Supporting Information

Chemical and Biological Incorporation of the Blue Florescent Amino Acid 4-Cyanotryptophan into Proteins: Application to Tune the Absorption and Emission Wavelengths of GFP

Manxi Wang and Feng Gai*

Beijing National Laboratory for Molecular Sciences, College of Chemistry and Molecular Engineering, Peking University, Beijing 100871, China. E-mail: fgai@pku.edu.cn

Table of Contents

- 1. General
- 2. Syntheses of 3B4CNI
- 3. Peptide Synthesis
- 4. Peptide Labeling via Cysteine Alkylation
- 5. Spectroscopic Measurement of YGGC*GG
- 6. Protein Labeling via Cysteine Alkylation
- 7. Tryptophan Auxotrophic E. coli. Fermentation and Protein Purification
- 8. MS-MS Results of 4CN-Trp-GFP and 4CN-Trp-Cerulean
- 9. Spectroscopic Measurement of 4CN-Trp-Cerulean
- 10. NMR Spectra
- 11. References

1. General

Reactions that require the removal of oxygen and water were carried out in a nitrogen atmosphere using Schlenk technology.

NMR measurements were performed on a Bruker 400 MHz NMR spectrometer at ambient temperature. Chemical shifts are reported in parts per million (ppm) in δ unit relative to that of chloroform-d (δ H = 7.26 ppm) or dimethyl sulfoxide-d6 (δ H = 2.50 ppm).

TLC was carried out on aluminum TLC plates with silica gel coated with fluorescent indicator F254 (40×80 mm, SIL G/UV254, 0.2 mm layer thickness) purchased from Merck. Detection was carried out using UV-light (254 nm or 366 nm).

High-resolution mass spectra were collected on a Quadrupole-TOF LC-MS/MS System (Waters Voin IMS QTof) in positive ion mode.

2. Synthesis of 3-(bromomethyl)-1H-indole-4-carbonitrile (3B4CNI)



In the first step, 2.18 g of 4-cyanoindole (1.0 eq, 14 mmol) was mixed with 12 mL of anhydrous DMF in a 100-mL three-necked flask at room temperature, and the solution was then cooled to 0 °C. After the solute was fully dissolved, an appropriate amount of POCl₃ (2.5 eq, in 3.2 mL DMF) was slowly added to the flask under N₂ protection. The reaction was then allowed to run overnight, followed by addition of 5 N NaOH solution (10 eq) to the flask and, subsequently, the temperature of the reaction solution was raised to 100 °C and heated under reflux for 10 min. After cooling the reaction mixture to room temperature, the brown solid was removed by filtration, and the filtrate was extracted with ethyl acetate (EA), dried over MgSO₄, and then rotary evaporated to yield the desired crude product, which was directly used in the next step.

In the second step, the crude product obtained above and 1.12 g of NaH (60% dispersion in mineral oil, 28 mmol) were dissolved in 40 mL of anhydrous DMF in a 100-mL three-necked flask under N₂ protection, then 3.73 g of tosyl chloride (TsCl, 19.6 mmol) was added in steps. After the reaction was carried out at room temperature for 4 hours, saturated brine was added to quench the reaction. The organic phase was extracted, filtered and evaporated. The residue was purified by chromatography on a silica gel column (elution with 3:2 hexane-EA) to yield a pale-yellow product (3-formyl-1-tosyl-1*H*-indole-4-carbonitrile). $R_f = 0.78$ (silica gel, EtOAc/hexane = 1/4).



In the third step, a suspension of 1.3 g (1.0 eq, 4 mmol) of the product obtained in the second step in 10 mL of EtOH was first mixed with 302 mg (2.0 eq, 8 mmol) of NaBH₄ and the mixture was stirred at room temperature for 3 h. Then, the reaction was quenched by NaHCO₃ aqueous solution, and the principate thus formed was filtered and dried under vacuum. The resultant white solid (3-hydroxymethyl-1-tosyl-1*H*-indole-4-carbonitrile) (its ¹H NMR spectrum is given in Fig S10) was directly used in the next step. $R_f = 0.40$ (silica gel, EtOAc/hexane = 1/3).



In the fourth step, 2.1 g (1.0 eq, 6.6 mmol) of the product obtained from the third step was first dissolved in 30 mL of anhydrous DCM under N_2 protection, and the solution was cooled to - 78 °C. Then, 744 µL (1.2 eq, 7.92 mmol) of PBr₃ was dropwise added , followed by stirring

for 5 h. The reaction mixture was allowed to slowly warm up to 20 °C. Subsequently, the solution was diluted with 30 mL of brine, and extracted with EA after removing DCM by rotary evaporation. The combined organic phase was then dried with Na₂SO₄ and concentrated under diminished pressure to yield the pale-yellow crude product (3-(bromomethyl)-1-tosyl-1*H*-indole-4-carbonitrile), which was further purified by chromatography on a silica gel column (elution with 3:1 hexane-EA). $R_f = 0.60$ (silica gel, EtOAc/hexane = 1/3).



Finally, 953 mg (1.0 eq, 2.45 mmol) of 3-(bromomethyl)-1-tosyl-1*H*-indole-4-carbonitrile and 1.6 g (2.0 eq, 4.9 mmol) of Cs_2CO_3 were mixed in a 150-mL flask, followed by addition of a 45 mL mixture of THF and MeOH (2:1). This reaction mixture was stirred for 3 h under room temperature and N₂ protection and then diluted with 30 mL of brine and extracted with EA after removing THF and MeOH. The pale-yellow crude product thus obtained was purified by chromatography on a silica gel column (elution 2:1 hexane-EA) to yield pure 3-(bromomethyl)-1*H*-indole-4-carbonitrile (3B4CNI; the corresponding ¹H NMR spectrum is given in Fig S11). R_f = 0.20 (silica gel, EtOAc/hexane = 1/3).

3. Peptide Synthesis

The YGGCGG peptide was synthesized via the standard 9-fluorenylmethoxycarbonyl (Fmoc) solid-phase synthesis protocol with Fmoc-protected amino acids from GL Biochem (Shanghai) Ltd on a CEM (Matthews, NC) Liberty Blue automated microwave peptide synthesizer. Peptide removal from the solid-phase Gly-Wang resin was achieved using the following cleavage cocktail: 92.5:2.5:2.5 TFA/DODT/water/TIPS (for 3 h). Crude peptide sample was isolated by precipitation using cold ether (10 volumes), purified by reverse-phase HPLC (Agilent Technologies 1260 Infinity) using a C18 column, and characterized by LC-MS (Quadrupole-TOF LC-MS/MS System). The molecular mass of the peptide product was measured to be 513.21, which is in agreement with the expected value (512.17).

4. Peptide Labeling via Cysteine Alkylation

An appropriate amount of the lyophilized peptide solid was first dissolved in sodium phosphate buffer (pH 8.0, 50 mM) to reach a final peptide concentration of 1 mM. Then, 200 μ L of this peptide solution was mixed with 20 μ L of a 100-mM 3B4CNI DMF solution to reach a 3B4CNI:peptide ratio of 10:1, allowing the cysteine alkylation reaction to begin under vigorous shaking condition at 37 °C. The progression of the reaction was monitored via mass spectrometry and, as shown (Fig. S1), the coupling reaction nearly reaches completion in ca. 3 h. Subsequently, any unreacted solids were removed from the reaction mixture by centrifugation (5 min, 4 °C, 10000 rpm). The supernatant was separated via HPLC, and the 4CNI-labeled peptide was collected and lyophilized for subsequent spectroscopic measurements. For the labeling reaction performed in PBS solution containing 8 M urea (pH 8.0), the same separation and purification procedures were used. As shown (Fig. S2), mass spectrometric measurements confirm the success of the cysteine alkylation, as the measured molecular mass (667.29) agree with the expected value (666.71).



Fig. S1 Progression of the reaction as assessed via mass spectrometry.



Fig. S2 Molecular mass of the labeled peptide obtained under reaction conditions where urea (8 M) was either absent (A) or present (B).

5. Spectroscopic Measurements of YGGC*GG

We measured the absorption and fluorescence spectra of the 4CNI-labelled peptide (i.e., YGGC*GG) in H₂O and the results are shown in Fig. S3, where the absorption and fluorescence spectra of $GW_{CN}G$ ($W_{CN} = 4$ -cyanotryptophan) in H₂O are also shown for comparison. Additionally, the absorption spectra of YGGC*GG in different solvents are shown in Fig S4.



Fig. S3 Normalized absorption (A) and fluorescence (B) spectra of $GW_{CN}G$ and $YGGC^*GG$ in H₂O, as indicated. For the fluorescence measurement, the excitation wavelength was 325 nm.



Fig. S4 Normalized absorption spectra of YGGC*GG in EtOH(A), IPA(B) and THF(C).

The fluorescence quantum yields (QYs) of YGGC*GG were determined using the methods described in the text. The I_F versus I_{abs} plots are shown in Fig. S5.



Fig. S5 I_F versus I_{abs} plots of YGGC*GG in H₂O (A), EtOH(B), IPA(C), THF(D) and DPA in cyclohexane.

Table S1 The maximum emission wavelength (λ_{em}) and fluorescence QY of YGGC*GG obtained in different solvents.

	λ_{em} (nm)	QY
H ₂ O	437	0.21 ± 0.03
EtOH	419	0.24 ± 0.07
IPA	415	0.20 ± 0.06
THF	401	0.23 ± 0.07

6. Protein Labeling via Cysteine Alkylation

BSA was prepared as a 13.2 mg/mL solution in phosphate buffer (pH 8.0, 50 mM). A 100 μ L aliquot of this protein solution was mixed with 4.6 mg of 3B4CNI (dissolved in 10 μ L DMF) in a 1.5 mL plastic tube. Then, this reaction mixture was vortexed and shaken vigorously at 20 °C for 3 h, followed by centrifugation (5 min, 4 °C, 10000 rpm) to remove any solids.



Fig. S6 Mass spectrometry results for labeled (above) and unlabeled (below) BSA.

7. Tryptophan (Trp) Auxotrophic E. coli. Fermentation and Protein Purification

Trp-auxotrophic *E. coli* strain was purchased from ATCC (catalog number 49980 [genotype WP2 uvrA]). All fermentation and expression experiments were performed in GMML supplemented with 200 μ M glucose (minimal media).¹ The expression host *E. coli* ATCC49980 was routinely transformed with either plasmid: pET22b-sfGFPwt (with a His-tag at the protein N-terminus) and pET22b-Cerulean (with a His-tag at the protein C-terminus using standard cloning methods.

The transformed Trp-auxotrophic strain was first grown in minimal media in the presence of 0.0075 mM Trp and 1 mM ampicillin. When the cell culture was grown to ~0.6 OD600, the cells were collected by low-speed centrifugation (25 min, 4 °C, 4000 rpm). Then, the harvested cells were resuspended in new minimal media that was saturated with 4CN-Trp, followed by vigorous shaking at 37 °C for 45 min. Then, 1 mM IPTG was added to this media to induce sfGFP (or Cerulean) expression.² After incubation at 30 °C overnight (for sfGFP) or at 27 °C for 24 h (for Cerulean), cells were collected again by low-speed centrifugation (25 min, 4 °C, 4000 rpm). Subsequently, the cell pellet was resuspended in PBS buffer (pH 7.4, 50 mM), which was then subject to ultrasonication for 6 min (1 s on and 2 s off) and subsequently highspeed centrifugation for 30 min (at 4 °C, 15000 rpm). Followingly, the resultant supernatant was purified using a HiTrap HP column attached to an AKTA FPLC to yield the desired protein solution. Further SDS-PAGE analysis of this solution indicates the existence of a ~27 kDa protein (Fig. S7), showing that sfGFP is successfully expressed. Finally, this protein solution was concentrated by ultrafiltration and desalted via a desalting column.

116.0 kDa	- Filmer	
66.2 kDa		
45.0 kDa •		
35.0 kDa		
	sfGEP	4CN-Trp-
25.0 kDa		sfGFP
18.4 kDa		3
14.4 kDa		

Fig. S7 SDS-PAGE analysis of purified 4CNI-GFP, which indicates the presence of a protein with a molecular weight of \sim 27 kDa.

8. MS-MS Results of 4CN-Trp-sfGFP and 4CN-Trp-Cerulean

To further confirm the incorporation selectivity of 4CN-Trp, the expressed 4CN-Trp-sfGFP and 4CN-Trp-Cerulean, which was pretreated by Chymotrypsin, was assessed by MS/MS. The residues mutated to 4CN-Trp are highlighted in bold.

Sequence	Modifications	Positions in master protein	Charge	m/z [Da]
EFVTAAGITHGMDELY	1xOxidation [M12]	[222-237]	2	885.4071
SKDPNEKRDHMVLLEFVTA AGITHGMDELY	1xOxidation [M11; M26]	[208-237]	3	1144.89111
MSKGEELFTGVVPILVELD GDVNGHKF	1xMet-loss [N- Term]	[1-27]	3	933.82391
ICTTGKLPVPWPTLVTTL	1xCarbamidomethy 1[C2]; 1xTrp- > C9H6N2 [W11]	[47-64]	2	1011.55865
ICTTGKLPVPWPTLVTTLTY	1xCarbamidomethy 1[C2]; 1xTrp- > C9H6N2 [W11]	[47-66]	2	1143.6145

Table S2 Mass spectral assignments of peptic digest of 4CN-Trp-sfGFP.

SVRGEGEGDATNGKLTLKF	[28-46]	2	990.01355
ITADKQKNGIKANF	[152-165]	3	516.6228
ITADKQKNGIKANFKIRHN VEDGSVQLADHY	[152-182]	4	878.21136
VQERTISF	[93-100]	2	490.26413
MSKGEELFTGVVPILVELD GDVNGHKF	[1-27]	3	977.5036
NFNSHNVY	[144-151]	2	497.72223
NFNSHNVYITADKQKNGIK ANF	[144-165]	3	841.76349
QQNTPIGDGPVLLPDNHY	[183-200]	2	989.48755
QQNTPIGDGPVLLPDNHYL	[183-201]	2	1046.02966
QQNTPIGDGPVLLPDNHYL STQSVL	[183-207]	2	1353.69214
SKDPNEKRDHMVLL	[208-221]	2	841.43842
LEFVTAAGITHGMDELY	[221-237]	2	933.95099
SKDPNEKRDHMVLLEFVTA AGITHGMDELY	[208-237]	3	1139.55945
VNRIELKGIDFKEDGNIL	[120-137]	3	691.71631
SRYPDHMKRHDF	[72-83]	3	530.25293
SRYPDHMKRHDFF	[72-84]	3	579.2757
SVRGEGEGDATNGKL	[28-42]	3	497.24619
SVRGEGEGDATNGKLTL	[28-44]	2	852.43176
TGVVPILVELDGDVNGHKF	[9-27]	2	1005.03961
SKDPNEKRDHMVLLEF	[208-223]	2	979.4942
ADHYQQNTPIGDGPVLLPD NHY	[179-200]	2	1232.5824

KTRAEVKFEGDTLVNRIEL KGIDFKEDGNIL	[107-137]	3	1183.30933
EFVTAAGITHGMDELY	[222-237]	2	877.40955
EYNFNSHNVY	[142-151]	2	643.77551
FKSAMPEGY	[84-92]	2	515.23914
FKSAMPEGYVQERTISF	[84-100]	2	995.49146
GHKLEYNFNSHNVY	[138-151]	2	861.40503
KDDGTYKTRAEVKFEGDTL VNRIEL	[101-125]	3	966.50555
KEDGNILGHKLEY	[131-143]	3	505.93155
KGIDFKEDGNIL	[126-137]	2	674.85919
KGIDFKEDGNILGHKL	[126-141]	3	595.32831
KGIDFKEDGNILGHKLEY	[126-143]	2	1038.54248
VQERTISFKDDGTY	[93-106]	3	553.60583
KGIDFKEDGNILGHKLEYNF	[126-145]	3	779.73444
KIRHNVEDGSVQL	[166-178]	2	747.90472
KIRHNVEDGSVQLADHY	[166-182]	2	990.99866
KIRHNVEDGSVQLADHYQQ NTPIGDGPVL	[166-194]	3	1067.54175
KIRHNVEDGSVQLADHYQQ NTPIGDGPVLLPDNHY	[166-200]	3	1313.98779
KSAMPEGY	[85-92]	2	441.7049
KSAMPEGYVQERTISF	[85-100]	2	921.95679
KSAMPEGYVQERTISFKDD GTY	[85-106]	2	1261.59888
KTRAEVKFEGDTL	[107-119]	2	747.40192

KTRAEVKFEGDTLVNRIEL	[107-125]	3	740.07831
KTRAEVKFEGDTLVNRIEL KGIDF	[107-130]	3	926.84387
KGSHHHHHH	[238-246]	3	371.84476
VTAAGITHGMDELY	[224-237]	2	739.35406

Table S3 Mass spectral assignments of peptic digest of 4CN-Trp-Cerulean.

Sequence	Modifications	Positions in master protein	Charge	m/z [Da]
ARYPDHMKQHDFFKSAMPE GYVQERTIF	2xOxidation [M7; M17]	[73-100]	4	865.90973
ARYPDHMKQHDF	1xOxidation [M7]	[73-84]	2	780.85419
ARYPDHMKQHDFF	1xOxidation [M7]	[73-85]	3	569.92914
ARYPDHMKQHDFFKSAMPE GY	1xOxidation [M7; M17]	[73-93]	3	857.72235
ARYPDHMKQHDFFKSAMPE GYVQERTIF	1xOxidation [M7; M17]	[73-100]	4	861.91138
FKSAMPEGYVQERTIF	1xOxidation [M5]	[85-100]	3	640.31671
KSAMPEGY	1xOxidation [M4]	[86-93]	2	449.70236
KSAMPEGYVQERTIF	1xOxidation [M4]	[86-100]	3	591.29254
KSAMPEGYVQERTIFFKDDG NY	1xOxidation [M4]	[86-107]	4	653.55817
KSAMPEGYVQERTIFF	1xOxidation [M4]	[86-101]	2	959.97174
EFVTAAGITLGMDELY	1xOxidation [M12]	[223-238]	2	873.41919
MVSKGEELF	1xOxidation [M1]	[1-9]	2	528.25751
MVSKGEELFTGVVPILVELD GDVNGHKF	1xOxidation [M1]	[1-28]	3	1015.85876
MVSKGEELFTGVVPILVELD GDVNGHKF	1xMet-loss [N- Term]	[1-28]	3	966.8468

MVSKGEEL	1xMet-loss [N- Term]	[1-8]	2	381.20551
MVSKGEELF	1xMet-loss [N- Term]	[1-9]	2	454.73987
KFICTTGKLPVPWPTLVTTL	1xCarbamidomethy 1 [C4]; 1xTrp- > C9H6N2 [W13]	[46-65]	2	1149.13953
ICTTGKLPVPWPTLVTTL	1xCarbamidomethy 1 [C2]; 1xTrp- > C9H6N2 [W11]	[48-65]	2	1011.55811
MVSKGEELF	1xAcetyl [N-Term]	[1-9]	2	541.26489
LSTQSALSKDPNEKRDHMVL		[202-221]	3	757.05927
ITADKQKNGIKANF		[153-166]	2	774.43109
ITADKQKNGIKANFKIRHNIE DGSVQL		[153-179]	3	1013.22656
ITADKQKNGIKANFKIRHNIE DGSVQLADHY		[153-183]	4	881.71594
KTRAEVKFEGDTL		[108-120]	2	747.40167
VNRIELKGIDFKEDGNILGHK LEY		[121-144]	3	934.1709
VQERTIF		[94-100]	2	446.74768
KTRAEVKFEGDTLVNRIEL		[108-126]	2	1109.61243
LSTQSALSKDPNEKRDHMVL LEF		[202-224]	3	886.79089
SVSGEGEGDATYGKLTLKF		[29-47]	2	979.98907
SVSGEGEGDATYGKLTL		[29-45]	2	842.40784
MVSKGEELFTGVVPILVELD GDVNGHKF		[1-28]	3	1010.5271
KSAMPEGYVQERTIFFKDDG NY		[86-107]	2	1298.11255
TGVVPILVELDGDVNGHKF		[10-28]	2	1005.03845

QQNTPIGDGPVLLPDNHY	[184-201]	2	989.48767
QQNTPIGDGPVLLPDNHYL	[184-202]	3	697.68829
QQNTPIGDGPVLLPDNHYLS TQSAL	[184-208]	3	893.45178
SKDPNEKRDHMVL	[209-221]	2	784.89661
SKDPNEKRDHMVLL	[209-222]	2	841.43829
STQSALSKDPNEKRDHMVL	[203-221]	3	719.36389
SVSGEGEGDATY	[29-40]	2	586.24091
SVSGEGEGDATYGKL	[29-43]	2	735.34149
MVSKGEELF	[1-9]	2	520.26007
NAISDNVYITADKQKNGIKA NFKIRHNIEDGSVQLADHY	[145-183]	4	1100.81482
ADHYQQNTPIGDGPVLLPDN HY	[180-201]	3	822.05627
KSAMPEGYVQERTIF	[86-100]	2	878.44
ADHYQQNTPIGDGPVLLPDN HYL	[180-202]	3	859.74969
ARYPDHMKQHDF	[73-84]	2	772.8573
ARYPDHMKQHDFF	[73-85]	3	564.59637
ARYPDHMKQHDFFKSAMPE GY	[73-93]	3	852.39191
ARYPDHMKQHDFFKSAMPE GYVQERTIF	[73-100]	4	857.91321
EFVTAAGITLGMDELY	[223-238]	2	865.42157
EGDTLVNRIEL	[116-126]	2	629.836

EYNAISDNVY	[143-152]	2	594.26428
FKDDGNYKTRAEVKF	[101-115]	2	909.46313
FKDDGNYKTRAEVKFEGDT L	[101-120]	3	778.38611
FKDDGNYKTRAEVKFEGDT LVNRIEL	[101-126]	3	1019.85974
FKSAMPEGY	[85-93]	2	515.23895
FKSAMPEGYVQERTIF	[85-100]	2	951.97443
GHKLEYNAISDNVY	[139-152]	2	811.89423
KSAMPEGYVQERTIFF	[86-101]	2	951.97437
GHKLEYNAISDNVYITADKQ KNGIKANF	[139-166]	3	1051.21143
KDDGNYKTRAEVKF	[102-115]	2	835.9295
KDDGNYKTRAEVKFEGDTL	[102-120]	3	729.36267
KDDGNYKTRAEVKFEGDTL VNRIEL	[102-126]	3	970.83685
KEDGNILGHKLEY	[132-144]	3	505.9317
KEDGNILGHKLEYNAISDNV YITADKQKNGIKANF	[132-166]	5	785.00873
KGIDFKEDGNIL	[127-138]	2	674.85913
KGIDFKEDGNILGHKL	[127-142]	2	892.48993
KGIDFKEDGNILGHKLEY	[127-144]	3	692.6969
KGIDFKEDGNILGHKLEYNAI SDNVY	[127-152]	3	984.82898
KIRHNIEDGSVQL	[167-179]	2	754.91315

KIRHNIEDGSVQLADHY	[167-183]	3	665.67468
KIRHNIEDGSVQLADHYQQN TPIGDGPVLLPDNHY	[167-201]	3	1318.65637
KSAMPEGY	[86-93]	2	441.70499
VQERTIFF	[94-101]	2	520.28223
GKLTLKF	[41-47]	2	403.76038
VQERTIFFKDDGNY	[94-107]	3	577.94983

*Due to the chemical reaction during the formation of the chromophore, the 4CN-Trp changes to the chromophore cannot be identified, but the successful incorporation of the unnatural amino acid is verified by the spectral changes.

9. Spectroscopic Measurement of 4CN-Trp-Cerulean

We measured the fluorescence QY and decay kinetic of 4CN-Trp-Cerulean in PBS buffer and the results are shown in Fig. S8 and Fig. S9. The QY of 4CN-Trp-Cerulean was determined using the aforementioned equation.



Fig. S8 I_{F} versus I_{abs} plots of fluorescein in EtOH (A) and 4CN-Trp-Cerulean in PBS buffer (B).



Fig. S9 Fluorescence decay kinetics of 4CN-Trp-Cerulean in PBS buffer with the residues of the double-exponential fits shown on the top.



10. NMR Spectra

Fig. S10 ¹H NMR spectrum of 3-hydroxymethyl-1-tosyl-1H-indole-4-carbonitrile.



Fig. S11 ¹H NMR spectrum of 3-(bromomethyl)-1*H*-indole-4-carbonitrile.

11. References

1. K. Boknevitz; J. S. Italia; B. Li; A. Chatterjee and S.-Y. Liu. *Chem. Sci.* 2019, **10**, 4994-4998.

2 .J.-D. Pédelacq; S. Cabantous; T. Tran; T. C. Terwilliger and G. S. Waldo. *Nat. Biotechnol.* 2006, **24**, 79-88.