

Supplemental data

Comparative investigation into formycin A and pyrazofurin A biosynthesis reveals branch pathways for the construction of C-nucleoside scaffolds

Meng Zhang,^{1,5} Peichao Zhang,^{1,2,5} Gudan Xu,^{1,5} Wenting Zhou,¹ Yaojie Gao,¹ Rong Gong,¹

You-Sheng Cai,¹ Hengjiang Cong,³ Zixin Deng,¹ Neil P. J. Price,⁴ Xiangzhao Mao^{2*}, Wenqing Chen^{1*}

¹Key Laboratory of Combinatorial Biosynthesis and Drug Discovery, Ministry of Education, and School of Pharmaceutical Sciences, Wuhan University, Wuhan 430071, China

²College of Food Science and Engineering, Ocean University of China, Qingdao 266003, China

³College of Chemistry and Molecular Sciences, Institute for Advanced Studies (IAS), Wuhan University, Wuhan 430072, China

⁴Agricultural Research Service, US Department of Agriculture, National Center for Agricultural Utilization Research, Peoria, IL, USA

⁵These authors contributed equally to this paper.

*For Correspondence: **Wenqing Chen**, School of Pharmaceutical Sciences, Wuhan University, Wuhan 430071, China. E-mail: wqchen@whu.edu.cn

*For Co-correspondence: **Xiangzhao Mao**, College of Food Science and Engineering, Ocean University of China, Qingdao 266003, China. E-mail: xzhmao@ouc.edu.cn

Table of Contents

1. Supplementary Figures.

Figure S1. Verification and genetic organizations of the FOR-A gene clusters

Figure S2. Enzymatic characterization of ForK/PrfK as a lysine N^6 -monooxygenase

Figure S3. *In vivo* functional characterization of *forT/prfT* in FOR-A/PRF-A biosynthesis

Figure S4. Biochemical characterization of PrfT/ForT as a β -RFA-P synthase-like enzyme

Figure S5. NMR analysis of compound **19**

Figure S6. Target-oriented genome mining of the biosynthetic gene clusters for *the* potential C-nucleoside antibiotics

2. Supplementary Tables.

Table S1. Deduced functions of the open reading frames in the *foc* and *for-cof* gene clusters

Table S2. Crystal data and structure refinement for ADCP and DCOP

Table S3. NMR data for **19**

2. Supplementary Figures.

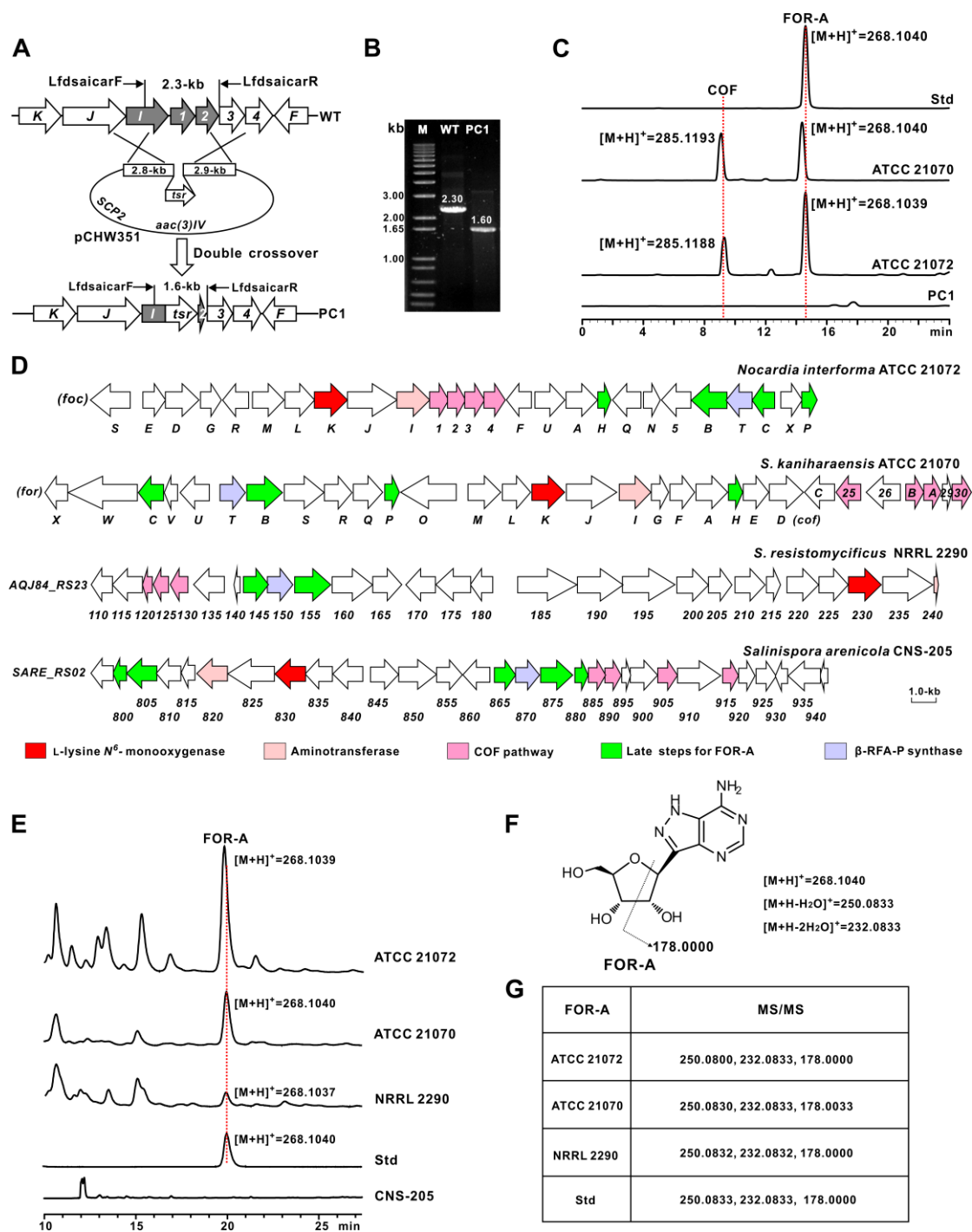


Figure S1. Verification and genetic organizations of the FOR-A gene clusters

(A) Schematic illustration for the construction of *Nocardia interformae* PC1 mutant. (B) PCR identification of the mutant PC1. M, 1 kb plus DNA ladder; WT, PCR product using genomic DNA of *Nocardia interformae* ATCC 21072 as template; PC1, PCR product using genomic DNA of *Nocardia interformae* PC1 as template. (C) LC-HRMS

analysis of the target metabolites produced by related strains. Std, the authentic standard of FOR-A; ATCC 21070, the metabolites (COF and FOR-A) produced by *S. kaniharaensis* ATCC 21070 as positive control; ATCC 21072, the metabolites (COF and FOR-A) produced by the strains of *Nocardia interforma* ATCC 21072; PC1, the metabolites produced by the strains of *Nocardia interforma* PC1 mutant. The main fragment ions of the metabolite COF produced by *S. kaniharaensis* ATCC 21070 are as follows: m/z 267.1104, 152.9898, and 134.9333. (D) Genetic organizations of the FOR-A and COF (-related) gene clusters. The *foc* gene cluster is from *Nocardia interforma* ATCC 21072, the *for-cof* gene cluster is from *S. kaniharaensis* ATCC 21070, the *AQJ84_RS23* gene cluster is from *S. resistomyces* NRRL 2290, and the *SARE_RS02* gene cluster (potentially responsible for FOR-A and COF biosynthesis) is from *Salinispora arenicola* CNS-205. (E) HPLC analysis of the metabolites individually produced by the strains of ATCC 21072, ATCC 21070, NRRL 2290, and CNS-205. ATCC 21072, the metabolites produced by *Nocardia interforma* ATCC 21072; ATCC 21070, the metabolites produced by *S. kaniharaensis* ATCC 21070; NRRL 2290, the metabolites of *S. resistomyces* NRRL 2290; Std, the authentic standard of FOR-A; CNS-205, the metabolites produced by *Salinispora arenicola* CNS-205. The fermentation conditions of *S. resistomyces* NRRL 2290 and *Salinispora arenicola* CNS-205 are identical to that of *S. kaniharaensis* ATCC 21070. (F) The fragmentation pattern of the FOR-A authentic standard. (G) The target main MS/MS fragments of the metabolite FOR-A produced by related strains. ATCC 21072, *Nocardia interforma* ATCC 21072; ATCC 21070, *S. kaniharaensis* ATCC 21070; NRRL 2290, *S. resistomyces* NRRL 2290; Std, the authentic standard of FOR-A. As we expected to emphasize the core conserved genes (whose functions could be easily assigned) related with the present manuscript, as a result, we just marked them with color, and left other genes unmarked (blank).

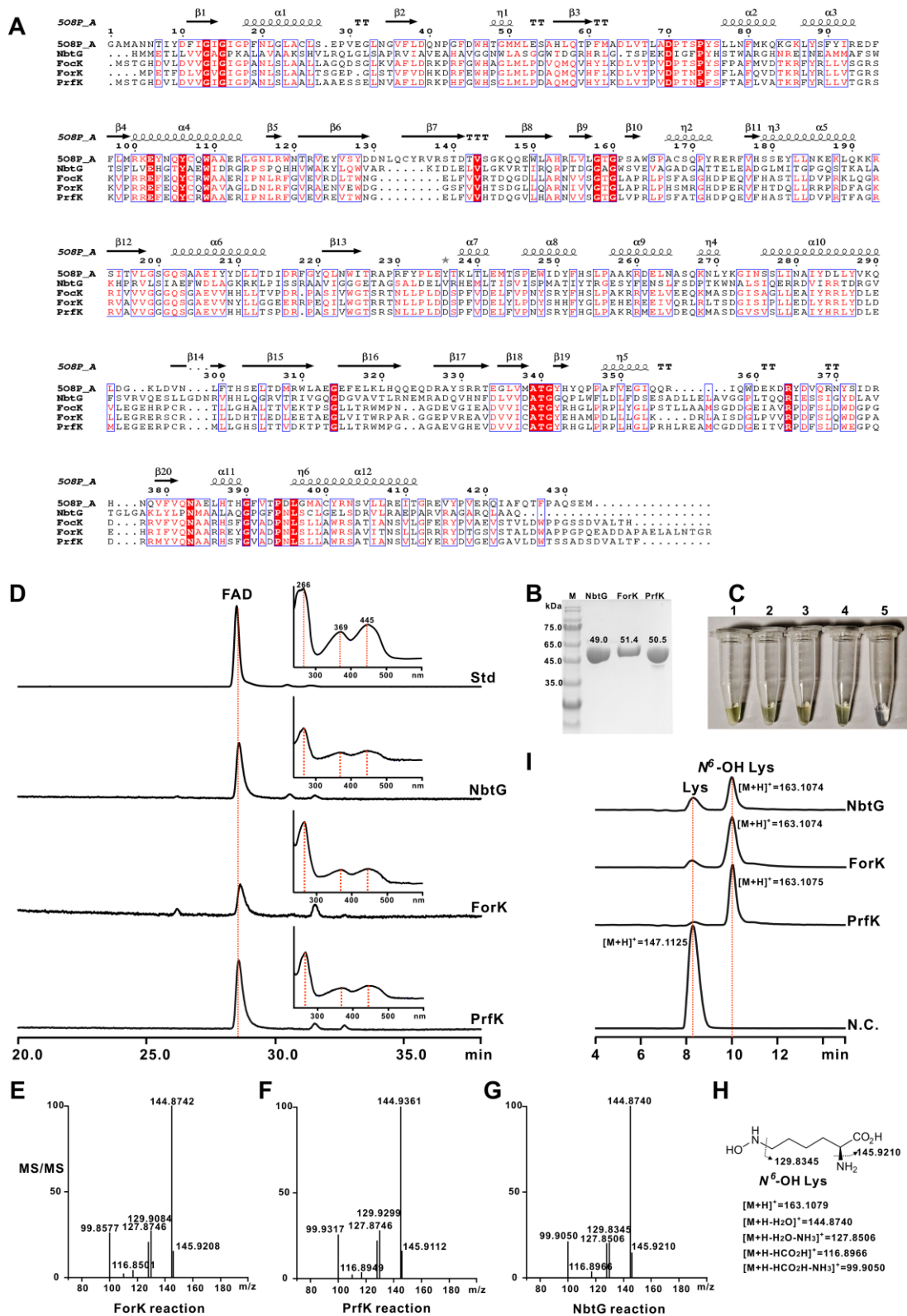


Figure S2. Enzymatic characterization of ForK/PrfK as a lysine N^6 -monooxygenase

(A) Sequence alignment of FocK/ForK/PrfK with its homologs was performed using ESPript 3.0 (<http://esprict.ibcp.fr/ESPript/ESPript/>). Secondary structure of 5O8P-A

(PDB: 5O8P-A), from *Erwinia amylovora* CFBP1430, is shown on the top. NbtG (GenBank: BAD55606) is from *Nocardia farcinica* IFM 10152. Helices represent alpha-helices and arrows signify beta-strands. Amino acids (white) in red background with blue boxes means conserved regions, while the amino acids (red) within blue boxes means less conserved regions. (B) SDS-PAGE analysis of NbtG, ForK, and PrfK. The predicted molecular weights of NbtG (49.0 kDa), ForK (51.4 kDa) and PrfK (50.5 kDa) are consistent with the individual migrations on SDS-PAGE. (C) The color of the target purified proteins. 1, the authentic FAD standard as positive control; 2, the purified NbtG; 3, the purified ForK; 4, the purified PrfK; 5, the protein stock buffer used as negative control. (D) HPLC analysis of FAD isolated from proteins. The corresponding UV spectrum was included in the related trace. Std, the authentic FAD standard; NbtG, the FAD isolated from NbtG; ForK, the FAD isolated from ForK; PrfK, the FAD in purified protein PrfK. (E-G) LC-HRMS/MS analysis of the product N^6 -OH Lys produced by the ForK/PrfK/NbtG reaction. (H) The fragmentation pattern of N^6 -OH Lys produced by the NbtG reaction. N^6 -OH Lys, N^6 -OH Lysine. (I) LC-MS analysis of the ForK/PrfK/NbtG reactions. NbtG, the reaction catalyzed by NbtG; ForK, the reaction catalyzed by ForK; PrfK, the reaction catalyzed by PrfK; N.C, the reaction without enzyme added as negative control.

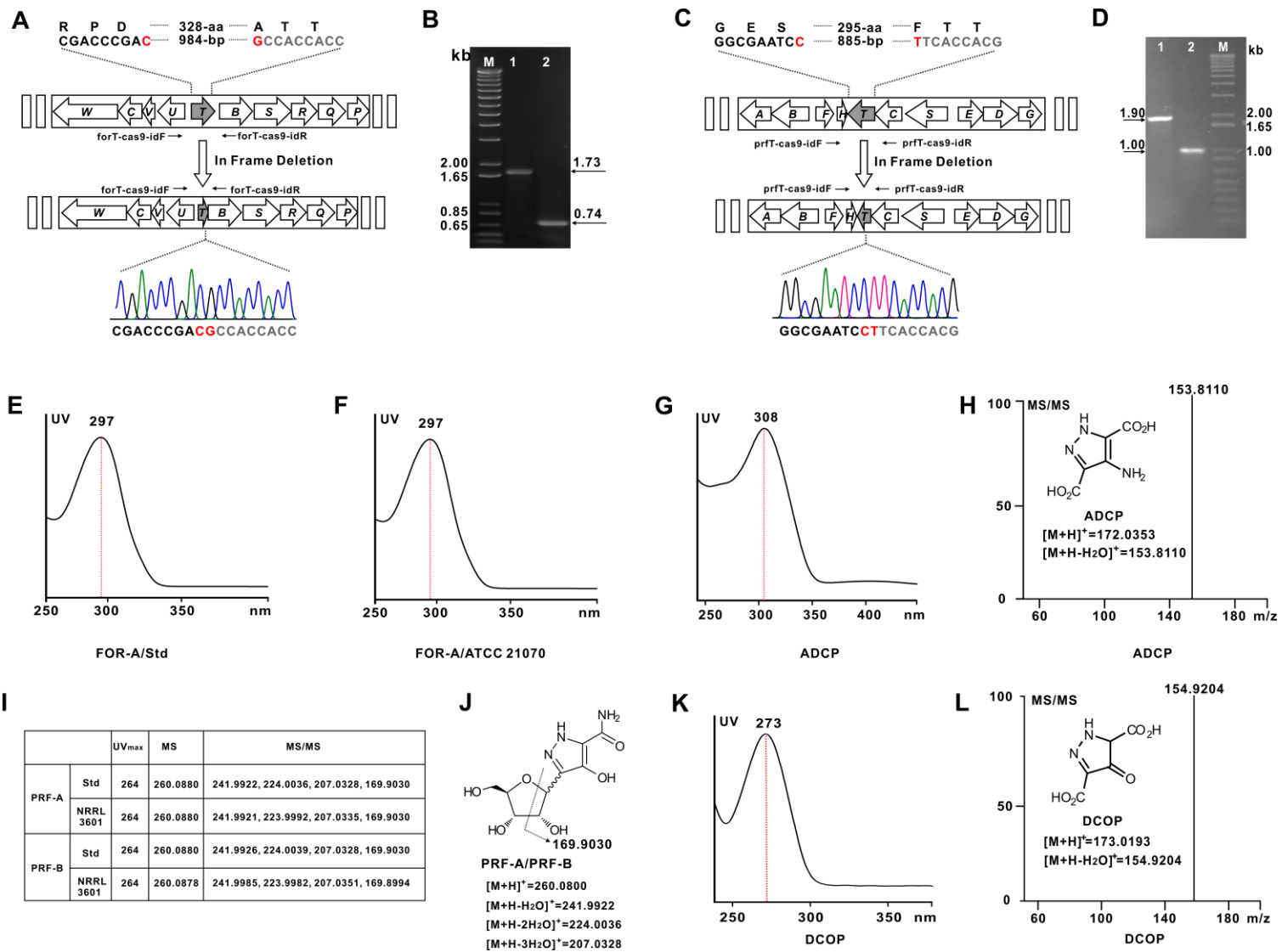


Figure S3. *In vivo* functional characterization of *forT/prfT* in FOR-A/PRF-A biosynthesis

(A) Schematic illustration for the construction of *forT* via CRISPR-Cas9 technology. (B) PCR identification of the in-frame deletion mutant. M, 1-kb Plus DNA ladder; 1, PCR product using genomic DNA of *S. kaniharaensis* ATCC 21070 as template; 2, PCR product using genomic DNA of *S. kaniharaensis* $\Delta forT$ mutant as template. (C) Schematic illustration for the construction of *prfT* via CRISPR-Cas9 technology. (D) PCR identification of the in-frame deletion mutant. M, 1-kb Plus DNA ladder; 1, PCR product using genomic DNA of *S. candidus* NRRL 3601 as template; 2, PCR product using genomic DNA of *S. candidus* $\Delta prfT$ mutant as template. (E) UV spectrum of the FOR-A standard. (F) UV spectrum of FOR-A produced by *S. kaniharaensis* ATCC 21070. (G) UV spectrum of 4-amino-3,5-dicarboxypyrazole (ADCP) produced by *S. kaniharaensis* $\Delta forT$ mutant. (H) MS/MS analysis of ADCP. (I) The UV spectrums and MS/MS fragments of PRF-A/PRF-B. Std, the authentic standard of PRF-A/PRF-B; NRRL 3601, the target metabolite PRF-A/PRF-B produced by *S. candidus* NRRL 3601. (J) The fragmentation pattern of PRF-A/PRF-B. (K) UV spectrum of 3,5-dicarboxy-4-oxo-4,5-dihydropyrazole (DCOP). (L) MS/MS analysis of DCOP. The fragmentation pattern of DCOP is also indicated in this panel.

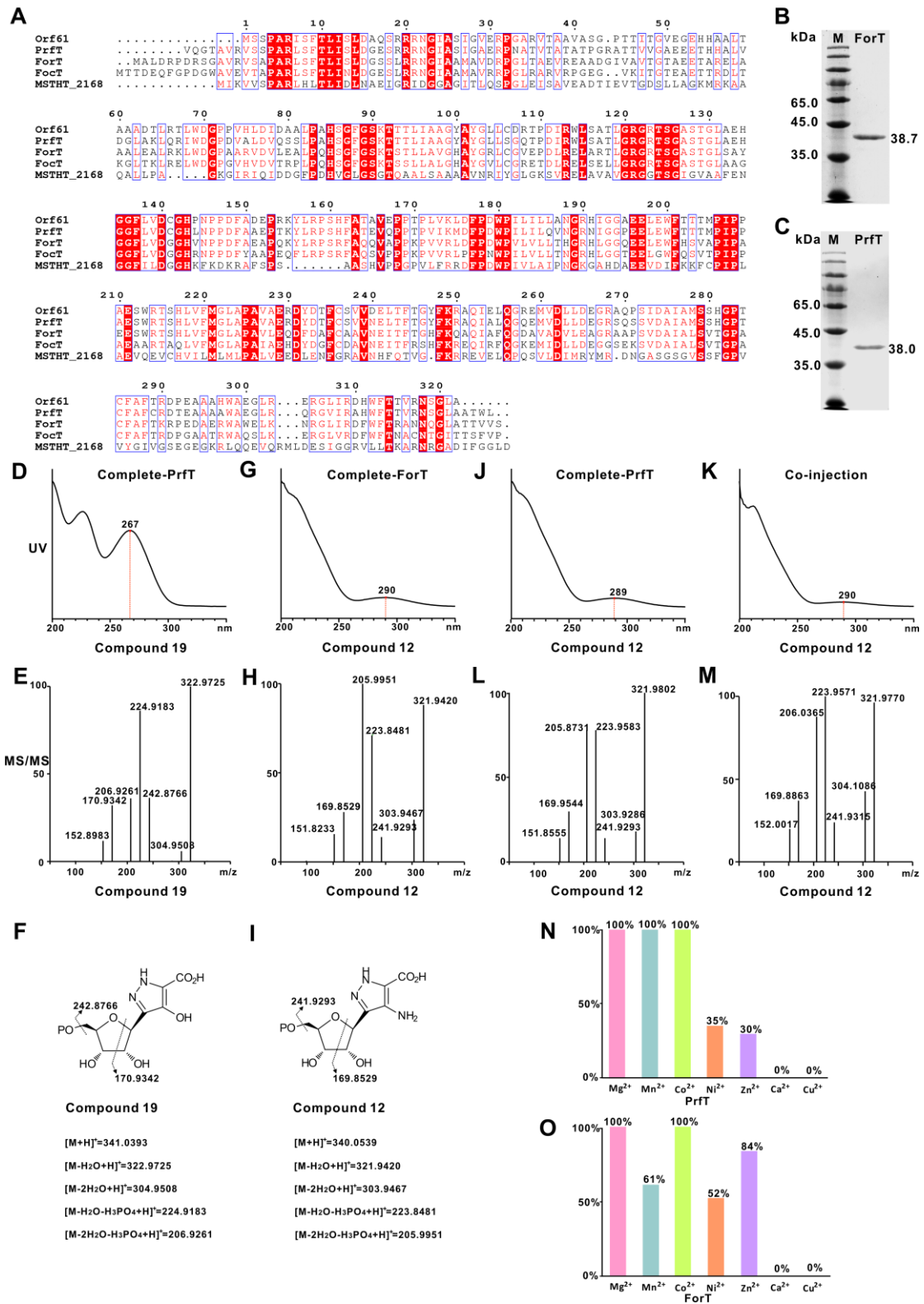


Figure S4. Biochemical characterization of PrfT/ForT as a β -RFA-P synthase-like enzyme

(A) Alignment of PrfT/ForT/FocT with relevant homologs. Orf61 (GenBank:

AJO72754.1) is from *Nocardia terpenica* IFM 0406. MSTHT_2168 (GenBank: AKB13926.1) is from *Methanosarcina thermophila* TM-1. White amino acids in red background with blue boxes means conserved regions, while the red amino acids with blue boxes means less conserved regions. (B) SDS-PAGE analysis of ForT. The predicted molecular weight of ForT (38.7 kDa) is consistent with the migration on SDS-PAGE. (C) SDS-PAGE analysis of PrfT. The predicted molecular weight of PrfT (38.0 kDa) is matched to the migration on SDS-PAGE. (D) UV spectrum of the PrfT-catalyzed product, compound **19**. (E) LC-HRMS/MS analysis of **19** produced by PrfT. (F) The fragmentation pattern of **19**. (G) UV spectrum of the ForT-catalyzed product, compound **12**. (H) LC-HRMS/MS analysis of **12** produced by ForT. (I) The fragmentation pattern of **12**. (J) UV spectrum of **12** produced by PrfT complete reaction. (K) UV spectrum of **12** (co-injection of the ForT and PrfT individual complete reactions). (L) LC-HRMS/MS analysis of **12** produced by PrfT. (M) LC-HRMS/MS analysis of **12** (co-injection of the ForT and PrfT individual complete reactions). (N) Relative efficiency of PrfT-catalyzed reaction with different divalent metal ions. (O) Relative efficiency of ForT-catalyzed reaction with different divalent metal ions.

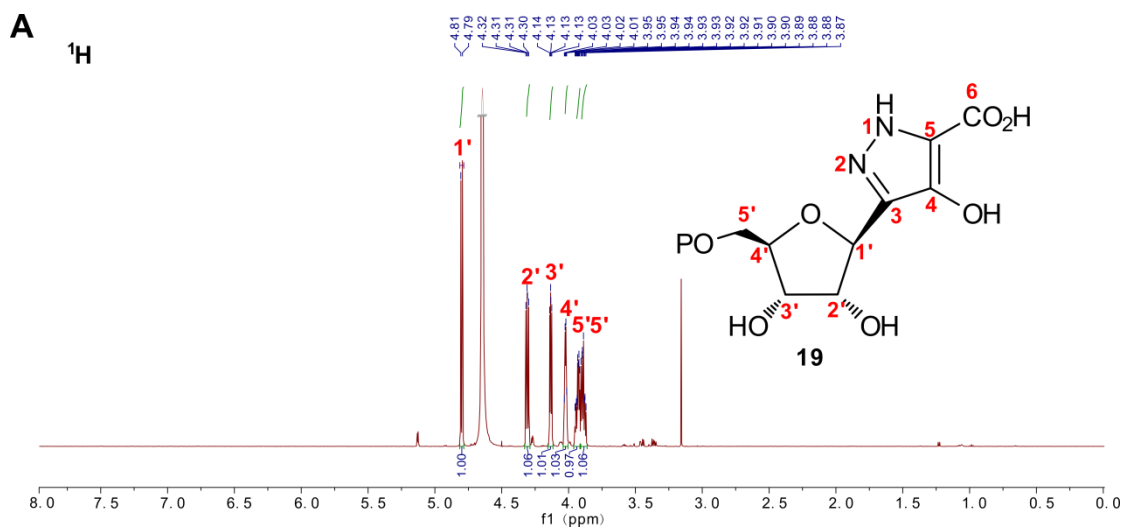


Figure S5A. ^1H NMR analysis of compound **19**

^1H NMR data of **19** (600 MHz, D_2O).

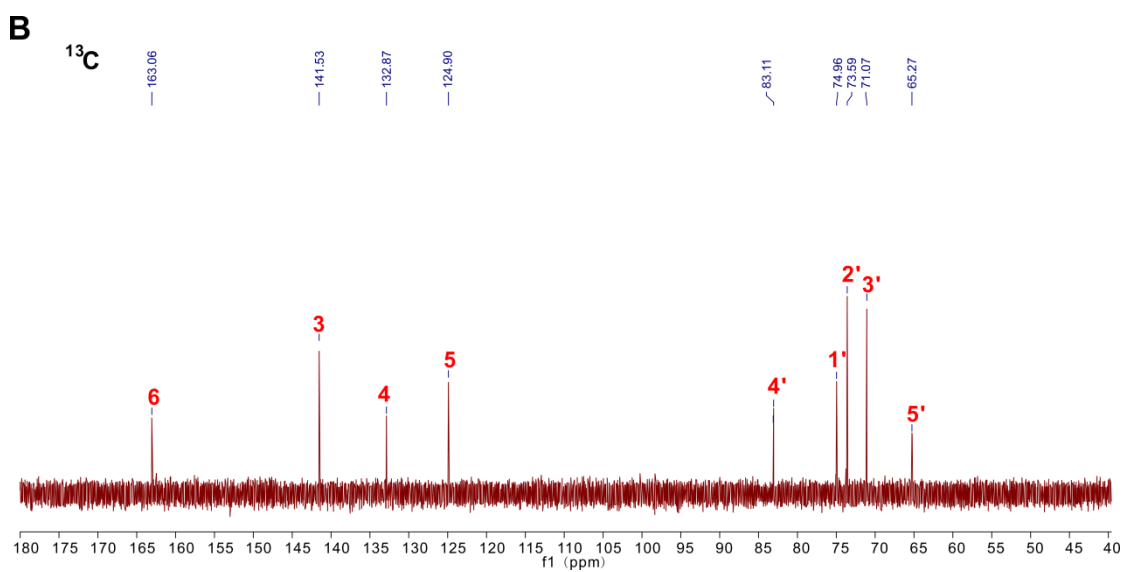


Figure S5B. ^{13}C NMR analysis of **19**

^{13}C NMR data of **19** (600 MHz, D_2O).

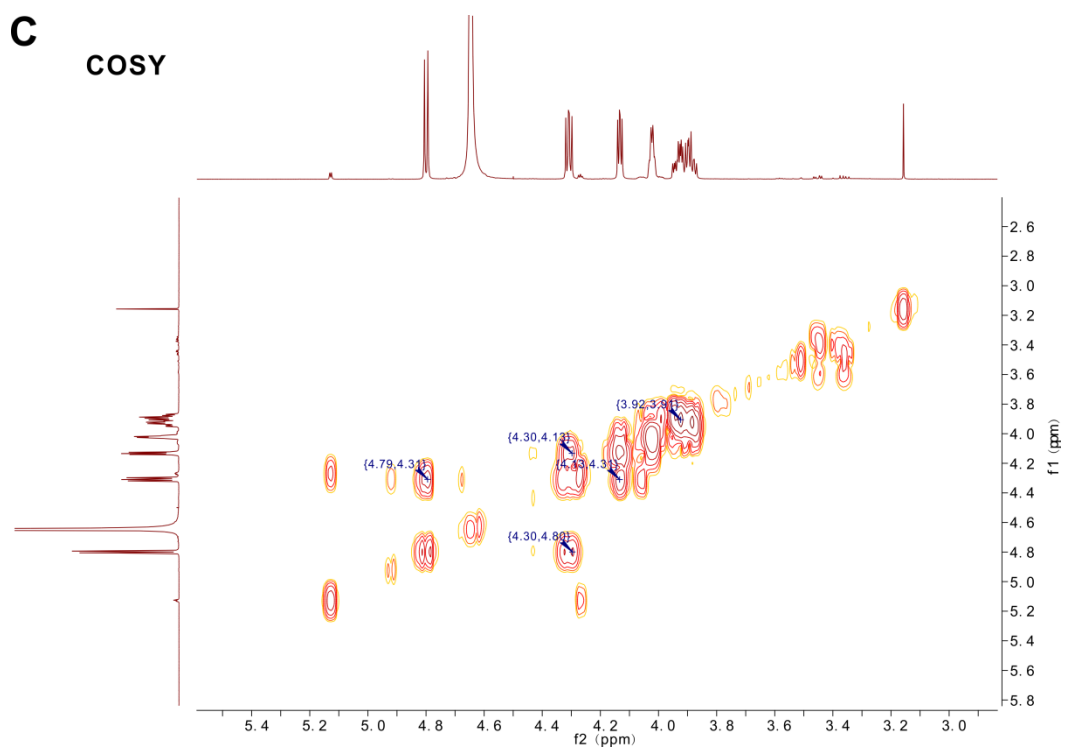


Figure S5C. COSY analysis of 19

COSY spectrum of **19** (600 MHz, D₂O).

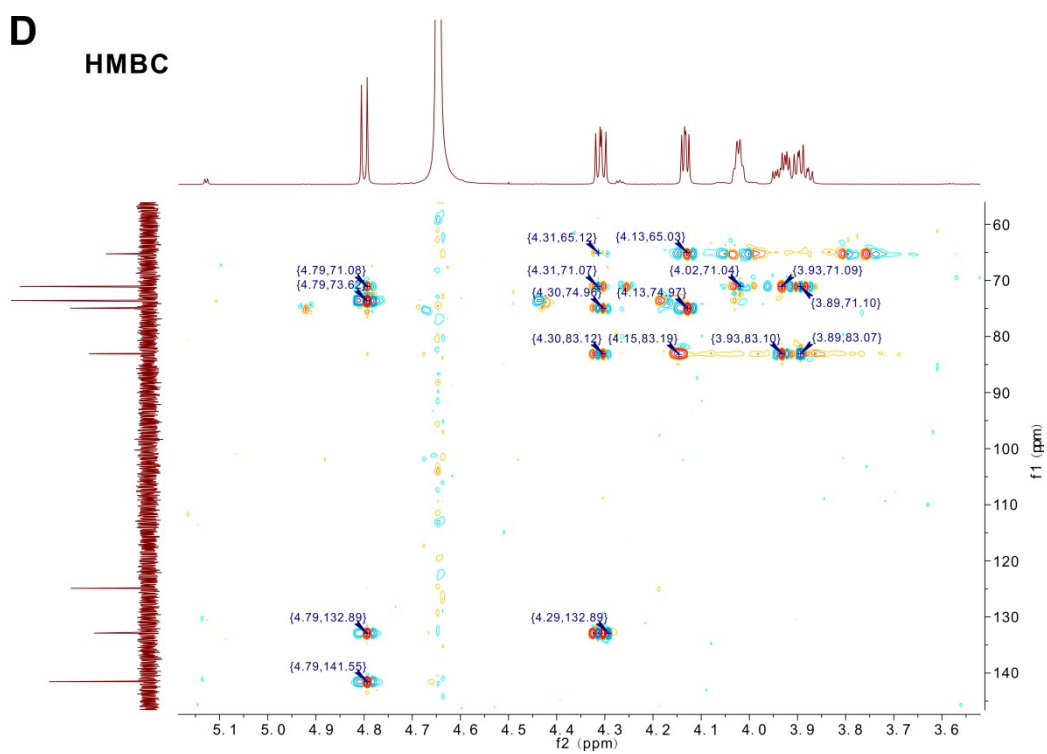


Figure S5D. HMBC analysis of 19

HMBC spectrum of **19** (600 MHz, D₂O).

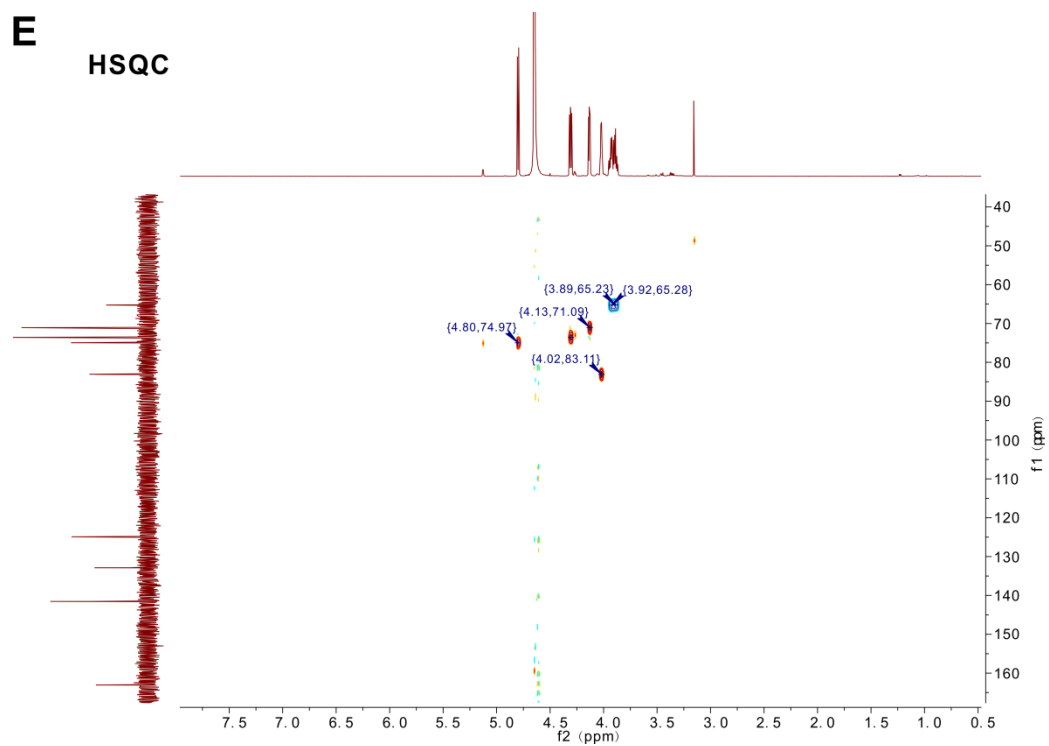


Figure S5E. HSQC analysis of **19**

HSQC spectrum of **19** (600 MHz, D_2O).

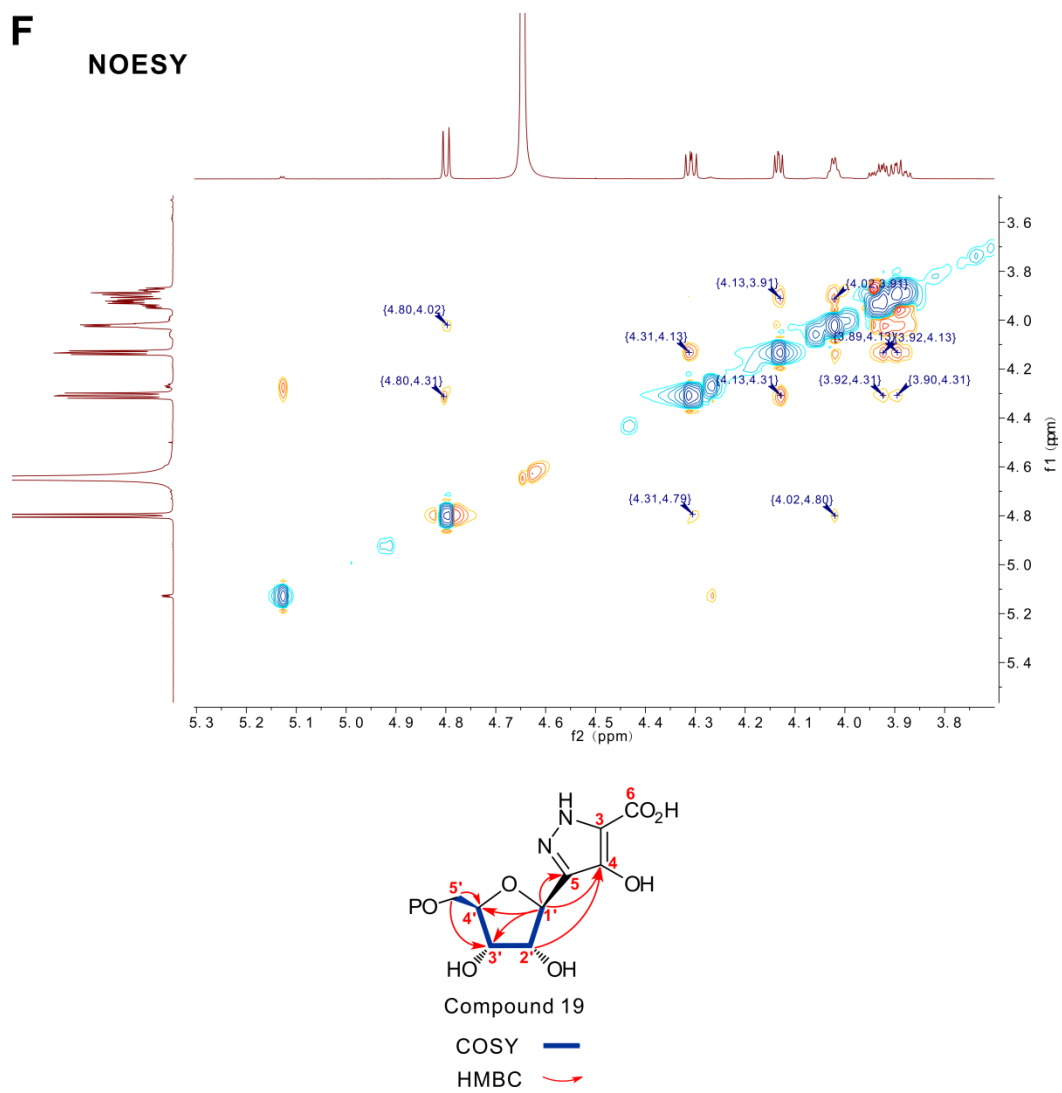


Figure S5F. NOESY analysis of 19

NOESY spectrum of **19** (600 MHz, D₂O).

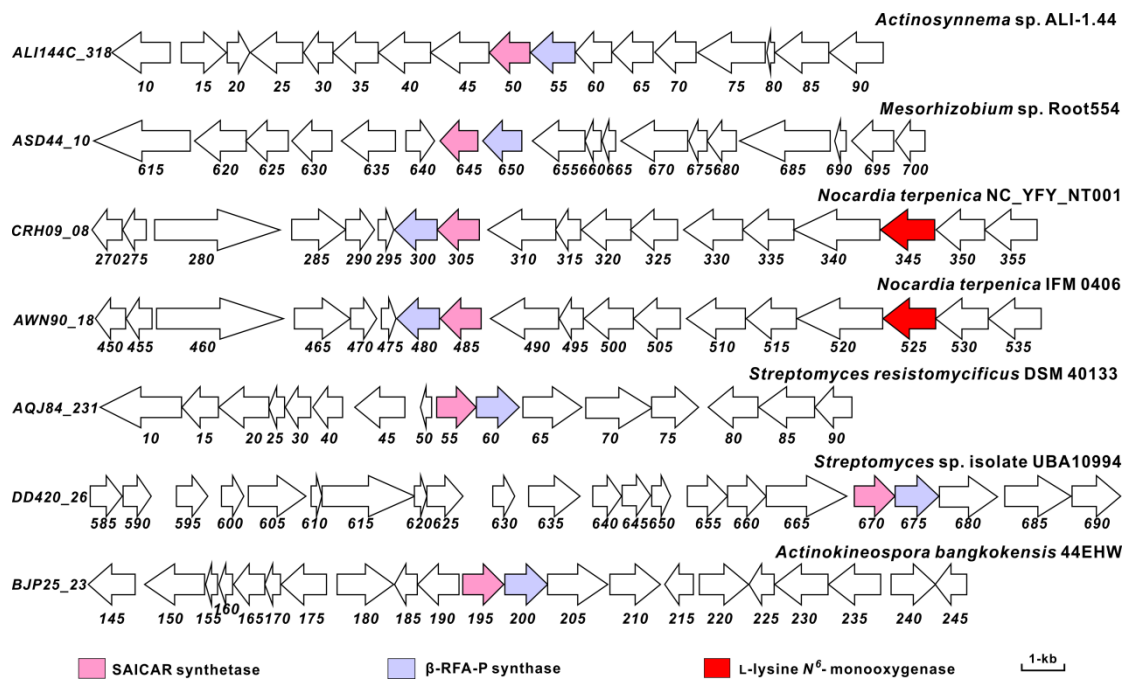


Figure S6. Target-oriented genome mining of the biosynthetic gene clusters for C-nucleoside antibiotics

The gene clusters for the potential C-nucleoside antibiotics were obtained using PrFT as enzyme probe. The proposed gene functions with highlighted colors were listed at the bottom. SAICAR, phosphoribosylaminoimidazole-succinocarboxamide.

2. Supplementary Tables.

Table S1. Deduced functions of the open reading frames in the *foc* and *for-cof* gene clusters*

Protein (aa)	Protein Function	Homolog, Origin	Identities, positives (%)	Accession no.
ForX FocX (295) (288) 55, 66	phosphofructokinase	OV450_4865, <i>Actinobacteria bacterium</i> OV450	58, 71	KPI01287
ForW (931)	FAD-binding oxidoreductase	MPHLEI_25011, <i>Mycobacterium phlei</i> RIVM601174	69, 77	EID09270
ForC FocC (333) (294) 58, 68	SAICAR synthetase	OV450_4874, <i>Actinobacteria bacterium</i> OV450	67, 75	KPI01296
ForV (185)	FMN reductase (NADPH)	BIV23_36585, <i>Streptomyces</i> sp. MUSC 1	73, 83	OIJ93994
ForU FocU (390) (407) 73, 84	phthalate 4,5-dioxygenase	ADK37_11500, <i>S. resistomycificus</i>	87, 92	KOG37051
ForT FocT (341) (341) 69, 77	C-glycosyltransferase	Orf61, <i>Nocardia terpenica</i> IFM 0406	64,75	AJO72754
ForB FocB (474) (480) 66, 80	adenylosuccinate lyase	ADL05_15155, <i>Nocardiosis</i> sp. NRRL B-16309	78, 86	KOX15518
ForS FocS (529) (528) 62, 69	fumarate reductase	CRH09_08310, <i>Nocardia terpenica</i>	64, 72	ATL71617
ForR FocR (374) (351) 49, 60	glycine/D-amino acid oxidase	Orf73, <i>Nocardia terpenica</i> IFM 0406	54, 65	AJO72766
ForQ FocQ (396) (387) 63, 75	amidohydrolase	A5717_10410, <i>Mycolicibacterium porcinum</i> ACS3670	62, 76	OCB14542
ForP FocP (187) (205) 65, 73	NUDIX hydrolase	BOQ63_27930, <i>S. viridifaciens</i> DSM 40239	82, 87	OJH68722
ForO (736)	S9 family peptidase	IX27_25180, <i>Streptomyces</i> sp. JS01	61, 70	KFK86867
ForM FocM (436) (413) 66, 79	phosphoribosylglycine amide synthetase	Orf72, <i>Nocardia terpenica</i> IFM 0406	69, 80	AJO72765
ForL FocL (384) (393) 57, 68	saccharopine dehydrogenase	Orf71, <i>Nocardia terpenica</i> IFM 0406	58, 67	AJO72764
ForK FocK (440) (436) 62, 72	L-lysine N6-monooxygenase	Spb38, <i>Streptomyces</i> sp. SoC090715LN-17	41, 55	BAW2770 2
ForJ FocJ (677) (668) 59, 68	methionine-tRNA ligase	Spb40, <i>Streptomyces</i> sp. SoC090715LN-17	36, 50	BAW2770 4
ForI FocI (423) (433)	aminotransferase	BJP25_23260, <i>Actinokineospora bangkokensis</i> 44EHW	64, 78	OLR92323

64, 78				
ForG (230)	RNA polymerase sigma-70 factor	GA0070616_0590, <i>Micromonospora nigra</i> DSM 43818	52, 67	SCL14959
ForF (333) 68, 81	FocF (334) phosphoglycerate dehydrogenase	BAL199_19391, <i>Alpha</i> <i>proteobacterium</i> BAL199	47, 63	EDP63049
ForA (423) 68, 79	FocA (426) adenylosuccinate synthetase	BJP25_23230, <i>Actinokineospora</i> <i>bangkokensis</i>	68, 77	OLR92236
ForH (196) 66, 73	FocH (187) AICAR transformylase	BJP25_23225, <i>Actinokineospora</i> <i>bangkokensis</i>	63, 72	OLR92235
ForE (342) 49, 60	FocE (312) putative monooxygenase	ASC89_11735, <i>Devosia</i> sp. Root413D1	41, 52	KQW7897 2
ForD (439) 78, 86	FocD (444) putative monooxygenase	Orf75, <i>Nocardia terpenica</i> IFM 0406	79, 86	AJO72768
CofC (404) 61, 74	Foc5 (394) MFS transporter	NbrT6, <i>Nocardia terpenica</i> IFM 0406	60, 74	AJO72759
Cof25 (335)	metal dependent phosphohydrolase	PenF, <i>S. antibioticus</i>	37, 49	AKA87335
Cof26 (452)	GntR family transcriptional regulator	SAZ_09915, <i>S. albulus</i> ZPM	71, 77	AKA02703
CofB (239) 61, 73	Foc1 (242) SAICAR synthetase	PenC, <i>S. antibioticus</i>	62, 75	AKA87338
CofA (233) 52, 66	Foc2 (231) short chain dehydrognase	PenB, <i>S. antibioticus</i>	49, 63	AKA87339
Cof29 (140)	hypothetic protein	GA0070616_0574, <i>Micromonospora nigra</i> DSM 43818	41, 52	SCL14901
Cof30 (264) 33, 42	Foc3 (257) HAD family phosphatase	Acel_2103, <i>Acidothermus</i> <i>cellulolyticus</i> 11B	35, 49	ABK53875
FocG (275)	transcriptional regulator	BJP25_23295, <i>Actinokineospora</i> <i>bangkokensis</i>	60,70	OLR92324
Foc4 (290)	ATP phosphoribosyltransf erase	BJP25_23245, <i>Actinokineospora</i> <i>bangkokensis</i>	63,73	OLR92239
FocN (233)	putative HAD superfamily hydrolase	BJP25_23215, <i>Actinokineospora</i> <i>bangkokensis</i>	71, 81	OLR92233

*The *for* genes are from *S. kaniharaensis* ATCC 21070, while the *foc* genes are from *Nocardia interforma* ATCC 21072. Taking ForX as example, it shares 55% identities (66% positives) to FocX, and it shows 58% identities (71% positives) to OV450_4865 of *Actinobacteria bacterium* OV450.

Table S2A. Crystal data and structure refinement for ADCP

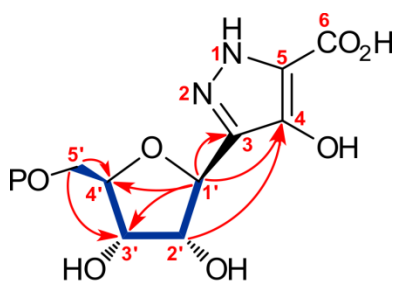
Identification code	x_a
Empirical formula	C5 H7 N3 O5
Formula weight	189.14
Temperature	200(2) K
Wavelength	1.54178 Å
Crystal system	Monoclinic
Space group	P2 ₁ /n
Unit cell dimensions	a = 12.8777(6) Å α = 90° b = 4.2342(2) Å β = 103.553(2)°.. c = 14.2564(6) Å γ = 90°..
Volume	755.71(6) Å ³
Z	4
Density (calculated)	1.662 Mg/m ³
Absorption coefficient	1.320 mm ⁻¹
F(000)	392
Crystal size	0.140 x 0.120 x 0.060 mm ³
Theta range for data collection	4.167 to 63.017°.
Index ranges	-14<=h<=12, -4<=k<=3, -16<=l<=16
Reflections collected	4581
Independent reflections	1213 [R(int) = 0.0217]
Completeness to theta = 63.017°	99.7 %
Absorption correction	Semi-empirical from equivalents
Max. and min. transmission	0.7530 and 0.5406
Refinement method	Full-matrix least-squares on F ²
Data / restraints / parameters	1213 / 0 / 146
Goodness-of-fit on F ²	1.108
Final R indices [I>2sigma(I)]	R1 = 0.0288, wR2 = 0.0729
R indices (all data)	R1 = 0.0300, wR2 = 0.0737
Extinction coefficient	n/a
Largest diff. peak and hole	0.179 and -0.240 e.Å ⁻³

Table S2B. Crystal data and structure refinement for DCOP

Identification code	x_a
Empirical formula	C5 H3 N2 O7
Formula weight	203.09
Temperature	296(2) K
Wavelength	1.54178 Å
Crystal system	Triclinic
Space group	P-1
Unit cell dimensions	a = 7.9342(7) Å α = 98.253(5)° b = 8.2041(8) Å β = 95.051(4)° c = 12.7997(12) Å γ = 97.782(4)°.
Volume	812.03(13) Å ³
Z	4
Density (calculated)	1.661 Mg/m ³
Absorption coefficient	1.441 mm ⁻¹
F(000)	412
Crystal size	0.120 x 0.080 x 0.060 mm ³
Theta range for data collection	3.510 to 61.155°.
Index ranges	-9<=h<=9, -9<=k<=9, -10<=l<=14
Reflections collected	5416
Independent reflections	2365 [R(int) = 0.0464]
Completeness to theta = 61.155°	95.0 %
Absorption correction	Semi-empirical from equivalents
Max. and min. transmission	0.7522 and 0.5691
Refinement method	Full-matrix least-squares on F ²
Data / restraints / parameters	2365 / 1 / 260
Goodness-of-fit on F ²	1.135
Final R indices [I>2sigma(I)]	R1 = 0.0991, wR2 = 0.2719
R indices (all data)	R1 = 0.1026, wR2 = 0.2753
Extinction coefficient	n/a
Largest diff. peak and hole	0.884 and -0.573 e.Å ⁻³

Table S3. NMR data for 19

	δ_C	$\delta_H, J(\text{Hz})$
C-3	141.53	
C-4	132.87	
C-5	124.90	
C-6	163.06	
C-1'	74.96	4.80, d, $J=7.24$
C-2'	73.59	4.31, dd, $J=5.40, 7.24$
C-3'	71.07	4.13, dd, $J=3.93, 5.22$
C-4'	83.11	4.02, m
C-5'	65.27	3.93, dd, $J=5.84, 11.37$
		3.89, dd, $J=6.45, 11.32$



19

COSY ———
 HMBC ———→