

**Structure, Volume 22**

**Supplemental Information**

**Random Single Amino Acid Deletion Sampling Unveils**

**Structural Tolerance and the Benefits of Helical**

**Registry Shift on GFP Folding and Structure**

**James A. J. Arpino, Samuel C. Reddington, Lisa M. Halliwell, Pierre J. Rizkallah, and D. Dafydd Jones**

## **Random single amino acid deletion sampling unveils structural tolerance and the benefits of helical registry shift on GFP folding and structure.**

James A. J. Arpino<sup>1</sup>, Sam C. Reddington<sup>1</sup>, Lisa M. Halliwell<sup>1</sup>, Pierre J. Rizkallah<sup>2</sup> & D. Dafydd Jones.<sup>1</sup>

### **Supporting Information.**

#### **Supporting Methods.**

##### ***EGFP TND library construction.***

Insertion of the engineered transposon MuDel into the *egfp* gene encoding enhanced green fluorescent protein (EGFP) residing within the pNOM-XP3 plasmid was performed using an *in vitro* transposition and selection procedure described previously (Baldwin et al., 2009) to generate the library *egfp* $\Delta^{2504}$ . *MlyI* restriction digestion was performed on *egfp* $\Delta^{2504}$  DNA (3  $\mu$ g) to remove MuDel from the pooled plasmid library and analysed by 1.0% (w/v) agarose gel electrophoresis. The linear library DNA was purified from the agarose gel using a QIAquick<sup>®</sup> gel purification kit (QIAGEN). The purified linear library DNA (50 ng) was recircularised by intramolecular ligation with Quick T4 DNA ligase and the reaction cleaned up with a MinElute reaction cleanup kit (QIAGEN). The ligation reaction mixture (1  $\mu$ l) was used to transform electrocompetent *E. coli* BL21-Gold (DE3) cells. The transformed cells were grown on LB agar plates supplemented with 100  $\mu$ g/ml ampicillin and 150  $\mu$ M IPTG and incubated at 37°C overnight then stored at 4°C. Colonies presenting a green colour phenotype upon illumination on a UV transilluminator and colonies with no colour phenotype were selected for a colony PCR screen with primers pEXP-F and DDJ013. The PCR products produced (2  $\mu$ l) were analysed by agarose gel electrophoresis and the rest (23  $\mu$ l) purified using a QIAquick PCR purification kit (QIAGEN) for DNA sequence analysis, to identify the nature of the triplet nucleotide deletions.

##### ***Protein production and purification***

The production and subsequent purification of EGFP and EGFP<sup>G4 $\Delta$</sup>  was performed as follows. LB Broth (15 ml) supplemented with 100  $\mu$ g/ml ampicillin was inoculated with a single *E. coli* BL21-Gold (DE3) colony containing a relevant plasmid (pNOM-XP3 (Baldwin et al., 2009) containing the *egfp* or *egfp*<sup>G4 $\Delta$</sup>  gene) to generate a starter culture and incubated overnight at 37°C. A 1/200 dilution of the starter culture was used to inoculate 1l of LB broth supplemented with 100  $\mu$ g/ml ampicillin and grown at 37 °C until an O.D.600 of 0.4-0.8 was achieved. Protein expression was induced by the addition of 1 mM IPTG and incubated for 24 hrs at 37 °C. The 1l culture was harvested by centrifugation (3000 x g for 20 mins) and the pellet resuspended in 25 ml 50 mM Tris-HCl, pH 8.0 (Buffer A) and supplemented with 1 mM phenylmethanesulfonylfluoride (PMSF) and 1 mM ethyldiaminetetraacetic acid (EDTA). The cells were lysed by French press using a chilled pressure cell. The lysate was then centrifuged (20000 rpm in a Beckman JA20 rotor for 30 mins) to pellet any cell debris and the supernatant was decanted and stored at 4°C. The cell lysate was subjected to fractionation with ammonium sulphate precipitation. An initial ammonium sulphate concentration of 45% (w/v) was used to precipitate unwanted proteins from solution. After clearance of unwanted precipitate by centrifugation (20000 rpm in a Beckman JA20 rotor for 40 mins) further addition of ammonium sulphate to a final concentration of 75% (w/v) was carried out to precipitate EGFP or EGFP<sup>G4 $\Delta$</sup> . The precipitate was resuspended in 5 ml Buffer A. The sample was buffer exchanged into fresh Buffer A by dialysis in a 10000 MWCO membrane to

remove any remaining ammonium sulphate. A precipitate formed during dialysis and was removed by centrifugation at 10,000 rpm in a Beckman JA-20 rotor for 20 min. The supernatant was applied to a Resource Q (GE Healthcare) anion exchange column (5 ml bed volume, flow rate 2 ml/min) equilibrated with Buffer A. Target proteins were eluted using a gradient from 0 mM to 500 mM NaCl in Buffer A over 5 column volumes with elution monitored at 280 nm and 488 nm. Pooled fractions were buffer exchanged into fresh Buffer A supplemented with 150 mM NaCl (Buffer B) with Amicon® Ultra centrifugal concentrators. Buffer exchanged protein samples were applied to a SP Superdex™ 200 gel filtration column (GE Healthcare) with elution monitored at 280 nm and 488 nm. The purified protein sample was finally stored in Buffer B. Protein concentration was determined with the DC Protein assay kit (Bio-Rad) using bovine serum albumin (BSA) as a protein standard. The assay was performed as to the manufactures guidelines for use in a microplate assay.

#### ***Size exclusion chromatography***

Gel filtration standards (Biorad) were applied to a Superdex™ 75 column (20 ml bed volume, 0.5 ml/min flow rate). As per the manufacturers guidelines with protein elution monitored at 280 nm. A standard curve was generated from the plot LogMw against  $K_{av}$ , where  $K_{av} = (V_e - V_o)/(V_t - V_o)$ ,  $V_e$  is the elution volume,  $V_t$  is the total volume and  $V_o$  is the void volume. Protein samples were prepared in Buffer B to final concentrations of 25, 50 or 100 uM and applied to a Superdex™ 75 column with protein elution monitored by absorbance at 488 nm. Elution volumes were determined for each sample and  $K_{av}$  values calculated. Using the standard curve estimated molecular weights could be determined for each protein sample.

#### ***Fit to 2 state unfolding.***

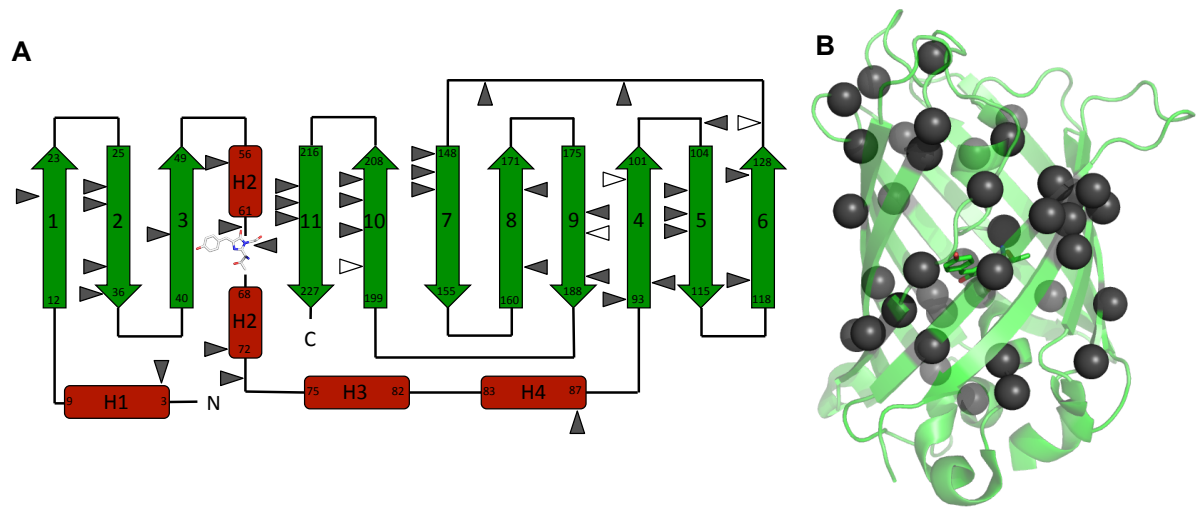
Equilibrium unfolding was fit to a 2-state model in the GraphPad Prism software (*equation 1*) to estimate approach to equilibrium (see Supporting Methods).

$$Y_N = \alpha_N + \beta_N [D], \quad Y_D = \alpha_D + \beta_D [D]$$

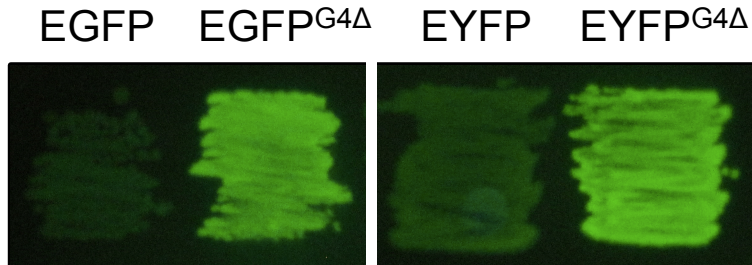
$$F = Y_N - (Y_N - Y_D) \frac{\exp\left(\frac{m_{N-D}([D] - [D]_{50\%})}{RT}\right)}{1 + \exp\left(\frac{m_{N-D}([D] - [D]_{50\%})}{RT}\right)} \quad \text{equation 1}$$

Where F is the fraction of folded protein,  $Y_N$  and  $Y_D$  are intensities of native and denatured states, respectively. To take into account sloping baselines for the fluorescence data,  $Y_N$  and  $Y_D$  are described as a function of  $\alpha_N$ ,  $\beta_N$ ,  $\alpha_D$  and  $\beta_D$ , respectively. Where  $\alpha_N$  and  $\alpha_D$  are the fluorescence intensities of the native and denatured states, respectively, and  $\beta_N$  and  $\beta_D$  are the slopes of the native and denatured baselines.  $m_{N-D}$  is a constant that describes the dependence of  $\Delta G$  on denaturant concentration, [D], between the native and denatured states.  $[D]_{50\%}$  is the estimated midpoint of the unfolding transition and represents the concentration of denaturant at which 50% of the protein is folded and 50% is unfolded.

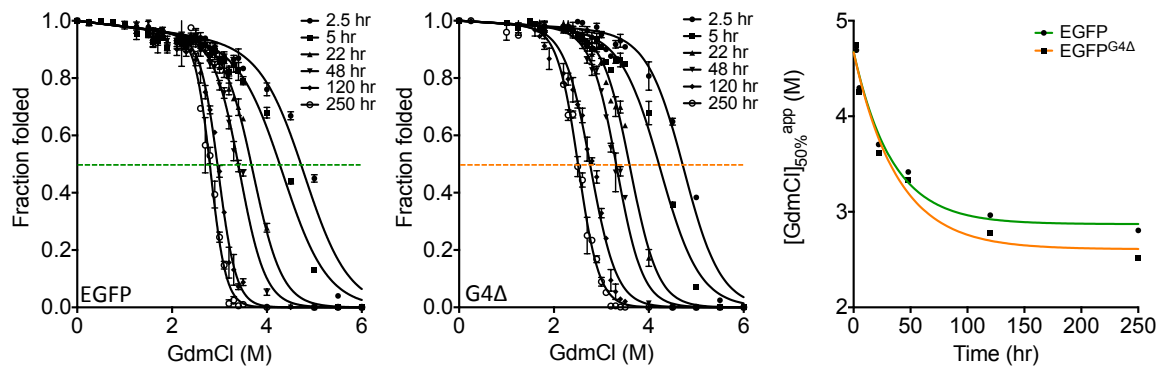
## Supporting Figures.



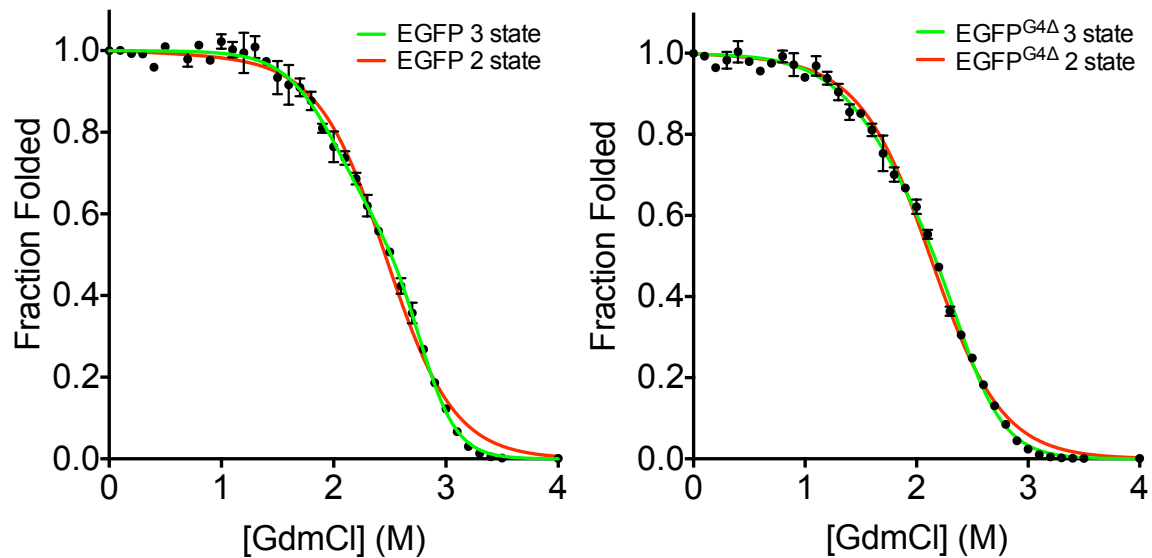
**Supporting Figure S1, related to Figure 1.** Mapping non-tolerated single amino acid deletion mutations with respect to EGFP (A) secondary and (B) tertiary structure. (A). The secondary structure arrangement and overall topology of EGFP shows the arrangement of  $\beta$ -strands (green),  $\alpha$ -helices (red) and loops (black). Disruptive single amino acid deletions identified in this study are indicated by black triangles and trinucleotide deletions generating stop codon are shown as white triangles. (B) Map of single amino acid deletions onto the tertiary structure of EGFP. Cartoon representation of EGFP (green) with disruptive deletions indicated by black spheres.



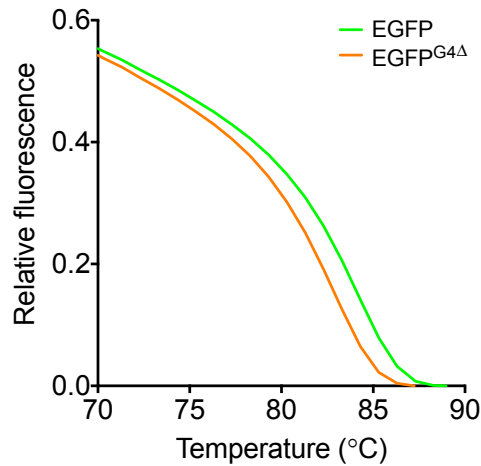
**Supporting Figure S2, related to Figure 3.** Colour version of cellular fluorescence of the EGFP and EYFP, and the corresponding G4 $\Delta$  variants presented in Figure 3 in the main manuscript.



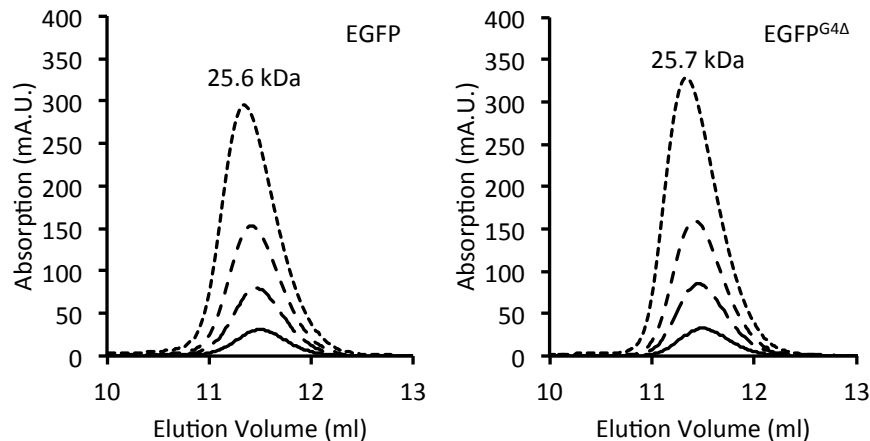
**Supporting Figure S3, related to Figure 4 and Table 1. Guanidinium chloride induced equilibrium unfolding and equilibrium kinetics.** Fluorescence emission at 520 nm after excitation at 480 nm was monitored for (A) EGFP and (B) EGFP<sup>G4Δ</sup>, over 250 hrs (as indicated in the figures) and data were fit to a two state model (GraphPad Prism). C, Apparent [GdmCl]<sub>50%</sub> values (the [GdmCl] at which 50% of the samples are in the native and 50% in the denatured states) were plot against time and fit to single exponential decay curves to assure close approach to equilibrium.



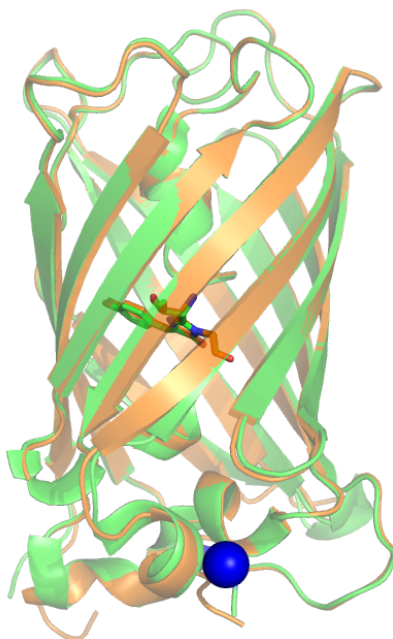
**Supporting Figure S4, related to Figure 4. Two state and three state model fits to equilibrium unfolding data.** Equilibrium unfolding data for EGFP (left panel) and EGFP<sup>G4Δ</sup> (right panel) fit to a two state (red) or three state (green) model highlights the poor fit of the data to a two state model.



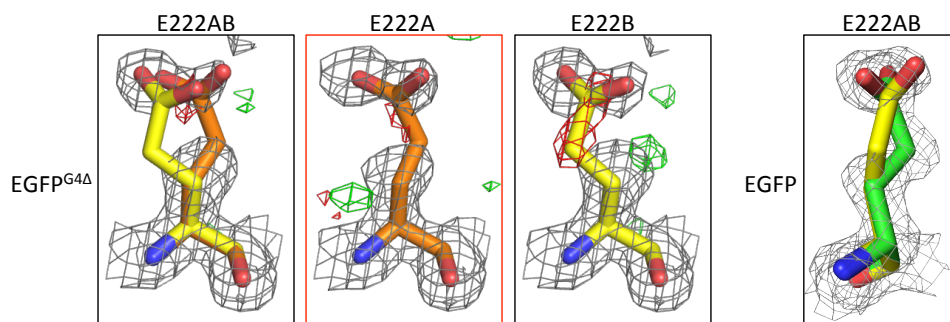
**Supporting Figure S5, related to Figure 4.** Thermal melting curves for EGFP and EGFP<sup>G4Δ</sup>. Melting temperatures ( $T_m$ ) of EGFP and EGFP<sup>G4Δ</sup> were determined by monitoring fluorescence with an Opticon 2 qPCR thermal cycler (MJ Research) while ramping the temperature from 25-98°C. Protein samples were diluted to a final concentration of 1  $\mu$ M in 50 mM sodium phosphate buffer pH 8.0 (total volume 50  $\mu$ l) and the temperature ramped at 1°C/min. MJ Research Software supplied with the qPCR machine was used to determine an apparent melting temperature.



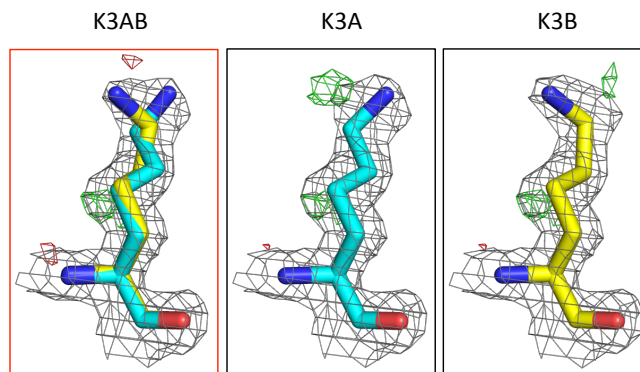
**Supporting Figure S6, related to Figure 5 and Table 3.** Size exclusion chromatography of EGFP<sup>G4Δ</sup>. The elution profiles of (A) EGFP and (B) EGFP<sup>G4Δ</sup> at 10  $\mu$ M (black line), 25  $\mu$ M (long dash), 50  $\mu$ M (medium dash) and 100  $\mu$ M (short dash). The estimated molecular weight based on the peak elution volume is shown on the graph.



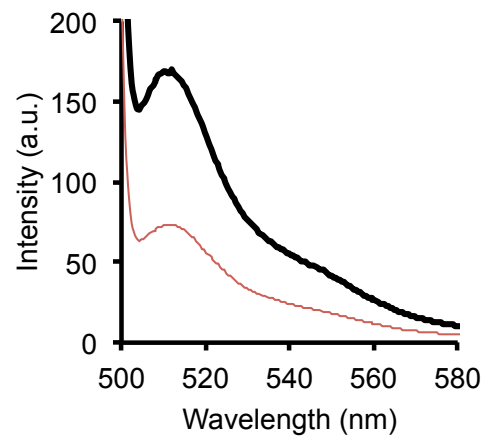
**Supporting Figure S7, related to Figure 5.** Overlap of EGFP (green) with EGFP<sup>G4Δ</sup> (orange) with the G4 residue in EGFP highlighted as a blue sphere and the chromophore shown as stick representation. The RMSDs between the two structures in terms of backbone and all atoms was 0.6Å and 1.2Å respectively.



**Supporting Figure S8, related to Figure 5. Rationale behind modelling of E222 as a single conformer in EGFP<sup>G4Δ</sup>.** Modelling of residue E222 as either the single conformer A (E222A), the single conformer B (E222B) or as a double conformer (E222AB). The electron density does not fully support the modelling of E222 in EGFP<sup>G4Δ</sup> as a double conformer. The model used in final crystal structure refinement is highlighted in the red box (E222A).



**Supporting Figure S9, related to Figure 5. Rationale behind modelling of K3 as a double conformer in EGFP<sup>G4Δ</sup>.** Modelling of residue K3 as either the single conformer A (K3A) or conformer B (K3B) does not fully satisfy the electron density. Modelling of residue K3 by both conformers does satisfy the electron density. The model used in final crystal structure refinement is highlighted in a red (K3AB) box.



**Supporting Figure S10, related to Figure 3.** Whole cell fluorescence emission (excited at 488 nm) spectra for cultures grown at 37°C expressing EGFP (black line) or EGFP<sup>K3N-G4Δ</sup>. Cell cultures were standardised to an OD<sub>600</sub> of 0.1.



**Supporting Table S1, related to Figure 1. Tolerated TNDs in *egfp* and subsequent amino acid mutations**

Nucleotide deletion <sup>a</sup>	Amino acid Mutation <sup>b</sup>	Frequency	Secondary structure <sup>c</sup>	SASA (Å <sup>2</sup> )	% SASA
<u>3</u> GTG AGC <sub>10</sub>	V1Δ S2G	2	N-terminus	ND	ND
<u>9</u> AAG GGC <sub>16</sub>	K3N G4Δ	4	H1	2.77	13
<u>12</u> GGC GAG <sub>19</sub>	G4Δ	8	H1	2.77	13
<u>12</u> GGC GAG <sub>19</sub>	E5Δ	2	H1	57.09	42
<u>18</u> GAG <sub>22</sub>	E6Δ	1	H1	84.96	42
<u>27</u> ACC GGG <sub>34</sub>	T9Δ G10R	6	H1	102.95	70
<u>27</u> ACC GGG <sub>34</sub>	G10Δ	2	Loop H1-S1	37.91	38
<u>36</u> GTG <sub>40</sub>	V12Δ	1	S1	9.20	12
<u>75</u> CAC <sub>79</sub>	H25Δ	2	S2	79.74	54
<u>114</u> ACC <sub>118</sub>	T38Δ	3	Loop S2-S3	65.55	37
<u>144</u> TGC <sub>148</sub>	C48Δ	1	S3	3.92	9
<u>147</u> ACC <sub>151</sub>	T50Δ	1	Loop S3-H2	81.50	50
<u>150</u> ACC GGC <sub>157</sub>	T50Δ G51S	2	Loop S3-H2	81.50	50
<u>159</u> CTG CCC <sub>166</sub>	L53Δ	1	Loop S3-H2	2.58	11
<u>225</u> CCC GAC <sub>232</sub>	P75Δ D76H	2	H3	21.59	17
<u>225</u> GAC <sub>229</sub>	D76Δ	2	H3	118.40	73
<u>237</u> AAG <sub>241</sub>	K79Δ	1	H3	58.66	24
<u>396</u> GAG GAC <sub>403</sub>	E132D D133Δ	1	Loop S6-S7	108.24	72
<u>411</u> GGG <sub>415</sub>	G138Δ	2	Loop S6-S7	26.72	21
<u>459</u> ATG GCC <sub>466</sub>	M153Δ A154T	2	S7	69.42	37
<u>462</u> GCC GAC <sub>469</sub>	A154Δ	5	S7	30.50	23
<u>465</u> GAC <sub>469</sub>	D155Δ	4	S7	22.16	22
<u>474</u> AAG AAC <sub>481</sub>	K158Δ	1	Loop S7-S8	106.96	57
<u>480</u> GGC <sub>484</sub>	G160Δ	1	S8	11.54	10
<u>513</u> ATC GAG <sub>520</sub>	I171M E172Δ	3	Loop S8-S9	88.73	39
<u>522</u> GGC <sub>526</sub>	G174Δ	2	Loop S8-S9	68.18	52
<u>525</u> AGC <sub>529</sub>	S175Δ	1	Loop S8-S9	59.04	34
<u>567</u> GGC GAC <sub>574</sub>	G189Δ	1	Loop S9-S10	22.96	36
<u>570</u> GAC GGC <sub>577</sub>	D190Δ	1	Loop S9-S10	152.83	100
<u>576</u> CCC GTG <sub>583</sub>	P192Δ V193L	3	Loop S9-S10	130.44	95
<u>588</u> CCC <sub>592</sub>	P196Δ	1	Loop S9-S10	5.11	16
<u>591</u> GAC <sub>595</sub>	D197Δ	1	Loop S9-S10	54.34	62
<u>594</u> AAC <sub>598</sub>	N198 Δ	1	Loop S9-S10	100.68	71
<u>633</u> CCC AAC <sub>640</sub>	P211Δ N212H	3	Loop S10-S11	112.07	58
<u>678</u> GCC GCC GGG <sub>687</sub>	A226Δ A227Δ	1	S11	30.10 / 28.62	12 / 20
<u>681</u> GCC GGG <sub>688</sub>	A227Δ	5	S11	28.62	20
<u>681</u> GCC GGG <sub>688</sub>	G228Δ	2	C-terminus	48.44	38
<u>690</u> ACT CTC <sub>697</sub>	L231Δ	1	C-terminus	178.68	93
<u>699</u> ATG GAC <sub>706</sub>	M233Δ D234N	2	C-terminus	ND	ND
<u>702</u> GAC GAG <sub>709</sub>	D234E E235Δ	2	C-terminus	ND	ND
<u>705</u> GAG <sub>709</sub>	E235Δ	1	C-terminus	ND	ND
<u>711</u> TAC <sub>715</sub>	Y237Δ	1	C-terminus	ND	ND

<sup>a</sup> Numbers refer to gene sequence numbering for *egfp* (GFPmut1)

<sup>b</sup> Δ after a residue number signifies that residue has been deleted, protein numbering as per wtGFP

<sup>c</sup> Secondary structure elements as defined by Fig 1, helices (H), strands (S).

**Supporting Table S2, related to Figure 1. Non-tolerated TNDs in *egfp* and subsequent amino acid mutations**

Nucleotide deletion <sup>a</sup>	Amino acid Mutation <sup>b</sup>	Frequency	Secondary structure <sup>c</sup>	SASA (Å <sup>2</sup> )
<sup>9</sup> <u>AAG GGC</u> <sub>16</sub>	K3Δ G4S	1	H1	178.25
<sup>60</sup> <u>GGC GAC</u> <sub>67</sub>	G20Δ	3	S1	5.93
<sup>81</sup> <u>TTC AGC</u> <sub>88</sub>	F27Δ S28C	1	S2	5.40
<sup>90</sup> <u>TCC GGC</u> <sub>97</sub>	S30Δ G31C	3	S2	31.28
<sup>99</sup> <u>GGC GAG</u> <sub>106</sub>	E34Δ	2	S2	89.14
<sup>105</sup> <u>GGC GAT</u> <sub>112</sub>	D36Δ	1	S2	26.72
<sup>135</sup> <u>AAG TTC</u> <sub>142</sub>	K45Δ F46I	1	S3	45.42
<sup>168</sup> <u>CCC TGG</u> <sub>175</sub>	W57Δ	1	H2	12.84
<sup>171</sup> <u>TGG</u> <sub>174</sub>	W57Δ	3	H2	12.84
<sup>189</sup> <u>ACC CTG</u> <sub>196</sub>	L64Δ	1	Loop H2-H3	0.00
<sup>192</sup> <u>CTG ACC</u> <sub>199</sub>	L64Δ T65P	2	Loop H2-H3/Cro	0.00
<sup>198</sup> <u>TAC GGC</u> <sub>205</sub>	Y66Δ G67C	1	Cro	ND
<sup>216</sup> <u>AGC</u> <sub>220</sub>	S72Δ	1	H3	2.38
<sup>219</sup> <u>CGC</u> <sub>223</sub>	R73Δ	1	Loop H3-H4	87.13
<sup>261</sup> <u>GCC</u> <sub>265</sub>	A87Δ	2	H5	5.30
<sup>279</sup> <u>GTC CAG</u> <sub>286</sub>	V93Δ Q94E	1	S4	19.40
<sup>282</sup> <u>CAG</u> <sub>286</sub>	Q94Δ	1	S4	5.31
<sup>300</sup> <u>TTC AAG</u> <sub>307</sub>	F100Δ K101STOP	1	S4	3.91
<sup>309</sup> <u>GAC GGC</u> <sub>316</sub>	D103Δ	1	Loop S4-S5	28.42
<sup>321</sup> <u>AAG ACC</u> <sub>328</sub>	K107Δ	1	S5	98.33
<sup>330</sup> <u>GCC GAG</u> <sub>337</sub>	A110Δ	3	S5	5.59
<sup>330</sup> <u>GCC GAG</u> <sub>337</sub>	E111Δ	1	S5	53.03
<sup>360</sup> <u>GTG</u> <sub>364</sub>	V120Δ	3	S6	8.67
<sup>360</sup> <u>GTG AAC</u> <sub>367</sub>	V120Δ N121D	1	S6	8.67
<sup>381</sup> <u>GGC ATC</u> <sub>388</sub>	G127Δ I128V	1	S6	0.42
<sup>390</sup> <u>TTC AAG</u> <sub>397</sub>	F130Δ K131STOP	1	Loop S6-S7	10.87
<sup>411</sup> <u>CTG</u> <sub>415</sub>	L137Δ	1	Loop S6-S7	22.36
<sup>435</sup> <u>TAC</u> <sub>439</sub>	Y145Δ	1	Loop S6-S7	23.93
<sup>444</sup> <u>CAC</u> <sub>448</sub>	H148Δ	3	S7	9.18
<sup>450</sup> <u>GTC TAT</u> <sub>457</sub>	V150Δ Y151D	3	S7	0.01
<sup>450</sup> <u>GTC TAT</u> <sub>457</sub>	Y151Δ	2	S7	103.92
<sup>486</sup> <u>AAG</u> <sub>490</sub>	K162Δ	1	S8	64.53
<sup>507</sup> <u>CAC</u> <sub>511</sub>	H169Δ	3	S8	8.20
<sup>510</sup> <u>AAC ATC</u> <sub>516</sub>	N170Δ	1	S8	50.12
<sup>540</sup> <u>GAC</u> <sub>544</sub>	D180Δ	1	S9	44.22
<sup>546</sup> <u>TAC CAG</u> <sub>553</sub>	Y182STOP Q183Δ	2	S9	0.00
<sup>561</sup> <u>CCC</u> <sub>565</sub>	P187Δ	1	S9	17.65
<sup>600</sup> <u>TAC CTG</u> <sub>607</sub>	Y200STOP L201Δ	1	S10	0.55
<sup>609</sup> <u>ACC CAG</u> <sub>616</sub>	Q204Δ	1	S10	101.34
<sup>615</sup> <u>TCC GCC</u> <sub>622</sub>	A206Δ	1	S10	55.05
<sup>618</sup> <u>GCC CTG</u> <sub>625</sub>	L207Δ	1	S10	22.63
<sup>621</sup> <u>CTG AGC</u> <sub>628</sub>	L207Δ S208R	1	S10	22.63
<sup>654</sup> <u>ATG GTC</u> <sub>661</sub>	M218I V219Δ	1	S11	27.30
<sup>660</sup> <u>CTG</u> <sub>664</sub>	L220Δ	1	S11	0.00
<sup>663</sup> <u>CTG</u> <sub>667</sub>	L221Δ	1	S11	65.26

<sup>a</sup> Numbers refer to gene sequence numbering for *egfp* (GFPmut1)

<sup>b</sup> Δ after a residue number signifies that residue has been deleted, protein numbering as per wtGFP

<sup>c</sup> Secondary structure elements as defined by Fig 1, helices (H), strands (S).

## **Supporting References**

Baldwin, A.J., Arpino, J.A., Edwards, W.R., Tippmann, E.M., and Jones, D.D. (2009). Expanded chemical diversity sampling through whole protein evolution. *Molecular BioSystems* 5, 764-766.