

SUPPORTING INFORMATION

Stabilizing proteins from sequence statistics: The interplay of conservation and correlation in triosephosphate isomerase stability

Brandon J. Sullivan¹, Tran Nguyen³, Venuka Durani², Deepti Mathur, Samantha Rojas, Miriam Thomas², Trixy Syu³ and Thomas J. Magliery^{2,3*}

¹Ohio State Biochemistry Program, The Ohio State University, Columbus, OH 43210

²Department of Chemistry, The Ohio State University, Columbus, OH 43210

³Department of Biochemistry, The Ohio State University, Columbus, OH 43210

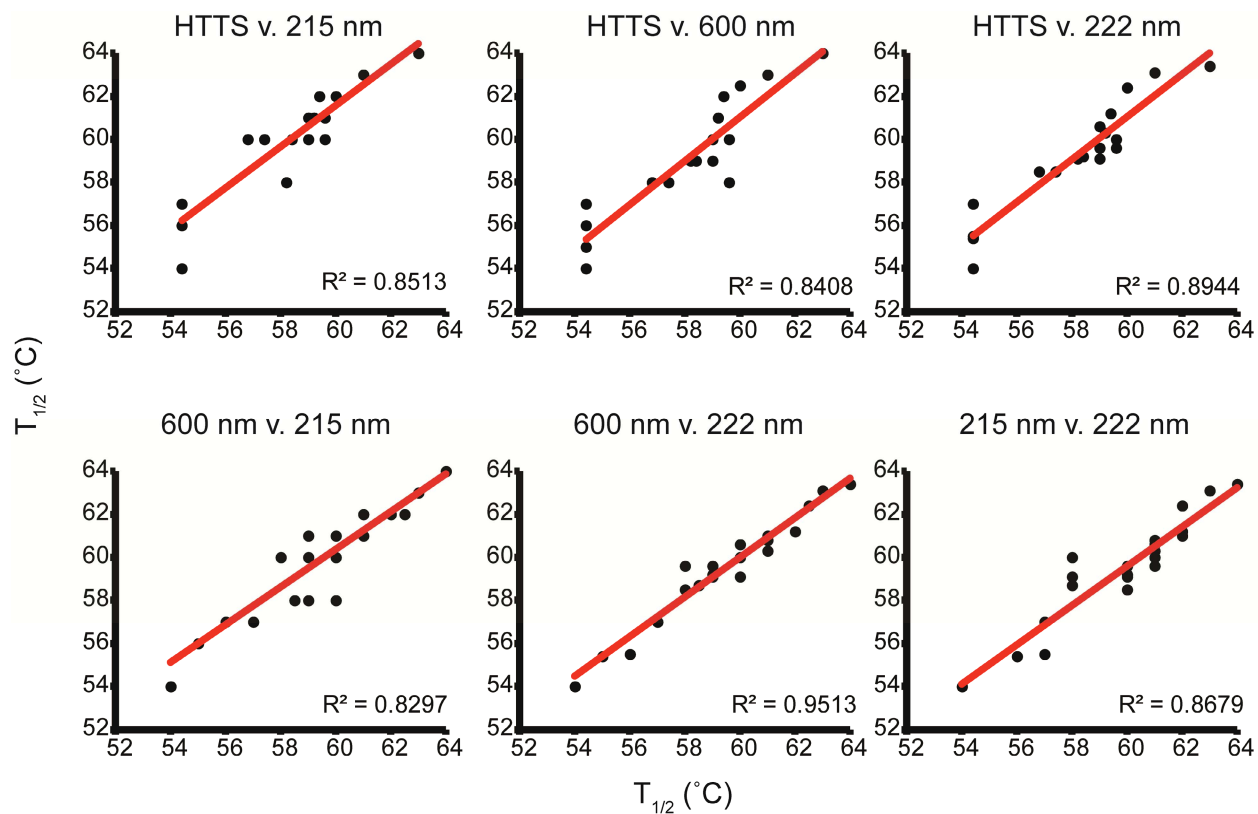
*Corresponding author. E-mail address: magliery.1@osu.edu. Departments of Chemistry and Biochemistry, The Ohio State University, 100 W. 18th Ave., Columbus, OH 43210. Phone (614) 247-8425. Fax (614) 292-1685.

Host protein sequence

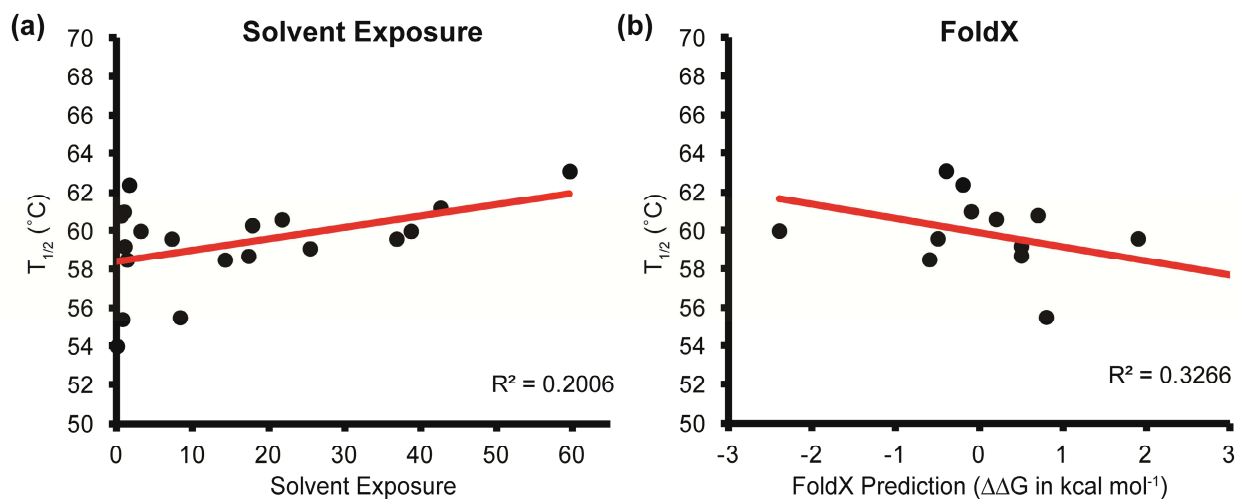
S.c. TIM sequence

2
MAHHHHHHGGENLYFQGSSGARTFFVGGNFKLNGSKQSIKEIVERLNTASIPENVEVVICPPATYLDYSVSLVKKPQ
VTVGAQNAYLKASGAFTGENSVDQIKDVGAKWVILGHSERRSYFHEDDKFIADKTKFALGQGVGVLICIGETLEEKK
AGKTLDVVERQLNAVLEEVKDWTNVVVAYEPVWAI GTGLAATPEDAQDIHASIRKFLASKLGDKAASELRILYGGSA
NGSNAVTFKDKADV D GFLVGGASLKPEFVDI INSRN*

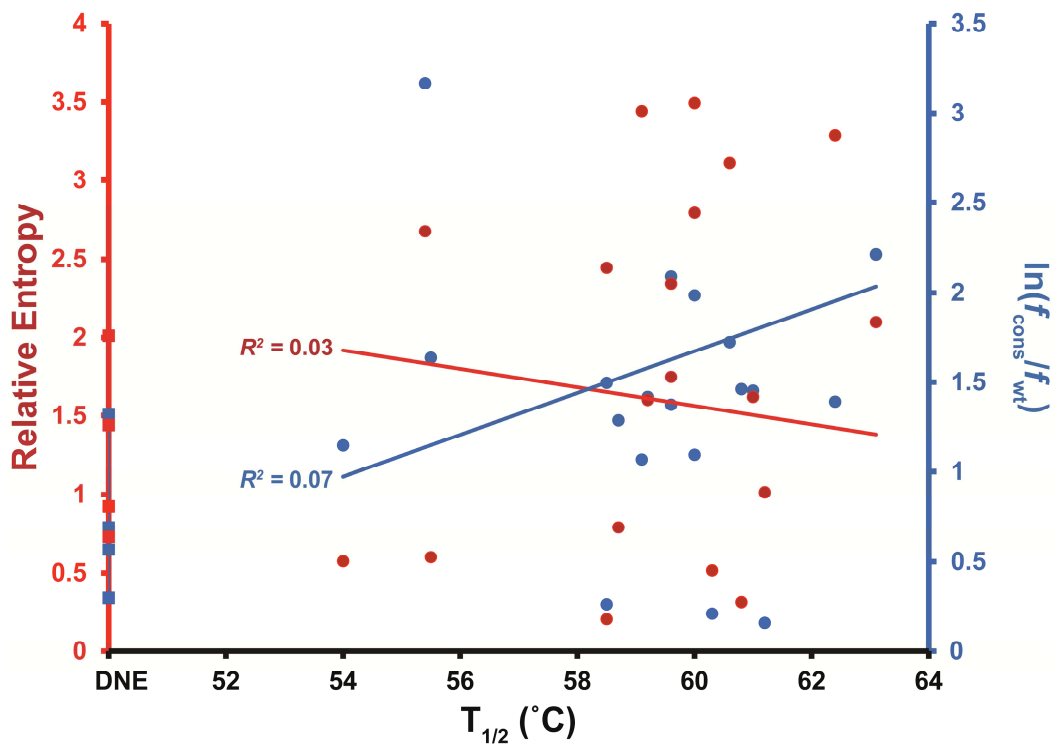
The TEV cleavable sequence is shown underlined. A GSSG linker was cloned downstream of the TEV recognition sequence to improve cleavage efficiency. The native *S.c.* TIM sequence begins MARTFFV..., where the initiating Met is labeled as residue 1. Our characterized variants begin with GSSGARTFFV..., where the second Gly is labeled as residue 1. Both sequences are numbered such that Ala at position 2 is consistent.



Supplemental Figure 1. The $T_{1/2}$ values are compared between four complementary methods: (1) Loss of CD ellipticity at 222 nm for α -helices, (2) Loss of CD ellipticity at 215 nm for β -strands, (3) Diffraction of light at 600 nm for detection of precipitation products and (4) High-Throughput Thermal Scanning. Here, the $T_{1/2}$ values are plotted for comparison.



Supplemental Figure 2. (a) The percent solvent exposure for each mutation is plotted against the $T_{1/2}$. (b) The computationally predicted $\Delta\Delta G$ from FoldX is plotted against the $T_{1/2}$.



Supplemental Figure 3. The relative entropy and $\ln(f_{\text{cons}}/f_{\text{wt}})$ are plotted as a function of the $T_{1/2}$. Neither metric is independently effective at deciphering stabilizing versus destabilizing mutations. The data points shown as squares did not express (DNE). These values were not considered when calculating the R^2 values.

(a) Positions > Noise	(b) Maximum MI	(c) Average MI	
V162I	0	V162I	0.09
V121L	0	V121L	0.09
A212V	2	I83L	0.10
I83L	2	A212V	0.11
N78I	2	N78I	0.11
I109V	4	A212V	0.11
K135E	8	I109V	0.12
I40V	11	I40V	0.12
I184V	12	F11W	0.14
F11W	18	I184V	0.14
I127V	28	I127V	0.15
V226I	36	K135E	0.15
A66C	60	I127V	0.15
L13M	71	V226I	0.16
V123P	79	A66C	0.17
Q82M	91	L13M	0.19
K134R	97	Q82M	0.20
C41A	100	K134R	0.20
D180Q	146	V123P	0.20
W90Y	172	C41A	0.21
		D180Q	0.24
		W90Y	0.29

Supplemental Figure 4. Sequence correlation details. Of the 103 possible consensus mutations to *S.c.* TIM only 20 have relative entropies greater than 1.42. (F11W is included here even though its relative entropy is greater than 3.) These 20 mutations are shown here arranged by several measures of correlation. We have individually characterized the effects of stability for 14 of these mutations and that data is shown as red and green for less and more stable, respectively. Mutual information scores for the TIM database range from 0 to 0.83, but MI values less than 0.23 are essentially zero (i.e., noise) based on the scrambled MSA control. (a) Many positions are correlated to each position of mutation above the noise threshold. The strongest site-to-site correlation (b) and the average correlation (c) value for all 240 positions resulting in similar rank-order lists. Note that W90Y, V123P, D18Q and C41A are at or near the bottom of all three lists (i.e., are most correlated).