

# THE DEVELOPMENT OF A NOVEL DATA MINING TOOL TO FIND *CIS*-ELEMENTS IN RICE GENE PROMOTER REGIONS

Koji Doi<sup>(1)</sup>, Aeni Hosaka<sup>(1)</sup>, Toshifumi Nagata<sup>(1)</sup>, Kouji Satoh<sup>(1)</sup>, Kohji Suzuki<sup>(2)</sup>, Ramil Mauleon<sup>(3)</sup>, Michael Jonathan Mendoza<sup>(3)</sup>, Richard Bruskiwich<sup>(3)</sup>, and Shoshi Kikuchi<sup>(1)</sup>  
 (1) National Institute of Agrobiological Sciences (2) Hitachi Software Engineering, Japan.Co., Ltd. (3) International Rice Research Institute

## Basic Feature of the Developed Tool

Information about over 35,000 full-length *Oryza sativa* cDNA and associated microarray gene expression data have enabled identification of conserved motifs in promoters of genes listed by microarray analysis. They are expected to act as *cis*-regulatory elements involved in key roles under the assayed experimental conditions.

We developed and have continued improvement of a novel tool to search *cis*-element candidates. Here we report the recent progress.

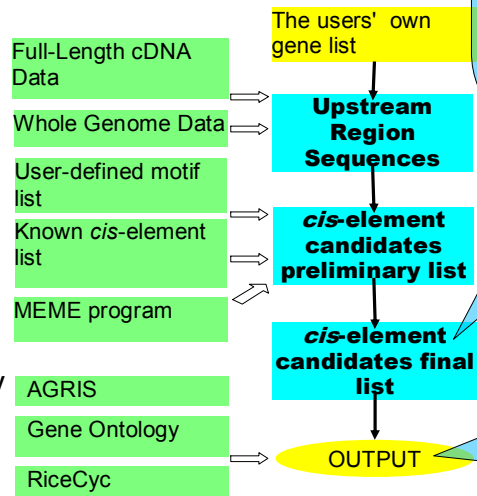
### The Important Features of the Developed *cis*-Element Search Tool

- It is publicly opened as a web application software.
- No special programming skill is required.
- Users can define own gene set. (e.g. by cluster analysis after microarray analysis)
- Motif search (MEME) and data mining method (association rule analysis) are performed by a pipeline system implemented in the host computer in IRRI.

**Basic Strategy for Search and Evaluation of *cis*-Element Candidates**

If motifs overrepresented in the upstream region of the focused genes, they could play specific roles on expression regulation of the genes.

### Flowchart of *cis*-Element Search



### Evaluation to Get Final *cis*-Element Candidate List

Candidates showing the highest likelihood (specificity) are retained in the final *cis*-element candidate list.

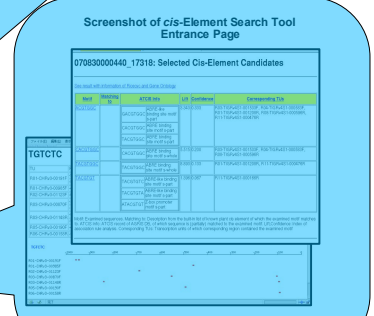
#### → Association Rule Analysis

X: The focused TU possess a *cis*-element candidate TGTTCT.  
 Y: The focused TU is corresponding to drought response.

#### A Fictional Example:

	Y		Total		
Yes	8	176	184	Support	8 / 20509 = 0.00039
No	1003	19322	20325	Confidence	8 / 184 = 0.04348
Total	1011	19498	20509	Lift	0.00039 / (184 * 20509 / 1011) = 0.88199

High lift value suggests strong relationship between X and Y. The developed tool defines default threshold lift value as 1.0, while users can change it. Thus, as for example above, relationships between TGTTCT *cis*-element and drought stress is not strongly supported.



## Recent Upgrade of the Tool

### Users can now select reference gene sets

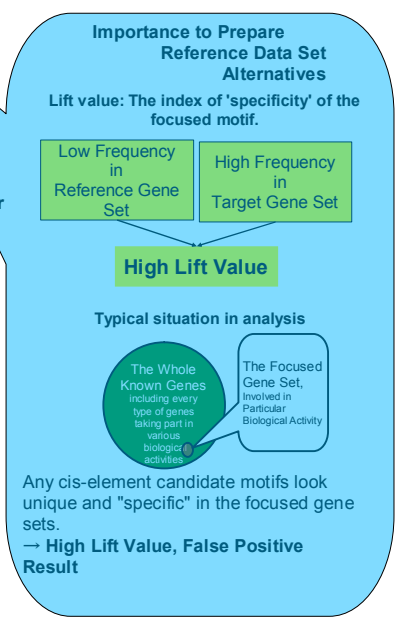
- Available or Preparing Reference Gene Sets:
- Whole set – including all genes in KOME database (default).
  - Sets of genes including particular motifs.
  - Sets of genes up-regulated in particular experimental condition (under construction).

We have confirmed that the contents of the reference have a great influence on result. Refinement of reference gene sets are still in progress.

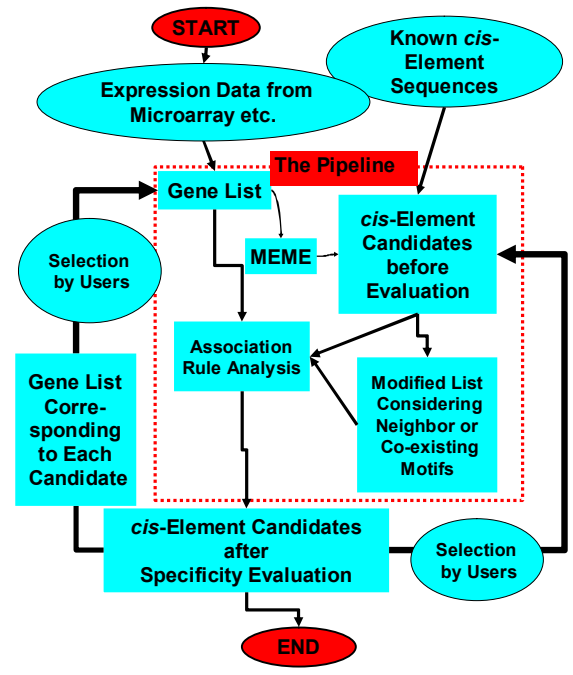
### Users can also specify *cis*-element candidate motifs

Perl-compatible regular expression powerfully support users to specify complex sequence patterns.

Example of regular expression for known *cis*-elements  
 Helix-turn-helix(HTH) (CTAATTG){2,3}  
 BBR/BPC ((GA)+(TC)+)  
 RAV CAACA[ACGT]\*CACCTG



### Advanced Data Mining Flowchart with Repeated Utilization of the *cis*-Element Search Tool



## Improvement of the method to obtain more accurate result

### Case Study for Feasibility Test of the Methods

#### Focused Topics:

**Drought Stress Response** – An important biological activity for survival of plants.

**Abscisic Acid (ABA) Response** – Closely related mechanism of drought stress response.

There are some known *cis*-elements that play roles on these activities, while it is considered that many unknown *cis*-elements still remain.

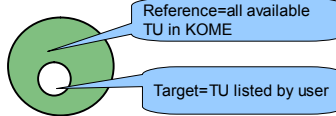
It is suitable for case study to evaluate feasibility of our method and developed tool.

We are considering new method to increase feasibility of the tool. Feasibility test is performed for the considered method.

### Current Issue and Method under consideration

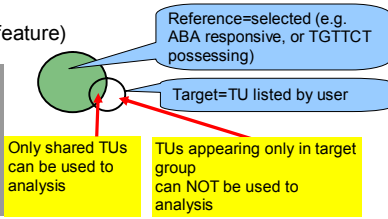
#### Default

Extraordinary large lift values tend to be obtained: Every candidate appears "unique"



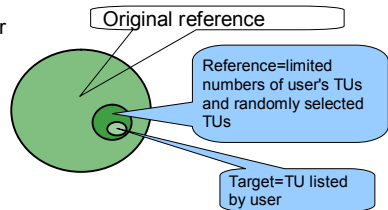
#### Switched Reference (newly implemented feature)

Very few samples can be investigated: Result with Low reliability



#### New approach under consideration

Clear result is expected



### Material and Method

#### Reference TU lists

- Two reference TU lists were prepared by literature survey.
- TUs containing ABRE motif (ACGTG[GT]C) in 1000bp of upstream region. (804TU).
- TUs containing ABA-responsive *cis*-elements motifs listed by literature survey in 1000bp of upstream region (11134TU)

#### User-defined TU lists

- Rabbani et al. (2003) showed a gene list (based on NCBI accession numbers) of some kinds of stress responsibility. Followings were from their work, that used as user-defined TU list for the feasibility test.
- ABA responsive genes (18TU)
- Drought responsive genes (24TU)

#### Analysis and evaluation of result

- All combination of reference and user-defined TU lists mentioned above were applied to analysis.
- Result was checked by simple negative control test.

### Negative Control Test

To check whether the found "positive" relationship can be considered true, following simple test was employed in this study.

- Select TUs randomly from the reference list, so that the size of result list is twice as large as user-defined TU list.
- Specificity of known *cis*-element motifs for user-defined TU list was evaluated with selected reference prepared in the former step by developed tool.
- Perform *cis*-element search.
- For randomly-selected TUs, lift values of *cis*-element listed in the examination should be negative, while they should be positive for real user-defined TUs.

An example of negative control test for Leucine zipper factors (bZIP: CACGTG)

cis-element	Lift Value	Negative Control Test using randomly selected TUs									
		50	100	500	1000	5000	10000	50000	100000	500000	1000000
LEA	2.81	0.12	0.13	0.11	0.11	0.13	0.11	0.11	0.11	0.11	0.11
WRKY	1.95	0.11	0.12	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11
MYB	1.81	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11
ERF1	1.91	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11
ERF2	1.81	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11
ERF3	1.81	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11

#### Previous Related Studies

As far as our literature survey, following studies are especially notable to consider the feasibility of our method.

**WRKY** – WRKY TFs are known as factors involving various biological reactions such as disease resistance. Otherwise, Mare et al. (2004) reported that a WRKY TF (Hv-WRKY38) is induced by drought stress in *Hordeum vulgare*.

**bZIP** - Responsibility of "water deficit stress" in *Phaseolus* is reported (Rodriguez-Urbe and O'Connell, 2006).

**JUMONJI** - Senthil-Kumar (2007) pointed out a JUMONJI TF showed relative drought tolerant phenotype in *Nicotiana benthamiana*.

As for LIM finger, we have not found reports discussing the relationships to stress response definitely.

### Result of feasibility test

Reference TU Set	ABA		Drought		ABA		Drought	
	ABRE	ABAresp	ABRE	ABAresp	ABRE	ABAresp	ABRE	ABAresp
Length of examined sequence	50	100	50	1000	50	1000	50	1000
AP2	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
ARF	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
ARF1	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
ARF2	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
BZR(BES1)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
CCATbox	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
C/7He/2c/4c zinc finger	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Cys2His2 zinc finger	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Cys2His2 zinc finger/Ringer	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Dof	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
DPE	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
ERF1	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
ERF2	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
ERF3	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
GATA-Factors	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Hella-loop-helix factors(HLH)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Hella-loop-helix factors(HLH)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Hella-loop-helix factors(HLH)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Hella-kum-helix(HKH)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Hella-kum-helix(HKH)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Hella-kum-helix(HKH)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
JUMONJI	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
JUMONJI	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
JUMONJI	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
LEAFY	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Leucine zipper factors(bZIP)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Leucine zipper factors(bZIP)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
LIM finger	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
MADS(CarG boxes)	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
Myb	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
RAW	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
TATAbox	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
YOZ-8	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
WRKY	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
WRKY	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11
YABBY	2.81	1.95	0.11	0.11	0.11	0.11	0.11	0.11

cis-elements showing positive correlation

lift values supported by negative control test

lift values not supported by negative control test

Lift values of known *cis*-element motifs for each examination are shown in the table. Values lower than 1.0, recognized as negative, are not shown.

### Discussion

Due to refinement on demand, now users can perform trial and error from various angles in using the *cis*-element search tool. The experimental version of the developed tool is still maintained to try refinement of data mining protocol.

The result of feasibility test is concordant with some previous reports. The number of detected *cis*-elements was limited, but the result suggested that plausible candidates can be obtained after slight trial and error.

We consider that the usefulness of the developed tool is ensured more due to update of the tool. We are going to reflect such attempt to the opened tool.

**Availability of the tool**  
<http://hpc.irri.cgiar.org/tool/nias/ces>

**Any comments are greatly appreciated**  
 Koji Doi (kdoi@affrc.go.jp)

### Future Work

- Refinement of the reference gene sets**
  - Literature search of known *cis*-elements, and preliminary search of found motifs to select genes containing them.
  - Collaboration with experimental researchers to create precise gene list responsive to particular biological phenomena.
- Biological discovery**
  - The project will successfully end when suggestion of novel *cis*-elements are persuasively achieved.