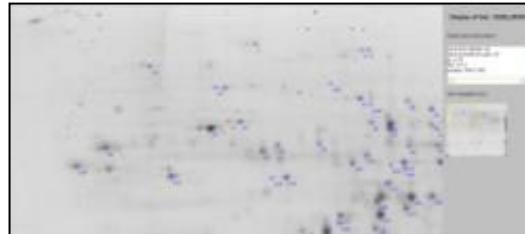
Hybridization results

Information

Results

Sample	Condition	Description	Replicates	Arrival date	Design status	Site identifier	Project	Plate	Well
1	14_pj000000	14_pj000000	1	2010-06-02	2010-06-02	14_pj000000	14_pj000000	14	1
2	14_pj000001	14_pj000001	1	2010-06-02	2010-06-02	14_pj000001	14_pj000001	14	2
3	14_pj000002	14_pj000002	1	2010-06-02	2010-06-02	14_pj000002	14_pj000002	14	3
4	14_pj000003	14_pj000003	1	2010-06-02	2010-06-02	14_pj000003	14_pj000003	14	4
5	14_pj000004	14_pj000004	1	2010-06-02	2010-06-02	14_pj000004	14_pj000004	14	5
6	14_pj000005	14_pj000005	1	2010-06-02	2010-06-02	14_pj000005	14_pj000005	14	6
7	14_pj000006	14_pj000006	1	2010-06-02	2010-06-02	14_pj000006	14_pj000006	14	7
8	14_pj000007	14_pj000007	1	2010-06-02	2010-06-02	14_pj000007	14_pj000007	14	8
9	14_pj000008	14_pj000008	1	2010-06-02	2010-06-02	14_pj000008	14_pj000008	14	9
10	14_pj000009	14_pj000009	1	2010-06-02	2010-06-02	14_pj000009	14_pj000009	14	10
11	14_pj000010	14_pj000010	1	2010-06-02	2010-06-02	14_pj000010	14_pj000010	14	11
12	14_pj000011	14_pj000011	1	2010-06-02	2010-06-02	14_pj000011	14_pj000011	14	12
13	14_pj000012	14_pj000012	1	2010-06-02	2010-06-02	14_pj000012	14_pj000012	14	13
14	14_pj000013	14_pj000013	1	2010-06-02	2010-06-02	14_pj000013	14_pj000013	14	14
15	14_pj000014	14_pj000014	1	2010-06-02	2010-06-02	14_pj000014	14_pj000014	14	15
16	14_pj000015	14_pj000015	1	2010-06-02	2010-06-02	14_pj000015	14_pj000015	14	16

Proteomic data





Perspectives БІОЗВЕСТЛЯГ

New Merge Model

New MergeGeneModel.pm

- Multiple input files
 - ✓ Integrate manually curated genes
- Filtering predictions
 - ✓ No overlapping evidences AND no BLAST hit against plants proteomes & Transposases
- Validation against evidences
- Scoring
 - Overlap with best blastp hit
 - % of identity
 - Validation score : n_{Valid}^2/n_{All}
 - Penalty for non canonical splice sites
- **Choice of best gene model**



S. Theil
GDEC



F. Choulet
GDEC

Two new independent code systems

Status

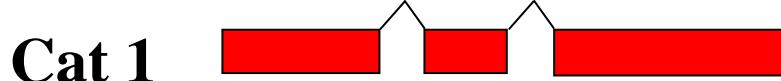
- ‘Full’
 - ✓ hit_coverage > 70%
- Pseudo
 - ✓ $50\% \leq \text{hit_coverage} \leq 70\%$
- Fragment
 - ✓ hit_coverage < 50%

Confidence

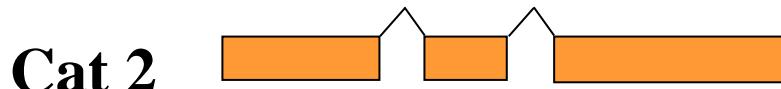
- High confidence (clear evidences: start/stop/intron-exon junctions)
- Low confidence (no clear evidences : start/stop/intron-exon junctions)

New color code system to assess the confidence of the structural annotation *

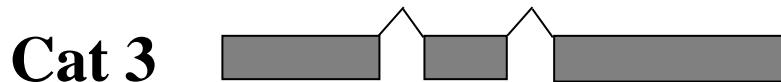
3 Categories



Full – *High Confidence* (1)



Full – *Low Confidence*



Pseudo & fragments

(1) This gene models do not need manual expertise

Wheat chromosome 3B annotation validation

Annotation on assembly V4.2 (expected overlap) – 8 789 Protein-coding gene models
After filtering

	high_confidence	low_confidence	
Full	4 136	2 477	6 613
Pseudo	769	692	1 461
Fragment	101	614	715
	56%	44%	

Manually expertised

Seq_problem : Ns | wrong splice site | frame shift

	Full	Pseudo	Fragment	
validated	3 029	470	25	
not_validated	622	305	184	
validated_seq_problem	19	212	3	
not_validated_seq_problem	30	87	9	
	3 700	1 074	221	4 995

Acknowledgements



P. Leroy
N. Guilhot
(S. Theil)
F. Choulet
C. Feuillet



C. Viseux
M. Alaix
H. Quesneville



H. Sakai
T. Itoh



C. Pozniack



C. Caron



L. Sterck



A. Mahul



D. Hill



Previous Collaborators

P. Dufour M. Seidel
L. Cerutti K. Mayer
B. Laubin H. Ohyanagi
F. Sabot
F. Legeai
E. Gicquelot
F. Giacomoni
C. Caron
I. Blanc-Lenfle
A. Mahul
A. Claude
M. Liauzu
T. Flutre
I. Luyten
C. Pelegrin
A. Bernard
O. Inizan
M. Reichstadt
S. Reboux
N. Amano
H. Numa
T. Tanaka

Thank You very much



Massif du Sancy
Joël Damase



Murol Castle

<http://www.clermont.inra.fr/triannot/>