

# THE CENTER FOR EUKARYOTIC STRUCTURAL GENOMICS AT THE UNIVERSITY OF WISCONSIN-MADISON



Frits Abildgaard, <sup>1</sup>David J. Aceti, <sup>1</sup>Craig Bingman, <sup>1</sup>Paul Blommel, <sup>1</sup>Heather Burch, <sup>1</sup>John Cao, <sup>1</sup>Claudia Cornilescu, <sup>1</sup>Jurgen Doreleijers, <sup>1</sup>David Dyer, <sup>1</sup>Hamid Eghbalnia, <sup>1</sup>Brian G. Fox, <sup>1</sup>Ronnie O. Frederick, <sup>1</sup>Adrian Hegeman, <sup>1</sup>Won Bae Jeon, <sup>1</sup>Todd Kimball, <sup>1</sup>Kelly Kjer, <sup>1</sup>Peter T. Lee, <sup>1</sup>Jing Li, <sup>2</sup>Betsy Lytle, <sup>2</sup>John L. Markley, <sup>1</sup>Ramya K. Narayana, <sup>1</sup>Craig S. Newman, <sup>1</sup>George N. Phillips, <sup>1</sup>Francis Peterson, <sup>1</sup>Ivan Rayment, <sup>1</sup>Bryan Ramirez, <sup>1</sup>Mike J. Runnels, <sup>1</sup>Kory Seder, <sup>1</sup>David Smith, <sup>1</sup>Mike Sussman, <sup>1</sup>Sandy Thao, <sup>1</sup>Eldon L. Ulrich, <sup>1</sup>Dmitry Vinarov, <sup>1</sup>Brian F. Volkman, <sup>2</sup>Gary Wesenberg, <sup>1</sup>W. Milo Westler, <sup>1</sup>Russell L. Wrobel, <sup>1</sup>Jianhua Zhang, <sup>1</sup>Qin (Kate) Zhao, <sup>1</sup>Zsolt Zolnai

<sup>1</sup>Department of Biochemistry, University of Wisconsin-Madison, 433 Babcock Drive, Madison, Wisconsin 53706, and <sup>2</sup>Department of Biochemistry, Medical College of Wisconsin, 8701

Watertown Plank Road, Milwaukee, Wisconsin 53226

## ABSTRACT

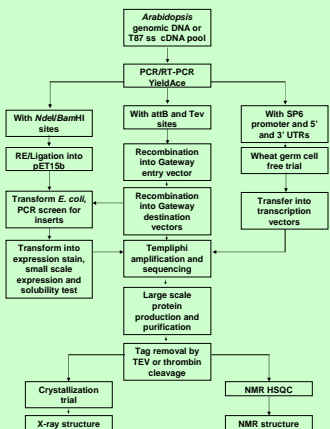
The major objectives of the University of Wisconsin Center for Eukaryotic Structural Genomics (CESG) are the development and critical analysis of methodologies for high-throughput, proteome-scale, eukaryotic protein production using *Escherichia coli* and cell-free translation, protein characterizations such as S-Tag detection, solubility, activity, and mass spectral analysis, and high resolution structure determination. The ultimate goal is to develop high-throughput structural biology methodologies to aid in the visualization of the complete diversity in naturally occurring protein structures. In this poster we summarize some of the strategies and rationale of our high-throughput protein production efforts as applied to the model photolithotropic eukaryote, *Arabidopsis thaliana*. In this system, protein targets were chosen from a prioritized list of about ~25,000 predicted *Arabidopsis* genes; we have focused on proteins whose primary sequences suggest novel folds or known folds accompanied by unusual biological functions. To meet the needs of this ambitious project, numerous supporting methodologies that allow for specialized gene cloning, gene expression evaluation, and protein purification are being developed. Finally, high-resolution structure determination of proteins is being solved using both NMR spectroscopy and X-ray crystallography methods.

## INTRODUCTION

A long-range goal in high-throughput structural biology will be to more completely map the naturally occurring diversity in protein structures. In this approach, the initial emphasis has been directed toward proteins whose sequences suggest a novel fold, proteins associated with novel functions, or proteins likely to have a known fold but a function not previously associated with that fold. In order to achieve the logistical demand of these ambitious programs, supporting methodologies that allow for specialized gene cloning, expression evaluation, and protein purification will be required. High resolution structure determination studies by NMR spectroscopy and X-ray crystallography will require the production of milligram amounts of soluble, native, and pure proteins. These efforts first become feasible by the wealth of information arising from genomic sequencing efforts, where the assumed proper identification of open reading frames has set the stage for massive scale gene cloning, followed by expression testing, and purification. The major aims of the University of Wisconsin Center for Eukaryotic Structural Genomics are the development and critical analysis of methods for high-throughput, proteome-scale, eukaryotic protein production, characterization, and structure determination. Here we describe some of the strategies and rationale that form the basis for our high-throughput protein production efforts.

## FLOWCHART OF THE CESG STRATEGY

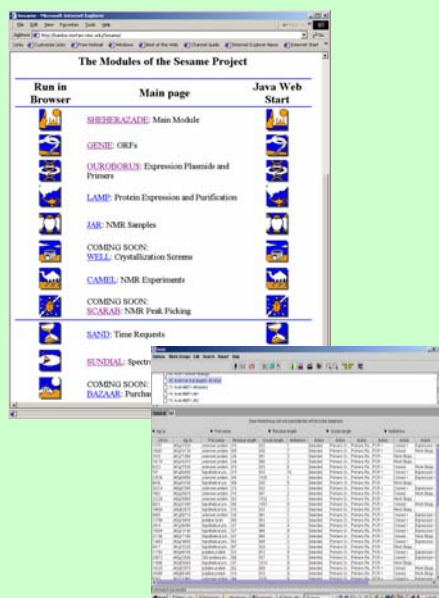
The following flowchart represents the strategy being developed at CESG. Points of interest include the use of *Arabidopsis thaliana* callus tissue cell culture as the transcript source, cloning by both restriction/ligation and recombinational methods, wheat germ cell-free protein synthesis as an alternative to *E. coli* heterologous expression, and Tempilphi (Amersham Biosciences Sunnyvale, California) to generate DNA for sequencing.



## SESAME: A LABORATORY INFORMATION MANAGEMENT SYSTEM

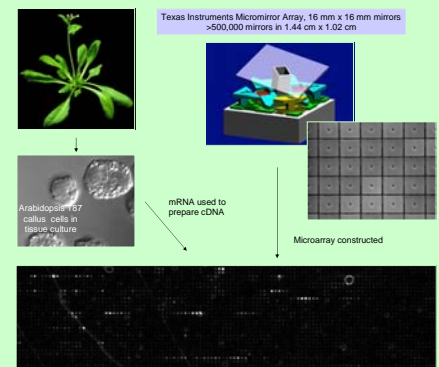
CESG utilizes the SESAME program developed in NMRFAM, UW-Madison, to track all experimental steps.

**Top:** From this homepage (<http://kamba.nmrfam.wisc.edu/Sesame/>), modules can be accessed for data storage and management of target ORFs, expression vectors, protein production, NMR samples, NMR experiments, scheduling of NMR spectrometers, and crystallization screens. **Bottom:** Genie is used for select targets, generating primers and tracking progress. This page of the Genie module is used for working with an individual "workgroup" (a set of target ORFs, usually numbering 96, that are processed as a unit). Actions and results describing the progress of individual ORFs from amplification from cDNA through structure determination can be entered.



## PROFILING GENE EXPRESSION IN ARABIDOPSIS T87 CELLS

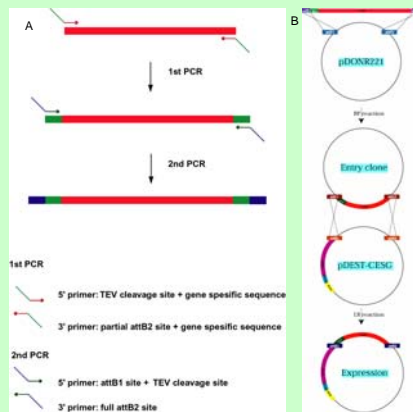
About 60-80% of all *Arabidopsis* genes may be expressed by the T87 callus cell line upon the basis of DNA microarray analyses. These results are being integrated into GENIE to optimize workgroup generation and to guide researchers in the selection of further experimental conditions.



## GATEWAY CLONING AT CESG

CESG uses Gateway technology (Invitrogen) to generate expression vectors that provide simple swapping of expression contexts for the target protein.

A. Amplification of an ORF from cDNA by PCR with the addition of recombination (*attB*) and TEV protease cleavage sites. B. Recombinational cloning of the PCR-generated insert to give entry and expression clones.



## EXPRESSION VECTORS DEVELOPED AT CESG

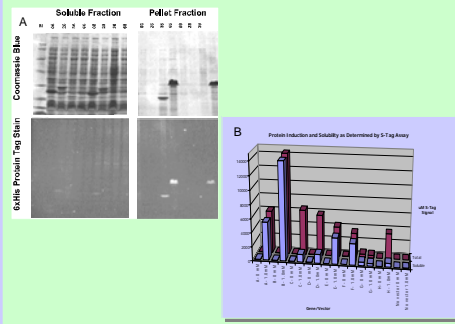
CESG has developed expression vectors that add S-tag (for detection), a His<sub>6</sub>-tag (for purification), and MBP (maltose-binding protein, for solubilization and purification) to the N-terminal of the target protein. When required, the entire fusion is cleavable from the protein target by TEV protease. These vectors are derived from pET (Novagen), pQE (Qiagen), and pBAD (Invitrogen) vector backbones.

Tag Name	Tag Size	Tag Identity	Location	Solid Phase	Insert Recovery	Recovery Conditions	Supplier	Removable Starter Phase	K <sub>d</sub>	Capacity (mg/mL)
His-tag	6aa	polyhistidine	N or C	N/A/T/A	Instalock	Mid	Novagen	Many	Yes	5-10
S-tag	15aa	PhosS	N or C	S-protein Agarose	His, Phos Agarose, NiNTA, Strep	Phenyl sepharose	Novagen	Yes	-1.6d	0.5
MBP	~40 Kda	Maltose binding protein	N	Amylose resin	Maltose	Mid	NEB	Yes		2-4

## DETECTION AND QUANTITATION WITH MULTIPLE TAGS

*E. coli* BL21-AI cells containing expression vectors encoding various multiple tag fusion proteins were grown to mid log phase and induced with IPTG and two different levels of arabinose.

A. Denaturing electrophoresis of soluble and pellet fractions with Coomassie Blue staining. B. Samples as shown in A treated with His-Tag stain. C. S-Tag detection analysis comparing the fluorescence obtained in the total cell lysate and in the soluble fraction.



## STATUS OF ARABIDOPSIS ORFS TARGETED BY CESG

Currently, ~1600 *Arabidopsis* ORFs have been targeted for investigation by CESG. Among these, ~1000 have been amplified by PCR, ~500 have been assembled into an expression vector, ~200 have been tested for expression, ~70 have given favorable solubility, ~30 have entered the purification stream, and 14 are currently in NMR or crystallization trials. The 367 "stopped" actions indicate a termination of workflow on an ORF due to a negative result, which may currently include decisions arising from the methods development process. The percentage of success reported for each action was calculated from the total number of ORFs tested in each category.

Status	# ORFs	Percentage	Status	# ORFs	Percentage
Selected targets	1627		Purified -	10	39
PCR +	979	75	Crystal +	4	57
PCR -	323	25	Crystal -	3	43
Cloned +	530	76	HSQC +	2	33
Cloned -	164	24	HSQC -	4	67
Expressed +	130	66	Quality X +	1	100
Expressed -	66	34	Quality X -	0	0
Soluble +	68	57	Structure +	0	
Soluble -	51	43	Structure -	0	
Purified +	16	61	Stopped	367	

## DISTRIBUTION OF CESG TARGETS

This map represents the approximately 25,000 ORFs targeted for *Arabidopsis thaliana*, arranged by priority of analysis by x-ray crystallography (determined by a number of factors, see <http://uwstructuralgenomics.org/target/target.htm>). Those targeted by CESG thus far are shown with status indicated by color.



## CONCLUSION

CESG is using *Arabidopsis thaliana* as a model system to develop methodologies for high-throughput eukaryotic protein production and high-resolution structure determination by NMR spectroscopy and X-ray crystallography. We are using both recombinational and restriction enzyme/ligation methods for gene cloning. Currently, ~1600 ORFs have been targeted by CESG. Greater than 95% of targeted intronless ORFs from genomic DNA of intronless genes and ~60% of intron-containing targets have been successfully amplified. Of these ~1000 ORFs, ~70 soluble protein products have been advanced to more detailed purification and screening efforts. Future plans revolve around increasing the availability of protein targets by systematization of the protocols described here.

CESG is supported by the National Institute of General Medical Sciences through the Protein Structure Initiative (P50 GM64598)