

## Supplementary Information for

Extreme stability in de novo designed repeat arrays is determined by unusually stable short-range interactions

Kathryn Geiger-Schuller<sup>1,2\*</sup>, Kevin Sforza<sup>1\*</sup>, Max Yuhas<sup>1,3</sup>, Fabio Parmeggiani<sup>4,5</sup>, David Baker<sup>4</sup>, & Doug Barrick<sup>1§</sup>

<sup>1</sup>Department of Biophysics and Program in Molecular Biophysics, Johns Hopkins University, 3400 N. Charles St, Baltimore MD 21218.

<sup>2</sup>Current address: Broad Institute of Harvard and Massachusetts Institute of Technology, Cambridge, MA 02142

<sup>3</sup>Current address: Applied Mathematics, Yale University, New Haven CN 06520

<sup>4</sup>Institute for Protein Design, Department of Biochemistry and Howard Hughes Medical Institute, University of Washington, Seattle, Washington 98195

<sup>5</sup>Current address: School of Chemistry and School of Biochemistry, University of Bristol. Life Sciences Building, 24 Tyndall Avenue, Bristol BS8 TQ1

\*These authors contributed equally to this work.

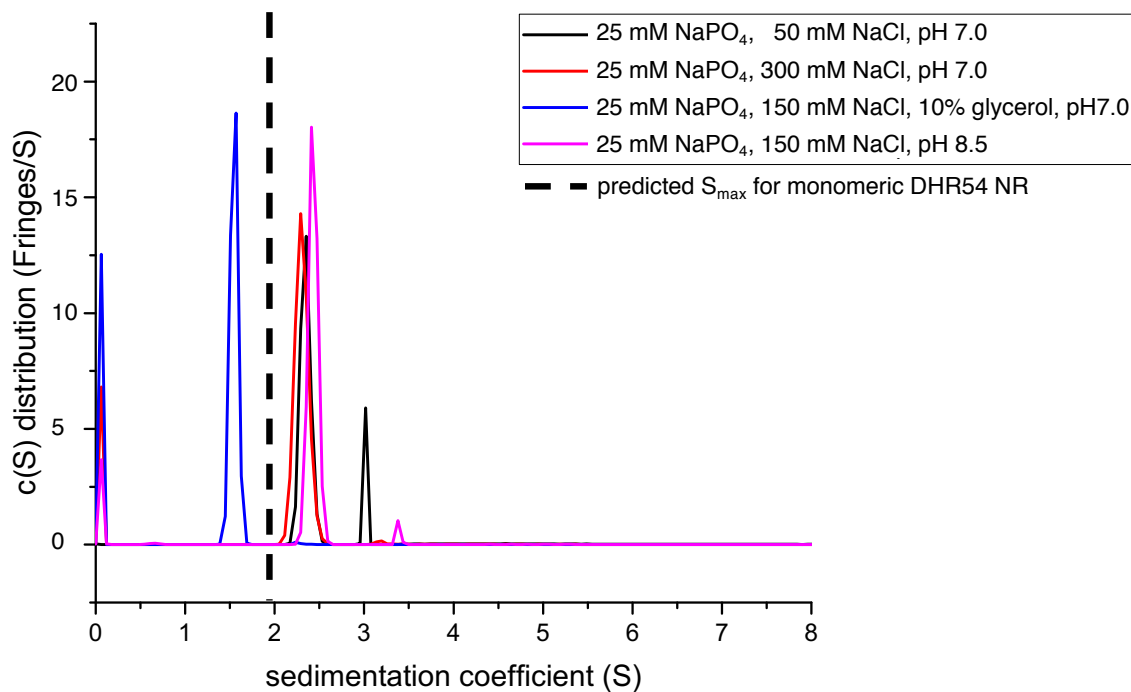
§Corresponding Author: barrick@jhu.edu, (410) 516-0409

### **This PDF file includes:**

Fig. S1

Tables S1 to S2

References for SI reference citations



**Fig. S1. Sedimentation Velocity  $c(S)$  plot for DHR54 NR in the absence and presence of glycerol.** Data were processed and fitted in Sedfit<sup>1</sup> as previously described<sup>2</sup>. The predicted  $S_{max}$  (dotted vertical line) was calculated for DHR54 NR using Sednterp<sup>3</sup>. In the presence of 10 % glycerol, the  $c(S)$  distributions are consistent with monomers.

**Table S1. Sequence features of designed helical repeat proteins used in this study.**

Construct	Sequence (NR <sub>2</sub> C)	n <sub>rep</sub> (n <sub>tot</sub> )	NCPR	Z	m <sub>Ising</sub> (m <sub>2s</sub> )	m <sub>pred</sub>
DHR10.2	SSE <b>KEELRE</b> RL <b>KEV</b> RENA <b>KRKGDDTEEARE</b> AAREAF <b>ERVREAAER</b> AGID SSEVLELAIRLIKEVVENAQREGYDISEAARAAAEAFKRVAEAAKRAGIT SSEVLELAIRLIKEVVENAQREGYDISEAARAAAEAFKRVAEAAKRAGIT SSE <b>TLKRAIEEIRKRV</b> E <b>EAQREGNDISEAARQA</b> AE <b>FRKKAELKRRGDG</b>	50 (200)	0.47	-0.02	4.92	4.75
DHR54	TTEDE <b>RR</b> ELEK <b>V</b> ARKAIEAAAREGNTDEVREQLQRALEIA <b>RESG</b> TTEAVKLAL <b>EV</b> VARVAIEAARRGNTDAVREALEVALEIA <b>RESG</b> TTEAVKLAL <b>EV</b> VARVAIEAARRGNTDAVREALEVALEIA <b>RESG</b> TTEAVRLAL <b>EV</b> VKRVSDEAKKQGNEDAVKEAE <b>EV</b> RKKIE <b>ESG</b>	43 (172)	0.41	-0.06	4.19	4.19
DHR71	DPEEIL <b>ER</b> AK <b>ES</b> LERAREASER <b>GD</b> EE <b>FR</b> KA <b>E</b> KALELAKRL <b>VE</b> QAK <b>KEG</b> DPELVLEAAKVALRVAELAAKNGDKEVFKKAAESALEVAKRL <b>VE</b> VAS <b>KEG</b> DPELVLEAAKVALRVAELAAKNGDKEVFKKAAESALEVAKRL <b>VE</b> VAS <b>KEG</b> DPELV <b>EE</b> AAKVAEEV <b>R</b> KLAKKQGD <b>EE</b> VY <b>E</b> KARETAREV <b>KE</b> ELK <b>R</b> V <b>EE</b> KG	50 (200)	0.475	0.045	4.75	4.75
DHR79	SSDEEEARELIERAKEAA <b>ERA</b> Q <b>EA</b> AERTGDPRVRELAREL <b>KRLA</b> Q <b>EA</b> AA <b>EE</b> V <b>KR</b> DPS SSDVNEALKLIVEAIEAAVRALEAA <b>ERT</b> GDPEVRELAREL <b>LVRLA</b> VEAA <b>EE</b> V <b>QR</b> NPS SSDVNEALKLIVEAIEAAVRALEAA <b>ERT</b> GDPEVRELAREL <b>LVRLA</b> VEAA <b>EE</b> V <b>QR</b> NPS SEEVNEALKKIVKAIQ <b>E</b> AVESL <b>REA</b> EE <b>SG</b> DPEKREKAREV <b>RE</b> AVE <b>RA</b> EEV <b>QR</b> DPS	56 (224)	0.44	-0.10	5.24	5.24

n<sub>rep</sub>: number of residues per repeat. n<sub>tot</sub>: number of residues in NR<sub>2</sub>C construct. NCPR: the net fractional charge per residue, (n<sub>K</sub>+n<sub>R</sub>+n<sub>D</sub>+n<sub>E</sub>)/n<sub>tot</sub>. Z: the total fractional charge, (n<sub>K</sub>+n<sub>R</sub>-n<sub>D</sub>-n<sub>E</sub>)/n<sub>tot</sub>. m<sub>Ising</sub>: experimental m-values from the fitted Ising parameters for NR<sub>2</sub>C (Table 2; 3xm<sub>Gdn, i</sub> + m<sub>Gdn, C</sub> for DHR71, 4xm<sub>Gdn, i</sub> for the other three constructs. m<sub>2s</sub>: experimental m-values from fitting a two-state model to the guanidine HCl-induced unfolding of NR<sub>2</sub>C constructs (Figure 2; units of kcal mol<sup>-1</sup> M<sup>-1</sup>). m<sub>calc</sub>: m-values estimated from the empirical correlation between m-values for guanidine HCl-induced folding and chain-length (Myers Pace Scholtz<sup>4</sup>; units of kcal mol<sup>-1</sup> M<sup>-1</sup>).

**Table S2. Structural features of designed helical repeat proteins used in this study.**

Construct	Contacts within repeats <sup>a</sup>	Contacts between repeats <sup>a</sup>	SASA per naked repeat <sup>b</sup>	SASA between repeats <sup>b</sup>	Twist <sup>c</sup> (radians)	Rise <sup>c</sup> (Å)	Radius <sup>c</sup> (Å)	$\Delta G_i$	$\Delta G_{i+1}$
DHR10.2	56.75	143.7	4,340	2,341	0.03	9.22	138.2	-2.51	-4.80
DHR54	46.0	92.3	3,654	1,874	-0.26	8.76	15.1	-2.04	-6.76
DHR71	66.5	108.2	4,090	1,904	0.23	5.12	49.02	-1.41	-9.93
DHR79	50.5	120.0	4,641	2,311	0.42	6.86	22.51	-3.48	-4.83

<sup>a</sup>Contacts were counted between non-hydrogen pairs closer than 4.2 Å, at a sequence separation of five or more residues from pdb files 5cwg, 5cwl, 5cwn, and 5cwp. For regions where multiple conformations were modeled, the major conformer was used in the contact calculation. Contacts within repeats were averaged over the four (NR<sub>2</sub>C) repeats. Contacts between adjacent repeats were averaged over the three (NR, RR, RC) interfaces.

<sup>b</sup>Solvent accessible surface areas (SASA) were calculated with the `get_area` command of `pymol` with the `dot_solvent` option set to one, using a solvent radius of 1.4 Å. To resolve SASA into naked-repeat (i.e., SASA for a single repeat fragment with no neighboring repeats) from SASA buried between repeats, SASA values were determined for one, two, and three repeat fragments, and for the full four-repeat constructs. For these ten values, SASA was plotted versus repeat number and fitted with a line. The intercept of the best-fit line represents the average SASA buried per interface (with contributions from both repeats forming the interface). The slope of the best-fit line represents the SASA per naked repeat minus the SASA per interface, and is rearranged to determine the former. This procedure gives very similar values to those obtained by summing the SASA of adjacent naked repeats and subtracting the SASA of the corresponding two-repeat fragment. SASA per internal repeat was averaged over the two internal (R) repeats, where each repeat was represented as a `pymol` "selection" within the four-repeat parent construct (rather than a separate "object").

<sup>c</sup>The geometric parameters twist, rise, and radius describe the superhelical geometry of each DHR, and are from Brunette *et al.*<sup>15</sup>.

## References

1. Schuck, P. Size-distribution analysis of macromolecules by sedimentation velocity ultracentrifugation and lamm equation modeling. *Biophys. J.* **78**, 1606–1619 (2000).
2. Marold, J. D., Kavran, J. M., Bowman, G. D. & Barrick, D. A Naturally Occurring Repeat Protein with High Internal Sequence Identity Defines a New Class of TPR-like Proteins. *Struct. Lond. Engl. 1993* (2015). doi:10.1016/j.str.2015.07.022
3. John Philo's Software Home Page. Available at: <http://www.jphilo.mailway.com/>. (Accessed: 6th October 2017)
4. Myers, J. K., Pace, C. N. & Scholtz, J. M. Denaturant m values and heat capacity changes: relation to changes in accessible surface areas of protein unfolding. *Protein Sci. Publ. Protein Soc.* **4**, 2138–2148 (1995).