

**Supporting Information for:
“Sequence determinants of protein phase behavior from a
coarse-grained model”**

Gregory L. Dignon¹, Wenwei Zheng², Young C. Kim³, Robert B. Best¹, and Jeetain Mittal^{2a}

¹Department of Chemical and Biomolecular Engineering, Lehigh University, Bethlehem, Pennsylvania 18015, United States

²National Institute of Diabetes and Digestive and Kidney Diseases,
National Institutes of Health, Bethesda, 20892, United States

³Center for Materials Physics and Technology, Naval Research Laboratory, Washington, DC 20375, United States

^aElectronic mail: jeetain@lehigh.edu

1. SUPPLEMENTARY METHODS

Modelling of LAF-1 helicase domain.

Since there is no solved structure for the helicase domain of LAF-1, we started by predicting its structure from its homologue VASA. Both LAF-1 and VASA belong to the DEAD-Box family and the structure of *Drosophila* VASA has been solved [1]. We first aligned the LAF-1 sequence to the structured part of VASA using the MUSCLE v3.8 web service [2]. The sequence similarity is 51% and the alignment of the structured part is shown below:

```

VASA      Residue 202-621
LAF-1     Residue 187-623
VASA      YIPPEPSNDAIEI-FSSGIASGIHFHFSKYNNIPVKVTGSDVQPPIQHFTSADLRDIIIDNV
LAF-1     WENRGARDERIEQELFSGQLSGINFDKYEEIPVEATGDDVQPPI SLFSDLSLHEWIEENI
          :   .  :: **  : **  *****:*****:*****. *:. .:: * :*:
VASA      NKSGYKIPTPIQKCSIPVISSGRDLMACAQTGSGKTA AFLLPILSKLLED-PHELEL---
LAF-1     KTAGYDRPTPVQKYSIPALQGGRDLMSCAQTGSGKTA AFLVPLVNAILQDGPDAVHRSVT
          :.::*. **:* ** .:*****:*****:*****:*****. :*: *  :
VASA      ---GR---PQVVIVSPTRELAIQIFNEARKFAFESYLKIGIVYGG-TSFRHQNECITRG
LAF-1     SSGGRKKQYPSALVLSPTRELSLQIFNESRKFAYRTPITSALLYGGRENYKDQIHKLRLLG
          **   * .:*****:*****:*****: : . :***** .:. *  : *
VASA      CHVVIATPGRLLDVDRTFITFEDTRFVVLDEADRMLDMGFSEDMRRIM--THVTMRPEH
LAF-1     CHILIIATPGRLIDVMDQGLIGMEGCRYLVLDEADRMLDMGFEPQIRQIVECNRMPSKEER
          **:*****:*.:. . * :. :*****:*****. :*: *  .:. . *
VASA      QTLMSATFPPEEIORMAGEFLK-NYVFVAIGIVGGACSDVKQTIYEVNKYAKRSKLIIEIL
LAF-1     ITAMFSATFPKEIQLLAQDFLKENYVFLAVGRV GSTSENIMQKIVVVEEDEKRSYLMDDL
          * *****:*** : * :*** *****:* * * .: .: * * : * : * : *
VASA      SEQADG--TIVFVETKRGADFLASFLSEKEFP TTSIHGDRLQSQRQALRDFKNGSMKVL
LAF-1     DATGDSSTLTVFVETKRGASDLAYYLN RQNYEVVTI HGDLKQFEREKHLDLFRTGTAPIL
          .  *. :*****. ** :. :. :***** * :*: *  .:. : *
VASA      IATSVASRGLDIKNIKHVINYDMP SKIDDYVHRIGRTGRVGNNGRATSF FDPEKDRAIAA
LAF-1     VATAVAARGLDIPNVKHVINYDLPSDVDEYVHRIGRTGRVGNVGLATSF FN-DKNRNIAR
          :*:*:***** * :*****:*. :*:*****:***** * *****: :*: * *
VASA      DLVKILEGSGQTVPDFLR
LAF-1     ELMDLIVEANQELPDWLE
          :*:*: :. * :*: *

```

We provided the VASA structure (PDB:2DB3) and the sequence alignment information shown above as the inputs to the Modeller software package v9.17 [3]. We then repeated the modelling process 100 times and picked the structure with the smallest energy as best model for the structure of the LAF-1 helicase domain (Figure S7).

Sequences of the proteins used in this work.

```

FUS WT
MASNDYTQQA TQSYGAYPTQ PGQYSQQSS QPYGQSSYSG YSQSTDTSGY QSSYSSYGO
SQNTGYGTQS TPQYGSTGG YGSSQSSQSS YGQSSYPGY GQPAPSSS S YGSSSSQSS

```

SYGQPQSGSY SQQPSYGGQQ QSYGQQQSYN PPQGYGQQNQ YNS

FUS 6E

MASNDYTQQA TQSYGAYPTQ PGQGYEQQSE QPYGQQSYSG YSQSTDTSGY GQSSYSSYGQ
SQNTGYGEQS TPQGYGSTGG YGSEQSEQSS YGQQSSYPGY GQQPAPSSTS GSYGSSEQSS
SYGQPQSGSY SQQPSYGGQQ QSYGQQQSYN PPQGYGQQNQ YNS

FUS 6E'

MASNDYTQQA TQSYGAYPTQ PGQGYEQQSE QPYGQQSYSG YSQSTDTSGY GQSSYSSYGQ
SQNTGYGEQS TPQGYGSTGG YGSEQSEQSS YGQQSSYPGY GQQPAPSSTS GSYGSSEQSS
SYGQPQSGSY SQQPSYGGQQ QSYGQQQSYN PPQGYGQQNQ YNS

FUS 6E*

MASNDYEQQA TQSYGAYPTQ PGQGYEQQSS QPYGQQSYSG YSQSTDTSGY GQSSYSSYGQ
SQNTGYGTQS TPQGYGSTGG YGSEQSEQSS YGQQSSYPGY GQQPAPSSTS GSYGSSEQSS
SYGQPQSGSY EQQPSYGGQQ QSYGQQQSYN PPQGYGQQNQ YNS

FUS 12E

MASNDYEQQA EQSYGAYPEQ PGQGYEQQSE QPYGQQSYSG YEQSTDTSGY GQSSYSSYGQ
EQNTGYGEQS TPQGYGSTGG YGSEQSEQSS YGQQSSYPGY GQQPAPSSTS GSYGSSEQSS
SYGQPQSGSY EQQPSYGGQQ QSYGQQQSYN PPQGYGQQNQ YNS

LAF-1 IDR

MESNQSNNGG SGNAALNRGG RYVPPHLRGG DGGAAAAASA GGDDRRGGAG GGGYRRGGGN
SGGGGGGGYD RGYNDNRDDR DNRGGSGGYG RDRNYEDRGY NGGGGGGGNR GYNNNRGGGG
GGYNRQDRGD GGSSNFSRGG YNNRDEGSDN RGSGRSYNND RRDNGGDG

LAF-1 full length

...QN TRWNNLDAPP
SRGTSKWENR GARDERIEQE LFSGQLSGIN FDKYEEIPVE ATGDDVPQPI SLFSDLSLHE
WIEENIKTAG YDRPTPVQKY SIPALQGGRD LMSCAQTGSG KTA AFLVPLV NAILQDGPDA
VHRSVTSSGG RKKQYPSALV LSP TRELSLQ IFNESRKFAY RTPITSALLY GGRENYKDQI
HKLRLGCHIL IATPGR LIDV MDQGLIGMEG CRYLVLDEAD RMLDMGFEPQ IRQIVECNRM
PSKEERITAM FSATFPKEIQ LLAQDFLKEN YVFLAVGRVG STSENIMQKI VWVEEDEKRS
YLMDLLDATG DSSLTLVFVE TKRGASDLAY YLNRQNYEVV TIHGDLKQFE REKHLDLFRT
GTAPILVATA VAARGLDIPN VKHVINYDLP SDVDEYVHRI GRTGRVGNVG LATSFFNDKN
RNIARELMDL IVEANQELPD WLE

CspIm

GPGMRGKVKW FDSKKGYGFI TKDEGGDVFV HWSAIEMEGF KTLKEGQVVE FEIQEGKKG
QAAHVKV

IN

GSHCFLDGID KAQEEHEKYH SNWRAMASDF NLPPVVAKEI VASCDKCQLK GEAMHGQVDC

Prothymosin alpha-N

GPSDAAVDTS SEITTKDLKE KKEVVEEAEN GRDAPANGNA ENEENGEQEA DNEVDEECE
GEEEEEEEE GDGEEEDGDE DEEAESATGK RAAEDDED DD VDTKKQKTDE DD

S4

Prothymosin alpha-C

MAHHHHHSA ALEVLFGQPM SDAAVDTSS E ITTKDLKEKK EVVEEAENGR DAPANGNANE
ENGEQEADNE VDEECEEGGE EEEEEEEGDG EEEDGDEDEE AESATGKRAA EDEDEDDVDT
KKQKTDEDD

R15

KLKEANKQQN FNTGIKDFDF WLSEVEALLA SEDYGKDLAS VNNLLKKHQL LEADISAHED
RLKDLNSQAD SLMTSSAFDT SQVKDKRETI NGRFQRIKSM AAARRAKLNE SHRL

R17

RLEESLEYQQ FVANVEEEEA WINEKMTLVA SEDYDGLAA IQGLLKKHEA FETDFTVHKD
RVNDVAANGE DLIKKNNHHV ENITAKMKGL KGKVSDEKA

hCyp

SSFHRIIPGF MSQGGDFTRH NGTGGKSIYG EKFEDEFIL KHTGPGILSM ANAGPNTNGS
QFFISTAKTE FLDGKHVVFV KVKEGMNIVE AMERFGSRNG KTSKKITIAD SGQLE

Protein L

MEEVTIKANL IFANGSTQTA EFKGTFEKAT SEAYAYADTL KKDNGEWTVD VADKGYTLNI
KFAQ

ACTR

GTQNRPLLRN SLDDLVGPPS NLEGQSDERA LLDQLHTLLS NTDATGLEEI DRALGIPELV
NQGQALEPKQ D

hNHE1cdt

MVPAHKLDSP TMSRARIGSD PLAYEPKEDL PVITIDPASP QSPESVDLVN EELKGKVLGL
SRDPAKVAEE DEDDDGGIMM RSKETSSPGT DDVFTPAKSD SPSSQRIQRC LSDPGPHPEP
GEGEPFFPKG Q

sNase

ATSTKKLHKE PATLIKAIDG DTVKLMYKQ PMTFRLLLVD TPETKHPKKG VEKYGPEASA
FTKMKVENAK KIEVEFDKQ RTDKYGRGLA YIYADGMVN EALVRQGLAK VAYVYKPNNT
HEQHLRSEA QAKKEK

alpha-synuclein

MDVFMKGLSK AKEGVVAAA E KTKQGVAAE GKTKEGVLYV GSKTKEGVVH GVATVAEKT
EQVTNVTGAV VTGVTAVAQK TVEGAGSIAA ATGFVKKDQL GKNEEGAPQE GILEDMPVDP
DNEAYEMPSE EGYQDYEPEA

Prothymosin alpha

GPSDAAVDTS SEITTKDLKE KKEVVEEAEN GRDAPANGNA ENEENGEQEA DNEVDEECE
GEEEEEEEEE GDGEEEDGDE DEEAESATGK RAAEDDEDDD VDTKKQKTDE DD

2. SUPPLEMENTARY FIGURES

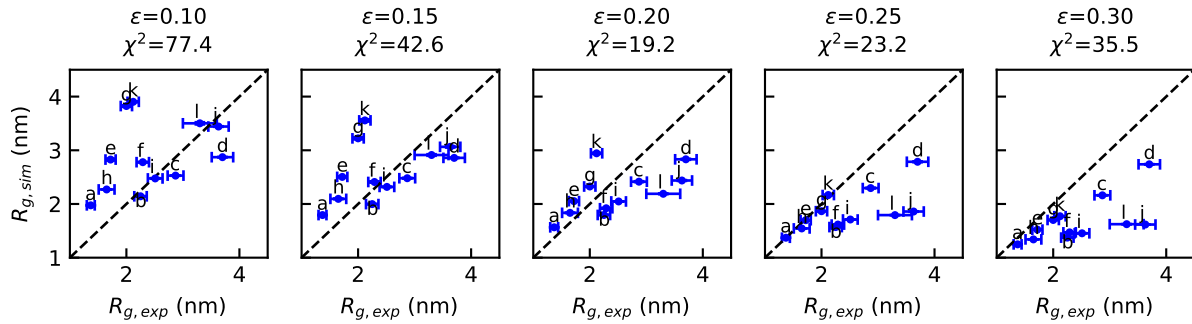


FIGURE S1. Comparison of R_g between simulations and experiments with different ϵ parameters for HPS model. The deviations χ^2 between the simulations and experiments are shown in the title. The list of the proteins and legends can be found in Table S2.

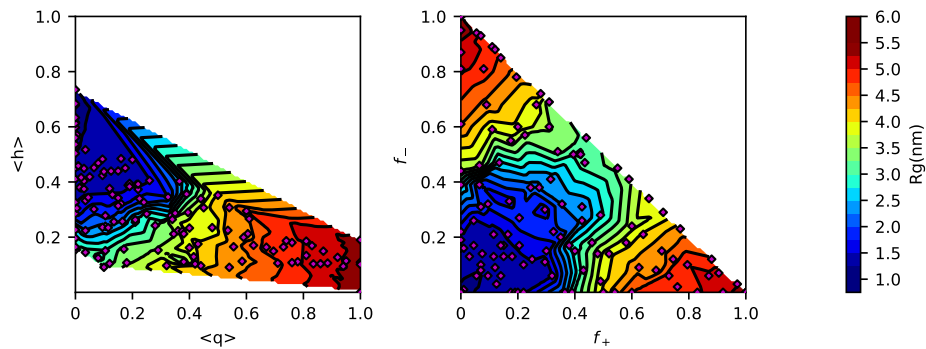


FIGURE S2. Randomly generated 100mers in a Uversky (left) and Pappu (right) plot to show the dependence of R_g on charge and hydrophobicity using the KH model.

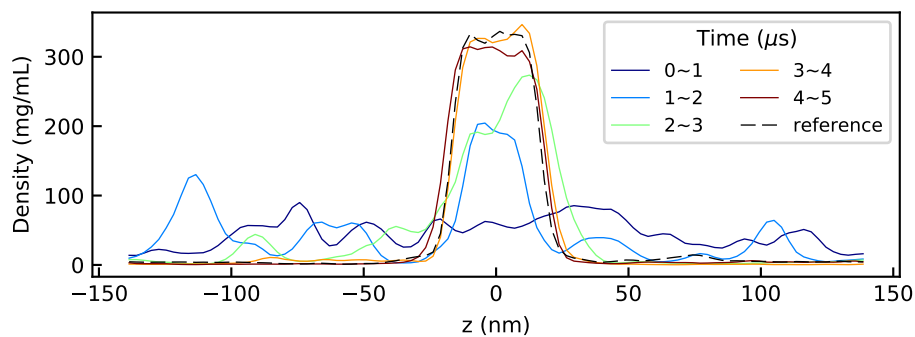


FIGURE S3. LAF-1 simulation started from dispersed state at 210K with KH-D model showing coalescence to a slab conformation after about $4 \mu\text{s}$. The colored lines show the density profile at different time ranges throughout the simulation. The black line shows the simulation starting from an initial slab configuration as a reference.

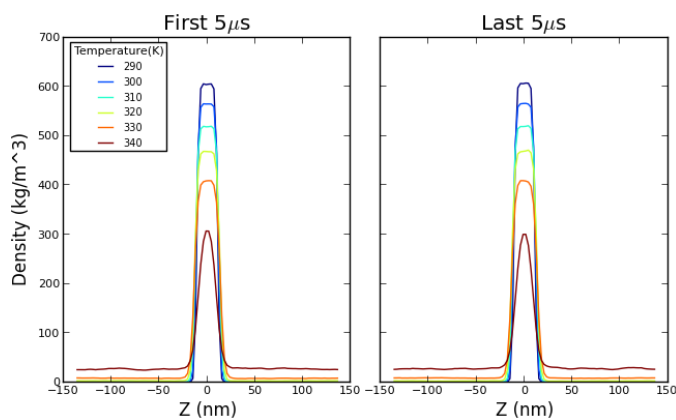


FIGURE S4. Comparison of density profiles between first $5 \mu\text{s}$ and last $5 \mu\text{s}$ of slab simulations of FUS WT.

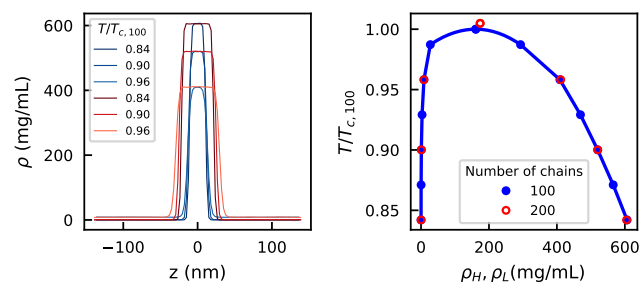


FIGURE S5. Comparison of FUS WT simulations with 100 (blue) and 200 (red) chains. Temperatures are scaled by the critical temperature of the simulations with 100 chains.

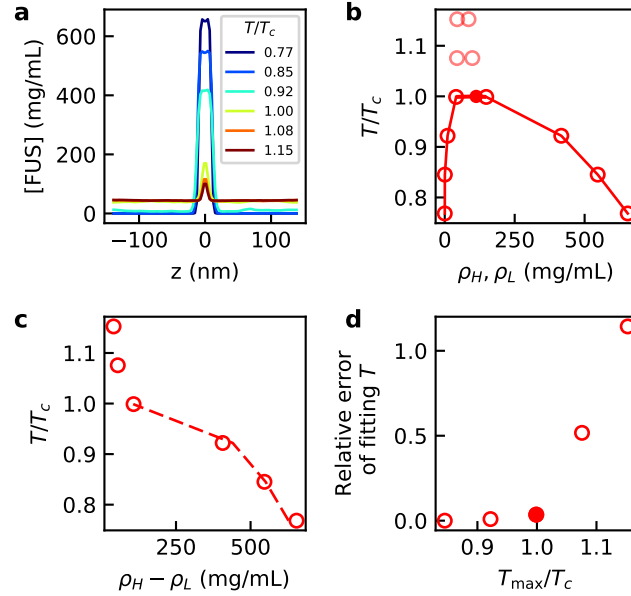


FIGURE S6. Methodology used to determine the range of temperatures to fit Eq. 6 in the main text. Temperatures are scaled by the critical temperature. a) Dependence of WT FUS concentration on the z-axis of the SLAB simulation with KH model. b) The concentrations of low- and high-density phases as a function of temperature. The empty symbols show the data and the solid symbol shows the critical temperature obtained from fitting to all except the two highest temperatures, where the system cannot be described by Eq. 6 any more. c) The difference of concentrations between the low- and high-density phases as a function of the temperature. Dashed line shows the fitting to Eq. 6. d) The relative error of fitting T using Eq. 6 as a function of the maximum temperature for fitting. The maximum temperature used for fitting to obtain the critical temperature is shown in solid symbol.

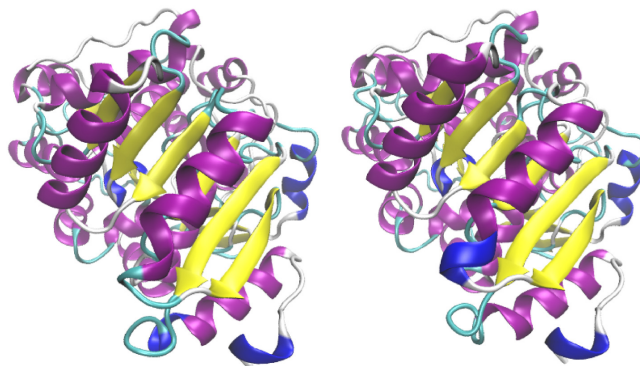


FIGURE S7. Homology modelling of helicase domain of LAF-1 using the structure of VASA. Left: the structure of VASA residue 202-621 (PDB:2DB3[1]); Right: the structure of LAF-1 helicase domain residue 187-623 from homology modelling.

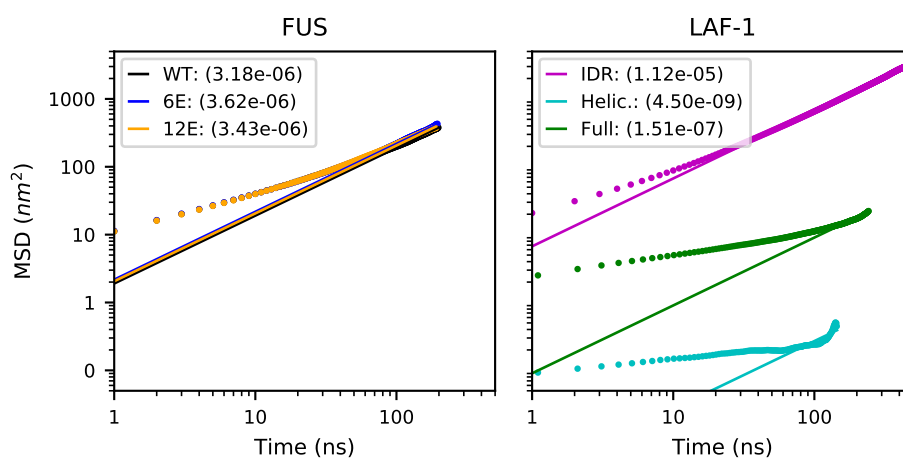


FIGURE S8. Mean squared displacement (MSD) as a function of time for FUS variants at 260K and 600 mg/mL (left), and LAF-1 at 210K and 260, 535 and 500 mg/mL for IDR, helicase, and full length respectively. Linear fits were calculated using all data points after 100 ns. Diffusion coefficients are included in parentheses in units of cm²/s .

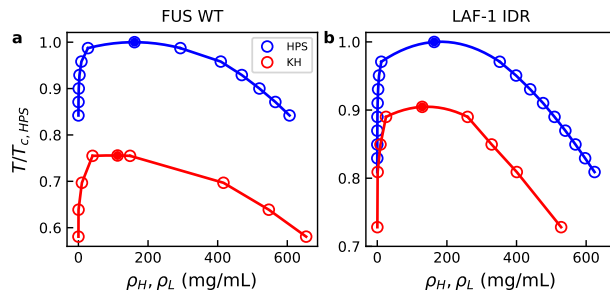


FIGURE S9. Comparison of the phase diagram generated with HPS (blue) and KH (red) model in FUS WT (left) and LAF-1 IDR (right). Temperatures are scaled by critical temperatures using HPS model.

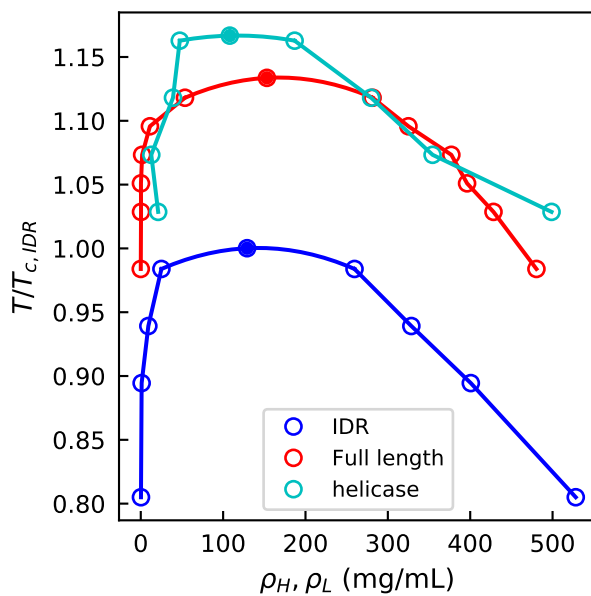


FIGURE S10. Phase diagram of IDR (blue), helicase (cyan) and full length (red) LAF-1. Temperatures are scaled by the critical temperature of IDR LAF-1.

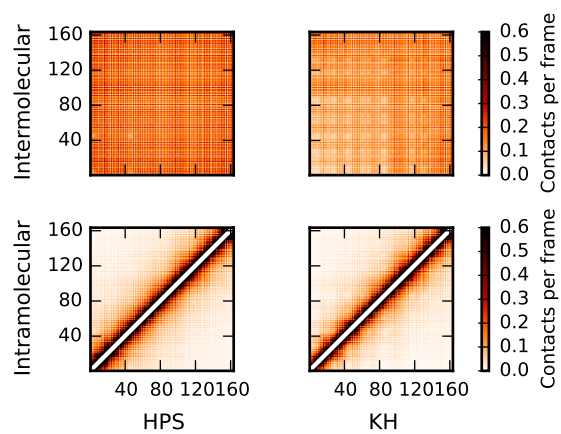


FIGURE S11. Inter- (upper) and intra-molecular (lower) contact maps for FUS WT at 260 K using HPS (left) and KH models (right).

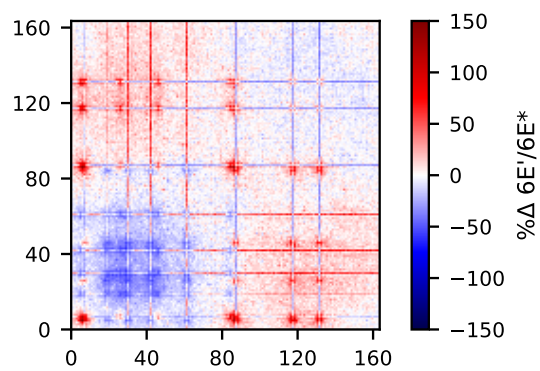


FIGURE S12. Intermolecular contacts for FUS 6E' divided by that of FUS 6E* showing how the overall number of contacts forming within the slab changes between the two sequences.

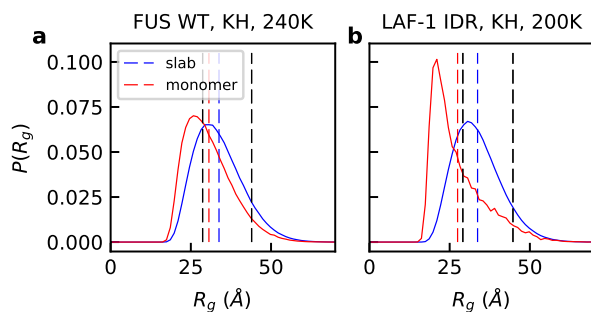


FIGURE S13. Radii of gyration of the disordered proteins inside (blue) and out of (red) the slab. a) FUS WT with KH model at 240K. b) LAF-1 IDR with KH model at 200K. Black lines show the R_g from the random coil or excluded volume chain with the same chain length.

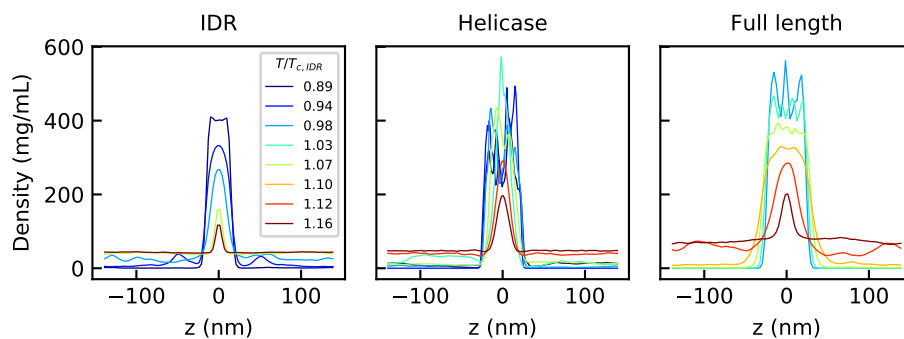


FIGURE S14. Slab density profiles of IDR (left), helicase (middle) or full length (right) LAF-1. Temperatures are scaled by the critical temperature of IDR LAF-1.

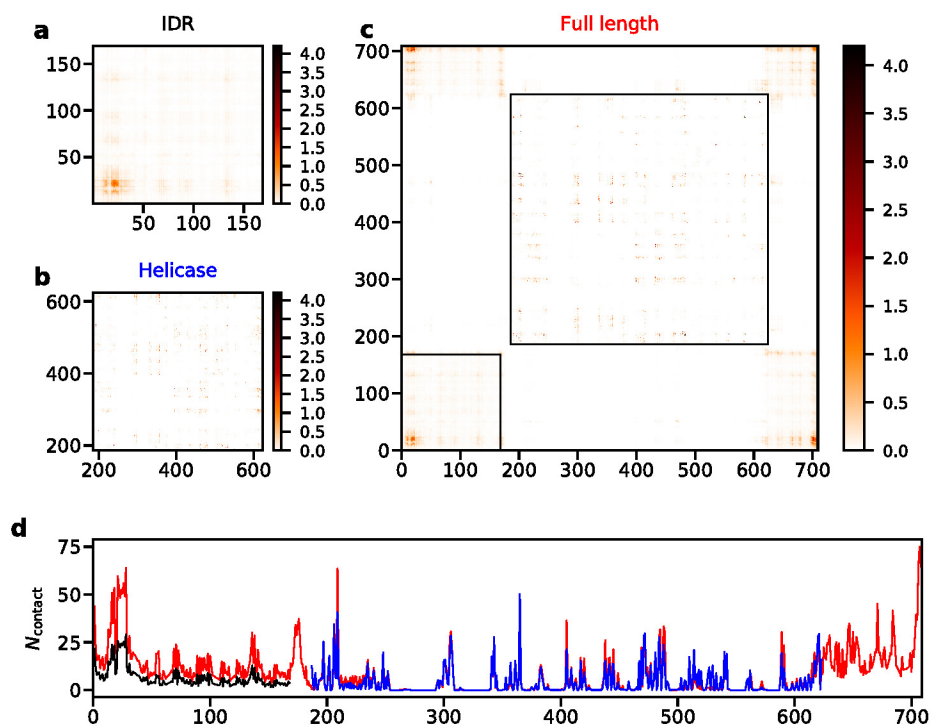


FIGURE S15. Number of intermolecular contacts per frame for different LAF-1 variants at 220K. a) Contact map of IDR LAF-1. b) Contact map of helicase LAF-1. c) Contact map of full length LAF-1. Black boxes illustrate the N-terminal IDR and the helicase domain. d) Number of intermolecular contacts per residue per frame for IDR LAF-1 (black), helicase LAF-1 (blue) and full length LAF-1 (red).

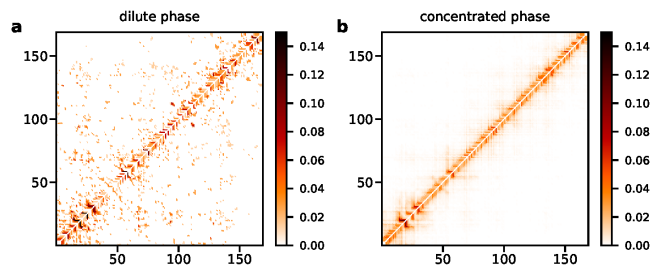


FIGURE S16. Number of intramolecular contacts per chain per frame for LAF1 IDR with KH model at 200K. The contact map for the chains in the dilute phase is shown on the left side and that in the concentrated phase is shown on the right side.

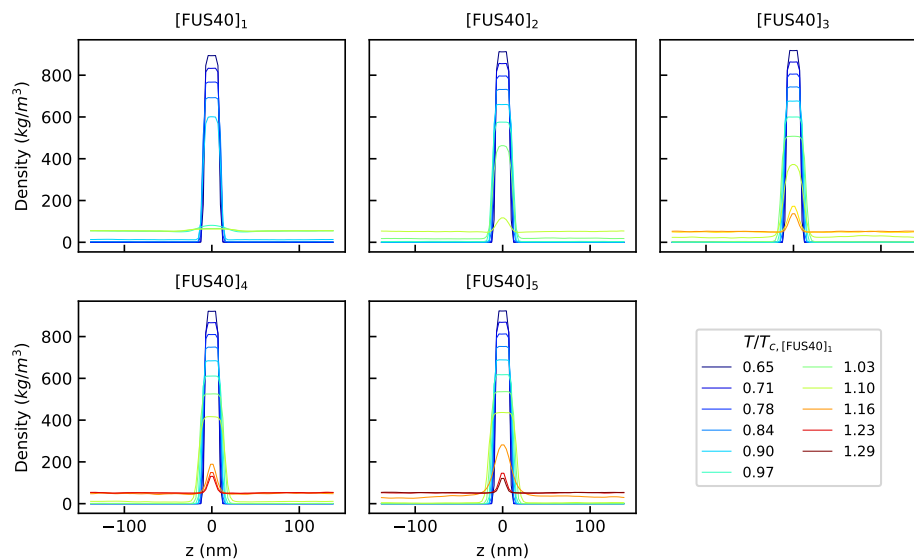


FIGURE S17. Slab density profiles of the repeated peptides of FUS fragment. Temperatures are scaled by the critical temperature of $[FUS40]_1$.

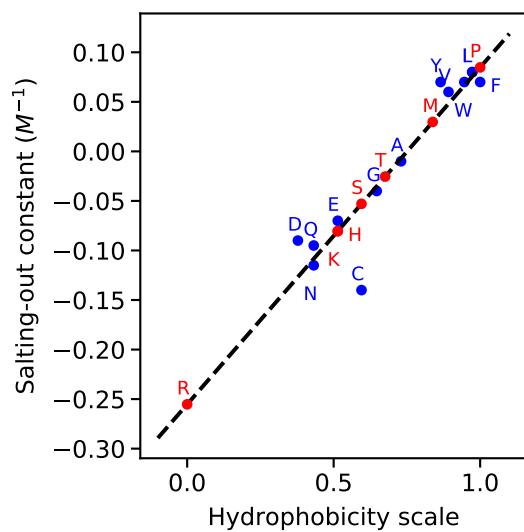


FIGURE S18. The correlation between salting-out constant and hydrophobicity scale. Black line shows the linear fitting curve between these two parameters. Blue dots show the data from literature [4, 5, 6, 7] and red dots show the estimate from linear inter- or extrapolation.

3. SUPPLEMENTARY TABLES

TABLE S1. The amino acid parameters used in the HPS model. σ is the diameter of the amino acid used in the short-ranged pair potential. λ is the scaled hydrophobicity from the literature [8].

Type	Mass (amu)	Charge	σ (Å)	λ
ALA	71.08	0	5.04	0.730
ARG	156.20	1	6.56	0.000
ASN	114.10	0	5.68	0.432
ASP	115.10	-1	5.58	0.378
CYS	103.10	0	5.48	0.595
GLN	128.10	0	6.02	0.514
GLU	129.10	-1	5.92	0.459
GLY	57.05	0	4.50	0.649
HIS	137.10	0.5	6.08	0.514
ILE	113.20	0	6.18	0.973
LEU	113.20	0	6.18	0.973
LYS	128.20	1	6.36	0.514
MET	131.20	0	6.18	0.838
PHE	147.20	0	6.36	1.000
PRO	97.12	0	5.56	1.000
SER	87.08	0	5.18	0.595
THR	101.10	0	5.62	0.676
TRP	186.20	0	6.78	0.946
TYR	163.20	0	6.46	0.865
VAL	99.07	0	5.86	0.892

TABLE S2. List of intrinsically disordered or unfolded proteins with experimentally determined R_g . The radii of gyration of ACTR and hNHE1cdt were measured at 5°C and 45°C in the experiment and have been interpolated to 25°C for comparison with the other proteins.

	Protein	Chain length	[Ion] (mM)	R_g , expt (nm)	Method	P_{charge}	Hydrophobicity
a	CspTm	67 (54)	42	1.37 (0.07)	FRET [9]	0.313	0.676
b	IN	60 (57)	50	2.25 (0.11)	FRET [9]	0.267	0.648
c	ProT α -N	112 (56)	42	2.87 (0.14)	FRET [9]	0.563	0.555
d	ProT α -C	129 (55)	42	3.70 (0.19)	FRET [9]	0.488	0.573
e	R15	114 (94)	128	1.72 (0.09)	FRET [10]	0.325	0.616
f	R17	100 (94)	128	2.29 (0.11)	FRET [10]	0.340	0.647
g	hCyp	167 (164)	85	2.00 (0.05)	FRET [10]	0.234	0.679
h	Protein-L	64 (64)	128	1.65 (0.14)	FRET [11]	0.266	0.682
i	ACTR	71	199	2.51 (0.13)	SAXS [12]	0.254	0.644
j	hNHE1cdt	131	199	3.63 (0.18)	SAXS [12]	0.298	0.671
k	sNase	136	17	2.12 (0.10)	SAXS [13]	0.331	0.659
l	α -synuclein	140	156	3.3 (0.3)	FRET [14]	0.279	0.678

TABLE S3. Summary of slab simulations and critical temperatures obtained.

System	Model	$N_{residue}$	N_{chain}	T_c (K)
FUS				
WT	KH	163	100	260.3
WT	HPS	163	100	344.4
WT	HPS	163	200	346.1
6E mutant	HPS	163	100	326.6
6Ep mutant	HPS	163	100	327.2
6Es mutant	HPS	163	100	326.4
12E mutant	HPS	163	100	280.3
LAF-1				
IDR	KH	168	100	223.6
IDR	HPS	168	100	247.2
Folded	KH	437	100	260.9
Full length	KH	708	100	253.5
Repeated fragment of FUS				
[FUS40] ₁	HPS	40	480	309.6
[FUS40] ₂	HPS	80	240	336.5
[FUS40] ₃	HPS	120	160	348.5
[FUS40] ₄	HPS	160	120	356.1
[FUS40] ₅	HPS	200	96	363.7

REFERENCES

- [1] Sengoku T, Nureki O, Nakamura A, Kobayashi S, Yokoyama S. Structural basis for RNA unwinding by the DEAD-box protein Drosophila Vasa. *Cell*. 2006;125(2):287–300.
- [2] Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*. 2004;32(5):1792–1797.
- [3] Šali A, Blundell TL. Comparative protein modelling by satisfaction of spatial restraints. *J Mol Biol*. 1993;234(3):779–815.
- [4] Schrier EE, Schrier EB. The salting-out behavior of amides and its relation to the denaturation of proteins by salts. *J Phys Chem*. 1967;71(6):1851–1860.
- [5] Nandi PK, Robinson DR. Effects of salts on the free energy of the peptide group. *J Am Chem Soc*. 1972;94(4):1299–1308.
- [6] Nandi PK, Robinson DR. Effects of salts on the free energies of nonpolar groups in model peptides. *J Am Chem Soc*. 1972;94(4):1308–1315.
- [7] Baldwin RL. How Hofmeister ion interactions affect protein stability. *Biophys J*. 1996;71(4):2056–2063.
- [8] Kapcha LH, Rossky PJ. A simple atomic-level hydrophobicity scale reveals protein interfacial structure. *J Mol Biol*. 2014;426(2):484–498.
- [9] Müller-Spāth S, Sorzano A, Hirschfeld V, Hofmann H, Rügger S, Reymond L, et al. Charge interactions can dominate the dimensions of intrinsically disordered proteins. *Proc Natl Acad Sci USA*. 2010;107:14609–14614.
- [10] Hofmann H, Sorzano A, Borgia A, Gast K, Nettels D, Schuler B. Polymer scaling laws of unfolded and intrinsically disordered proteins quantified with single-molecule spectroscopy. *Proc Natl Acad Sci USA*. 2012;109:16155–16160.
- [11] Sherman E, Haran G. Coil-globule transition in the denatured state of a small protein. *Proc Natl Acad Sci USA*. 2006;103:11539–11543.
- [12] Kjaergaard M, Norholm AB, Hendus-Altenburger R, Pedersen SF, Poulsen FM, Kragelund BB. Temperature-dependent structural changes in intrinsically-disordered proteins: formation of α -helices or loss of polyproline II? *Protein Sci*. 2010;19:1555–1564.
- [13] Flanagan JM, Kataoka M, Shortle D, Engelman DM. Truncated staphylococcal nuclease is compact but disordered. *Proc Natl Acad Sci U S A*. 1992;89(2):748–752.
- [14] Nath A, Sammalkorpi M, DeWitt DC, Trexler AJ, Elbaum-Garfinkle S, O'Hern CS, et al. The conformational ensembles of α -synuclein and tau: Combining single-molecule FRET and simulations. *Biophys J*. 2012;103(9):1940–1949.
- [15] Kim YC, Hummer G. Coarse-grained models for simulation of multiprotein complexes: application to ubiquitin binding. *J Mol Biol*. 2008;375:1416–1433.