

### Bioinformatics at the Center for Eukaryotic Structural Genomics (CESG)

Craig S. Newman, Zhaohui Sun, Chris Oldfield, Ryan Bannen, Ip Kei Sam, Kaushik Raghunath, Frank C. Vojtki, Zsolt Zolnai, Brian G. Fox, George N. Phillips Jr., John L. Markley, and Eldon L. Ulrich

University of Wisconsin-Madison, 433 Babcock Drive, Madison, Wisconsin, USA 53706-1549, <http://www.uwstructuralgenomics.org>

### Abstract

CESG is one of nine pilot centers supported by the NIH Protein Structure Initiative (PSI). The aim of CESG is to expand our understanding of fold space and to elucidate novel structure-function relationships by solving 3D structures of a large number of individual proteins in a high throughput manner. Projects such as this generate tremendous amounts of data and require the support of strong bioinformatics and data management efforts. The Bioinformatics Team is charged with evaluating potential targets on a genome-wide basis and with prioritizing the predicted ORFs on the basis of a variety of parameters. We present our scoring protocol and illustrate its use in analyzing data. The Bioinformatics Team also harvests data from the Sesame laboratory information management system and other sources to analyze protein, identify bottlenecks, and uncover factors that influence the success of individual steps in the pipeline.

### CESG target scoring process

Potential targets are scored on the basis of a variety of quantitative and predicted characteristics identified by software packages developed in house, by collaborators or licensed from academic sources. Characteristics (and the relevant tools) used to score targets include:

- homology to proteins in the PDB and in the TargetDB that are at nearing structure completion (BLAST)
- presence of target "fragments" which lack homology to proteins in the PDB and in the TargetDB (BLAST)
- colled coil domains (COLIS)
- transmembrane segments (SCOP classification, TMHMM, and HMMTOP)
- signal peptides (SignalP and TargetP)
- number of Cys residues
- percent low complexity (seg)
- microarray analysis of *in vivo* transcription (where available)
- predicted disorder (PONDR)
- similarity to transposable elements (Censor)
- general suitability (on the basis of qualitative annotations, literature searches and knowledge and experience from within the CESG team)

Based on the results of these analyses, each ORF is allocated a 14 digit "target score" (see Figure 1). In each category, low digits are considered to be desirable. These scores are then used to divide the ORFs into priority tiers. An additional digit, representing the overall tier into which the ORF has been placed, is then added to the beginning of the target score. ORFs in the first, second, or third tiers are considered prime targets.

To date, 69,886 unique targets have been evaluated with the CESG target scoring protocol. This includes the entire *Arabidopsis thaliana* genome (version 5.0), non-redundant clones from the Mammalian Gene Collection (including sequences from human, mouse, rat, zebrafish, and *Xenopus*) and a number of community requested targets. As the new proteome predictions become available, additional species can easily be analyzed and prioritized for entry into the CESG protein production pipelines.

### Summary: criteria used in prioritizing CESG targets

Homology to proteins in the PDB and in the TargetDB that are at nearing structure completion (BLAST)	0	1	2	3	4	5	6	7	8	9
Presence of target "fragments" which lack homology to proteins in the PDB and in the TargetDB (BLAST)	0	1	2	3	4	5	6	7	8	9
Colled coil domains (COLIS)	0	1	2	3	4	5	6	7	8	9
Transmembrane segments (SCOP classification, TMHMM, and HMMTOP)	0	1	2	3	4	5	6	7	8	9
Signal peptides (SignalP and TargetP)	0	1	2	3	4	5	6	7	8	9
Number of Cys residues	0	1	2	3	4	5	6	7	8	9
Percent low complexity (seg)	0	1	2	3	4	5	6	7	8	9
Microarray analysis of <i>in vivo</i> transcription (where available)	0	1	2	3	4	5	6	7	8	9
Predicted disorder (PONDR)	0	1	2	3	4	5	6	7	8	9
Similarity to transposable elements (Censor)	0	1	2	3	4	5	6	7	8	9
General suitability (on the basis of qualitative annotations, literature searches and knowledge and experience from within the CESG team)	0	1	2	3	4	5	6	7	8	9
Overall Target Score	0	1	2	3	4	5	6	7	8	9
Priority Tier	0	1	2	3	4	5	6	7	8	9

Figure 1. Each scoring criteria (columns 2 through 14) is given a score ranging from 0 to 9. Low scores are considered to be desirable. The overall target score is limited by the individual criteria scores. The "Flam Score," "New Fold Prediction" and "Solubility Prediction" categories are not currently in use. Gene chip results are currently only available for targets from *Arabidopsis*. For examples of target scores, see Table 1.

Table 1. Target Scores for the Most Recent *E. coli* Pipeline Produced Structures

Gene ID	Source/Protein	Weight Score	Homology to proteins in the PDB and in the TargetDB that are at nearing structure completion (BLAST)	Presence of target "fragments" which lack homology to proteins in the PDB and in the TargetDB (BLAST)	Colled coil domains (COLIS)	Transmembrane segments (SCOP classification, TMHMM, and HMMTOP)	Signal peptides (SignalP and TargetP)	Number of Cys residues	Percent low complexity (seg)	Microarray analysis of <i>in vivo</i> transcription (where available)	Predicted disorder (PONDR)	Similarity to transposable elements (Censor)	General suitability (on the basis of qualitative annotations, literature searches and knowledge and experience from within the CESG team)	Overall Target Score	Priority Tier
B000094	Mouse	1000001000000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG6262	Arabidopsis	2000001000000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG6263	Arabidopsis	2000001000000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG6272	Arabidopsis	2000001000000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG6280	Arabidopsis	2000001000000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG6175	Arabidopsis	2000001000000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG5000	Arabidopsis	30000000004000	3	0	0	0	0	0	0	0	0	0	0	0	0
ASG1460	Arabidopsis	20000100100000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG1470	Arabidopsis	20000100100000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG1754	Arabidopsis	60000000000000	6	0	0	0	0	0	0	0	0	0	0	0	0
ASG1400	Arabidopsis	20000100100000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG1420	Arabidopsis	20000100100000	2	0	0	0	0	0	0	0	0	0	0	0	0
ASG14000	Arabidopsis	99990001000000	9	9	9	9	9	9	9	9	9	9	9	9	9
ASG14100	Arabidopsis	99990001000000	9	9	9	9	9	9	9	9	9	9	9	9	9
ASG14000	Arabidopsis	99990001000000	9	9	9	9	9	9	9	9	9	9	9	9	9
ASG14000	Arabidopsis	99990001000000	9	9	9	9	9	9	9	9	9	9	9	9	9

Table 2 uses the final target scores for the most recent CESG structures to illustrate our target scoring protocol. For each target in a lower tier (2 through 9), the criterion that determined the overall Selection Tier is highlighted; in some cases, two different criteria relegate a target to the same lower tier. Note that scores are updated routinely and that the target score at the time of selection may not be the same as the score at the time of structure completion.

### Current scoring criteria

Some of the key CESG scoring criteria are as follows:

**Selection Tier** represents the overall ranking of the target as determined by the values in the remaining 14 categories

**Homologous Structure** reflects homology to a previously solved structure deposited in the Protein Data Bank (PDB)

- 0 - expectation value greater than E-6
- 9 - expectation value less than E-6 and percent identity >29%

**Fragment Score** identifies targets with significant stretches of sequence not homologous to a known structure

- 0 - percent identity <29% over at least 100 residues against the PDB and an expectation value greater than E-3
- 9 - with this score, the Homologous Structure score is ignored in determining the overall Selection Tier

**Transmembrane** reflects the number of predicted transmembrane segments

- 0 - no predicted segments or no data
- 1 - one predicted segment
- 9 - more than one predicted segment

**Cysteine Count** reflects the number of cysteine residues present in a protein

- 0 - 0 cysteine residues
- 1 - 1 cysteine residue
- 2 - 2 cysteine residues
- 3 - 3 cysteine residues
- 4 - 4 cysteine residues
- 5 - 5 or 6 cysteine residues
- 6 - 7 or 8 cysteine residues
- 7 - 9 or 10 cysteine residues
- 9 - greater than 10 cysteine residues

**Low Complexity** reflects the amino acid complexity of the target

- 0 - less than 5% of the target is determined to be of low complexity
- 1 - between 5% and 10% of the target is determined to be of low complexity
- 2 - between 10% and 15% of the target is determined to be of low complexity
- 3 - between 15% and 20% of the target is determined to be of low complexity
- 9 - greater than 20% of the target is determined to be of low complexity

**Intrinsic Protein Disorder Prediction** reflects the predicted intrinsic disorder for a target

- 0 - less than 30% disorder predicted
- 4 - between 30% and 35% disorder predicted
- 5 - between 35% and 40% disorder predicted
- 6 - between 40% and 45% disorder predicted
- 7 - between 45% and 50% disorder predicted
- 9 - greater than 50% disorder predicted or greater than 30% disorder predicted and the longest disordered segment greater than or equal to 40 residues

### Analysis of experimental results

To evaluate the efficacy of the CESG score scheme, results from the CESG *E. coli* protein production pipeline have been examined. Presented below is a comparison of data from the protein production and structure determination pipeline versus key characteristics used in generating target scores. Our analysis shows that the CESG target score can predict the potential success of a target and that the success of the protein purification step seems to be most reflective of the target score. Table 2 shows that targets from the top three tiers yield 10 mg or more of purified protein at almost twice the frequency of targets in tiers 4, 7, or 8. A number of the specific criteria used to assign the Selection Tier seem to impact the purification success rate (Table 3). The cysteine count of a protein, number of transmembrane segments, amino acid complexity and intrinsic disorder all seem to have an effect on the ability to recover sufficient amounts of purified protein. Success at other steps of the protein production and structure determination pipeline (i.e., cloning, protein expression, NMR spectroscopy, and X-ray crystallography) show less correlation with target score. However, some key characteristics do appear to predict success at specific stages. The intrinsic protein disorder score of a target correlates well with the production of a folded protein for NMR spectroscopy (see below for a more detailed discussion) while an increasing number of transmembrane domains in a target adversely affects cloning efficiency and protein expression (Tables 4 and 5).

Table 2. Selection Tier Analysis for the Protein Purification Step

selection tier	total	purified >10mg	% purified >10mg	purified between 0 and 10 mg	% purified between 0 and 10 mg	failed purification	% failed purification
1	26	14	53.8%	10	38.5%	2	7.7%
2	181	78	43.1%	20	11.0%	83	45.9%
3	189	87	45.8%	31	16.4%	71	37.4%
4	112	23	20.5%	9	8.0%	79	70.5%
7	56	19	33.9%	9	16.1%	27	48.2%
8	23	5	21.7%	2	8.7%	16	69.6%
9	103	41	39.8%	14	13.6%	48	46.6%
total	729	247	33.9%	90	12.3%	392	53.8%

Table 3. Protein Purification Data

Cys score	total	purified >10mg	% purified >10mg	purified between 0 and 10 mg	% purified between 0 and 10 mg	failed purification	% failed purification
0	80	28	35.0%	11	13.8%	41	51.2%
1	90	33	36.7%	11	12.2%	46	51.1%
2	98	36	36.7%	17	17.3%	45	45.9%
3	97	38	39.1%	11	11.3%	50	51.5%
4	91	33	36.3%	5	5.5%	53	58.1%
5	105	48	45.7%	27	25.7%	32	30.5%
6	93	30	32.3%	12	12.9%	52	55.9%
7	13	4	30.8%	1	7.7%	8	61.5%
8	20	4	20.0%	1	5.0%	15	75.0%
9	103	41	39.8%	14	13.6%	48	46.6%
total	729	247	33.9%	90	12.3%	392	53.8%

Table 4. Cloning Data

# of transmembrane segments	total	cloning success	cloning failure	% cloning success	% cloning failure
0	2188	1028	1160	47.0%	53.0%
1	119	43	76	36.1%	63.9%
2	157	59	98	37.6%	62.4%
3	109	39	70	35.8%	64.2%
4	13	4	9	30.8%	69.2%
5	11	3	8	27.3%	72.7%
6	8	4	4	50.0%	50.0%
7	4	1	3	25.0%	75.0%
8	5	1	4	20.0%	80.0%
9	53	15	38	28.3%	71.7%
total	2291	1028	1263	44.9%	55.1%

Table 5. Protein Expression Data

# of transmembrane segments	total	expression success	expression failure	% expression success	% expression failure
0	81	37	44	45.7%	54.3%
1	274	109	165	39.8%	60.2%
2	14	9	5	64.3%	35.7%
3	2	1	1	50.0%	50.0%
4	1	0	1	0.0%	100.0%
5	1	0	1	0.0%	100.0%
6	1	0	1	0.0%	100.0%
7	1	0	1	0.0%	100.0%
8	1	0	1	0.0%	100.0%
9	201	147	54	73.1%	26.9%
total	2291	1472	819	64.3%	35.7%

Table 6. Disorder Predictions and HSQC Results

disorder prediction criteria	% proteins predicted to be ordered	% proteins predicted to be disordered	% predictions that proved to be correct by HSQC	% folded by HSQC
Total	71	56%	44%	84%
None	42	40%	60%	84%
Weighted	29	79%	21%	48%

### Protein disorder predicts HSQC results

In solution, intrinsically disordered proteins lack a stable three-dimensional conformation or contain a combination of well-structured (i.e., ordered) and disordered regions. Conformational heterogeneity makes intrinsically disordered proteins inherently antagonistic to high throughput solution-based structure determination (for example, 44 of 71 proteins examined were found to be disordered and thus not amenable to NMR-based structure determination). The avoidance of sequences with such properties should increase the efficiency of efforts aimed at high-throughput structure determination. CESG uses the Predictor of Naturally Disordered Regions

(PONDR<sup>®</sup> - developed by Dr. A.K. Dunker and Molecular Kinetics, Inc.) VLXT algorithm (a collection of neural networks derived from known ordered and disordered residues characterized by X-ray crystallography, NMR, or circular dichroism) to evaluate targets as part of its target scoring scheme. A detailed investigation of the correlation of PONDR analysis with experimental HSQC data has been undertaken (see Oldfield et al., Proteins, in press).

This study (Table 6) looked at a total of 71 proteins experimentally determined to be either ordered (and thus suitable for structure determination) or disordered (and thus unsuitable for continued experimentation) by HSQC NMR spectroscopy. Of those 71 targets, 42 were selected before PONDR disorder predictions were incorporated into the CESG scoring scheme and 29 were selected after disorder predictions began being used.

Table 6. Disorder Predictions and HSQC Results

disorder prediction criteria	% proteins predicted to be ordered	% proteins predicted to be disordered	% predictions that proved to be correct by HSQC	% folded by HSQC
Total	71	56%	44%	84%
None	42	40%	60%	84%
Weighted	29	79%	21%	48%

CESG has found a strong correlation between proteins predicted to be disordered and those found to be experimentally disordered (84% success rate). Incorporation of this structural genomics center's (SG) differences has resulted in an increase in the fraction of targets found to be folded (36% prior to incorporation versus 41% after incorporation). While the current target scoring scheme successfully uses a disorder prediction to weight the overall tier score, our data suggests that we may be able to set a hard threshold for disorder score, eliminating a significant fraction of targets, most of which are unlikely to yield samples suitable for NMR structure determination. For this data set, such a filtering of targets would have resulted in a 12% increase in success.

### Comparing structures solved by structural genomics groups vs. PDB

An analysis of the structures deposited in PDB by structural genomics group reveals several interesting trends. Figure 2 shows the percentage of proteins containing various degrees of low complexity regions in eukaryotic proteins solved by X-ray crystallography from the entire Protein Data Bank (PDB) and from just the structural genomics centers (SG). The differences between the PDB and SG for the various categories indicate that the proteins solved by structural genomics centers contain more regions of low complexity. This shows that the structural genomics centers are solving structures that were previously under-represented in the PDB in terms of presence of low-complexity (see Hunter and Godling, Proteins, 48:134 (2002)). Thus, the SG centers seem to be fulfilling an important role in structural biology.

When the PDB entries are categorized according to Figure 3, there is a clear difference between the structures in the PDB and the structures solved by the structural genomics groups (SG). The majority of the structural genomics groups have clearly been focusing on prokaryotic proteins (~78% of all SG structures have come from prokaryotes). However, the prokaryotes account for less than 40% of all structures in the PDB. Additional focus on eukaryotic structures must be a priority for the future. The vast majority of CESG structures are eukaryotic and provide a necessary balance with the focus of the structural genomics initiatives. The data in this graph comes from all structures solved by both X-ray crystallography and NMR. Archaea and viruses were included as additional categories and factor into the displayed percentages. However, they are not shown on the graph because the percentages were small compared to the prokaryotes and eukaryotes.

