# Conserved control signals in the transcriptome of higher plants

**Khanh Tran**

**Thesis submitted for the degree of Doctor of Philosophy**

**May 2010**

**Discipline of Plant and Pest Science**

**School of Agriculture, Food, and Wine**

**The University of Adelaide**

# TABLE OF CONTENTS

# ABSTRACT

Understanding the mechanisms that regulate gene expression is an important goal in bioinformatic research. There are two major levels of gene regulation: transcriptional and post-transcriptional control. Much attention has been directed to transcriptional control, but it is now clear that the untranslated regions (UTRs) of messenger RNA (mRNA) also play an important role in post-transcriptional control of gene expression. Two important control signals found in 5′-UTRs of both animal and plant mRNAs are stem-loop motifs and upstream open reading frames (uORFs).

One strategy for identifying functional uORFs in plants is to use a comparative approach (Crowe et al. 2006; Hayden and Jorgensen 2007; Pavesi et al. 2007). There are extensive EST datasets for five important cereal crops (rice, wheat, barley, maize, and sorghum). Rice is the best characterised of these cereals with a sequenced genome (Yu et al. 2002) and a cDNA database containing 32,000 clones that are enriched for 5′ full-length sequences (Kikuchi et al. 2003). In this research, comparative R-nomics was used to identify conserved stem-loop motifs and uORFs in cereals using publicly available assembled EST data.

To determine the prevalence of 5′-UTR stem-loop structures in plants a bioinformatics pipeline was developed to predict secondary structures. The pipeline used a program called RNAProfile to predict stem-loops that are conserved in both sequence and structure. The findings from this study concluded that conserved 5′-UTR stem-loops in long 5′-UTRs (200 to 1200 nt) are rare (~8%) in the cereal transcriptome, the genes themselves that contain conserved 5′-UTR stem-loop motifs are spread across different functions, and appear to have a biological role based on higher structure than sequence conservation in at least three out of four cereal species.

Another control signal that is involved in post-transcriptional control is the uORF. A recent study in distantly related plants, such as rice and

Arabidopsis, found that uORFs are rare in these transcriptomes (Hayden and Jorgensen 2007), but it is unclear how prevalent uORFs are in closely related plants. To address this question, the bioinformatics pipeline was modified to use a program called uORFSCAN to find conserved uORFs in five cereals that could potentially regulate translation. Major conclusions from this study are that the identified uORFs are highly conserved (50% median amino acid sequence similarity), are rare in cereal transcriptomes (<150 loci contain them), are generally short (less than 100 nt), position independent in their 5′-UTRs, and their start codon context and the usage of rare codons do not appear to be important for translation.

Two candidate uORFs were selected for mutational analyses, and a quantitative *in vitro* transcription and translation system was used to determine if they function in translational control. The rice *SAMDC* small and *S6K* long uORFs were shown to be capable of down-regulating translation of a luciferase reporter gene. This study has provided evidence, for the first time, that the *S6K* uORF is involved in controlling translation. In conclusion, this study has identified new genes that may be controlled at the level of translation by stem-loop motifs and conserved uORFs.

# DECLARATION

This work contains no material which has been accepted for the award of any other degree or diploma in any university or other tertiary institution to Michael Khanh Tran and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text.

I give consent to this copy of my thesis when deposited in the University Library, being made available for loan and photocopying, subject to the provisions of the Copyright Act 1968.

The author acknowledges that copyright of published works contained within this thesis (as shown below) resides with the copyright holder/s of those works.

I also give permission for the digital version of my thesis to be made available on the web, via the University's digital research repository, the Library catalogue, the Australasian Digital Theses Program (ADTP) and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Chapter 3 contains material from the following publication:

Tran, M.K., C.J. Schultz, and U. Baumann. 2008. Conserved upstream open reading frames in higher plants. *BMC Genomics* **9:** 361-378.

Signature:

Date:

# ACKNOWLEDGEMENTS

# LIST OF TABLES

**CHAPTER 3**

# LIST OF FIGURES

## CHAPTER 1

## CHAPTER 2

# ABBREVIATIONS

| | |
|---|---|
| aa | amino acid |
| ADH1 | alcohol dehydrogenase-1 |
| AGRF | Australian Genome Research Facility |
| ARF | auxin response factor |
| Avg. | average |
| BDT | Big Dye Terminator |
| BLAST | Basic Local Align Search Tool |
| bp | base pair |
| bZIP | basic region leucine zipper |
| CBL | calcineurin B-like |
| CDS | coding sequence |
| CNS | conserved non-coding sequence |
| CPA1 | carbamoyl phosphate synthetase |
| CST | conserved sequence tag |
| DFCI | Dana Farber Cancer Institute |
| eIF4F | eukaryotic initiation factor 4F |
| EST | expressed sequence tag |
| *ETT* | *ETTIN* |
| FL | full-length |
| GO | gene ontology |
| GPRM | genetic programming for RNA motifs |
| GUS | β-glucuronidase |
| Hb | *Hordeum bulbosum* |
| Hv | *Hordeum vulgare* |
| INDEL | insertions and deletions |
| IRE | iron responsive element |
| IRES | internal ribosome entry site |
| IRP | iron responsive protein |
| KOME | Knowledge-Based Oryza Molecular Biological Encyclopedia |

| | |
|---|---|
| Lp | *Lolium perenne* |
| LUC | luciferase |
| MOPS | 3-[N-Morpholino]propanesulfonic acid |
| mORF | main open reading frame |
| *MP* | *MONOPTEROS* |
| NELF | negative elongation factor |
| nt | nucleotide |
| NMD | nonsense-mediated decay |
| PC | *Phalaris coerulescens* |
| POS | positive control |
| rbcs | ribulose biphosphate small subunit |
| rbh | reciprocal best hit |
| RBP | RNA-binding protein |
| RNA | ribonucleic acid |
| RNApromo | RNA prediction of motifs |
| RNase P | ribonuclease P |
| *S6K* | S6 ribosomal kinase |
| *SAMDC* | S-adenosylmethionine decarboxylase |
| SAPAC | South Australian Partnership for Advanced Computing |
| Sc | *Secale cereale* |
| SCFG | stochastic context-free grammars |
| SECIS | selenocysteine insertion sequence |
| SEM | standard error of mean |
| SLASH | Stem-Loop Align SearcH |
| Taes | *Triticum aestivum* |
| TC | tentative contig |
| TIGR | The Institute for Genomic Research |
| TIS | translation initiation site |
| *Trx4* | *thioredoxin4* |
| UniProtKB | UniProt Knowledgebase |
| uATG | upstream start codon |
| uORF | upstream open reading frame |

| URL | universal resource locator |
|-----|----------------------------|
| UTR | untranslated region |