

## Supplementary Materials for

### Structure of *Drosophila melanogaster* ARC1 reveals a repurposed molecule with characteristics of retroviral Gag

Matthew A. Cottee, Suzanne C. Letham, George R. Young, Jonathan P. Stoye, Ian A. Taylor\*

\*Corresponding author. Email: [ian.taylor@crick.ac.uk](mailto:ian.taylor@crick.ac.uk)

Published 1 January 2020, *Sci. Adv.* **6**, eaay6354 (2020)

DOI: [10.1126/sciadv.aay6354](https://doi.org/10.1126/sciadv.aay6354)

#### **This PDF file includes:**

Table S1. dARC1 CA statistics of data collection, phasing, and refinement.

Table S2. Hydrodynamic parameters of dARC1 CA.

Fig. S1. Crystal structures of dARC1 CA.

Fig. S2. dARC1 CA dimer.

Fig. S3. Comparison of dARC1 and mam-ARC CA-NTDs.

Fig. S4. Comparison of dARC1 and Ty3 CA.

**Table S1. dARC1 CA statistics of data collection, phasing, and refinement.**

	dARC1 CA oP (Se-Met, peak)	dARC1 CA oP (Se-Met, Irem)	dARC1 CA oP (Se-Met, Irem STARANISO)	dARC1 CA hP (Se- Met, anom)
<b>Data collection</b>				
Space group	P2 <sub>1</sub> 2 <sub>1</sub> 2 <sub>1</sub>	P2 <sub>1</sub> 2 <sub>1</sub> 2 <sub>1</sub>	P2 <sub>1</sub> 2 <sub>1</sub> 2 <sub>1</sub>	P6 <sub>1</sub> 22
Cell dimensions				
a, b, c (Å)	62.43, 70.715,	62.35, 70.39,	62.33, 70.39,	62.18, 62.18,
α, β, γ (°)	142.07 90, 90, 90	142.84 90, 90, 90	142.82 90, 90, 90	401.02 90, 90, 120
Wavelength (Å)	0.9794	0.9840	0.9840	0.9794
Resolution (Å)	47.36-2.06 (2.12-2.06)	142.84-1.70 (1.74-1.70)	71.41-1.55 (1.70-1.55)	57.29-2.14 (2.20- 2.14)
Anisotropic diff. limits (a*, b*, c*)	-	-	1.55, 2.40, 1.62	
Unique reflections	39409 (2872) <sup>†</sup>	70026 (5125)	48232 (2412)	26972 (1886)
R <sub>meas</sub> (%)	21.5 (254.6) <sup>‡</sup>	13.5 (105.8)	9.3 (150.1)	12.2 (111.0) <sup>‡</sup>
R <sub>pim</sub> (%)	8.2 (106.1) <sup>‡</sup>	4.3 (33.9)	2.9 (47.8)	2.5 (31.6) <sup>‡</sup>
CC <sub>1/2</sub>	0.998 (0.595)	0.998 (0.568)	0.999 (0.557)	1.000 (0.456)
I/σ(I)	10.3 (1.1)	9.1 (1.5)	16.4 (1.5)	22.0 (2.7)
Completeness (spherical, %)	99.5 (95.4)	100.0 (100.0)	52.9 (11.5)	99.9 (99.3)
Completeness (ellipsoidal, %)	-	-	84.1 (66.0)	-
Multiplicity	12.9 (10.5)	9.6 (9.5)	10.0 (9.7)	40.8 (22.4)
Anom Multiplicity	6.8 (5.5)	-	-	22.4 (11.8)
<b>Phasing</b>				
Significant anomalous signal to: (Aimless)	2.82 Å	-		2.59 Å
No. sites (found/expected)	7/4	-		10/4
Mean FOM (Phenix)	0.28	-		0.324
Local rms density corr (Phenix)	0.89			0.83
<b>Refinement</b>				
Resolution (Å)			71.41-1.70 (1.74-1.70)	53.86-2.30 (2.36- 2.30)
refl working/free			41243/2241 (1680/69)	36685/1844 <sup>§</sup> (2798/121)
R <sub>work</sub> /R <sub>free</sub> /Test set size (%)			18.1/21.3/4.9 (25.2/28.2/3.9)	24.2/27.8/4.8 (31.6/31.8/4.1)
<i>No residues/atoms</i>				
Protein			328/2730	321/2658
Ligands			5	
Water			391	20
<i>B-factors (Å<sup>2</sup>)</i>				
Wilson			19.74	44.15
Protein			23.40	66.04
Ligands			33.17	
Water			31.52	46.10
Overall			24.43	65.89
<i>Geometry</i>				
Bond lengths RMSD (Å)			0.010	0.004
Bond angles RMSD (deg)			1.71	0.697
Ramachandran Outliers			0 %	0 %
Ramachandran Favoured			99.38 %	97.79 %
Molprobit score (N number/percentile)			0.92 (9248/100 <sup>th</sup> )	1.44 (8909/99 <sup>th</sup> )

<sup>†</sup>Values in parenthesis refer to the highest resolution shell, <sup>‡</sup>Values quoted for within I<sup>†</sup>/I, <sup>§</sup>Friedel pairs separated.

**Table S2. Hydrodynamic parameters of dARC1 CA.**

<b>Protein</b>	dARC1 CA
<b>Protein Parameter</b>	
$v$ (mL.g <sup>-1</sup> )	0.731
$\rho$ (g.mL <sup>-1</sup> )	1.005
<sup>a</sup> $M_r$	19,633
$\epsilon_{280}$ (M <sup>-1</sup> cm <sup>-1</sup> )	23,600
<b>Sed velocity</b>	
$C_{\text{range}}$ ( $\mu\text{M}$ )	25-100
<sup>b</sup> $S_{20,w}$ ( $\times 10^{13}$ ) sec	2.92 $\pm$ 0.03
<sup>c</sup> $S_{20,w}$ (calc) ( $\times 10^{13}$ ) sec	2.93
<sup>d</sup> $D_{20,w}$ ( $\times 10^7$ ) cm <sup>2</sup> sec <sup>-1</sup>	6.95 $\pm$ 0.05
<sup>e</sup> $D_{20,w}$ (calc) ( $\times 10^7$ ) cm <sup>2</sup> sec <sup>-1</sup>	6.77
<sup>f</sup> $M_w$ (S/D) kD	38.2 $\pm$ 1
<sup>g</sup> $f/f_0$ ( $S_{20,w}$ )	1.41
<sup>h</sup> $f/f_0$ ( $D_{20,w}$ )	1.37
<sup>i</sup> $M_w$ C(S) kD	38.7 $\pm$ 3
<sup>j</sup> $f/f_0$ C(S)	1.41
<sup>k</sup> RMSD C(S)	0.005 – 0.008
<b>Sed eqm</b>	
$C_{\text{range}}$ ( $\mu\text{M}$ )	25-70
<sup>l</sup> $M_w$ kD	38.9 $\pm$ 1
<sup>m</sup> RMSD	0.002 – 0.004
<sup>n</sup> $\chi^2$	0.31

<sup>a</sup>Molar mass calculated from the protein sequence after removal of the C-terminal His-tag

<sup>b</sup>Sedimentation coefficient at standard conditions (water @ 20° C) determined by combining all data from discrete component and C(S) analysis.

<sup>c</sup>Sedimentation coefficient calculated using HYDROPRO-10 from the co-ordinates of the dARC1 CA crystal structure using an atomic element radius (AER) of 2.84 Å.

<sup>d</sup>Translational diffusion coefficient for dARC1 CA determined from combined data from discrete component analysis.

<sup>e</sup>Translational diffusion coefficient calculated from the crystal structure using HYDROPRO-10, AER = 2.84 Å.

<sup>f</sup>The weight averaged molecular weight from discrete component analysis. Error is the range of three measurements

<sup>g</sup>The frictional ratio calculated from  $S_{20,w}$  from discrete component analysis.

<sup>h</sup>The frictional ratio calculated from  $D_{20,w}$  from discrete component analysis.

<sup>i</sup>The weight averaged molecular weight from C(S) analysis. Error is the range of three measurements

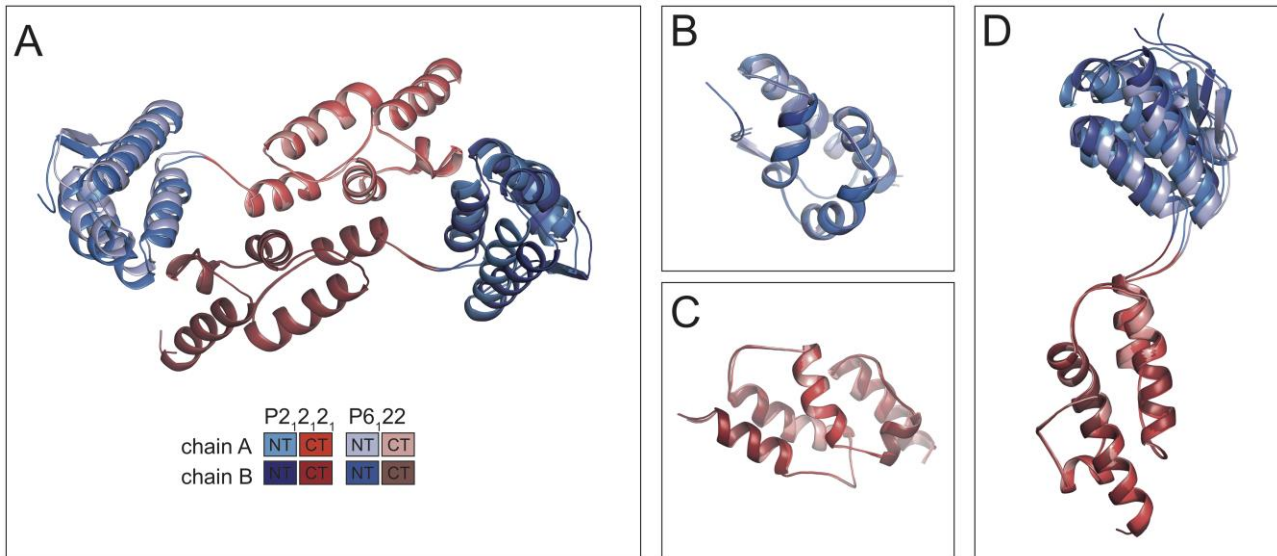
<sup>j</sup>The weight-averaged frictional ratio from the best fit C(S) distribution function.

<sup>k</sup>The range of the rms deviations observed when data were fitted using a continuous sedimentation coefficient distribution model.

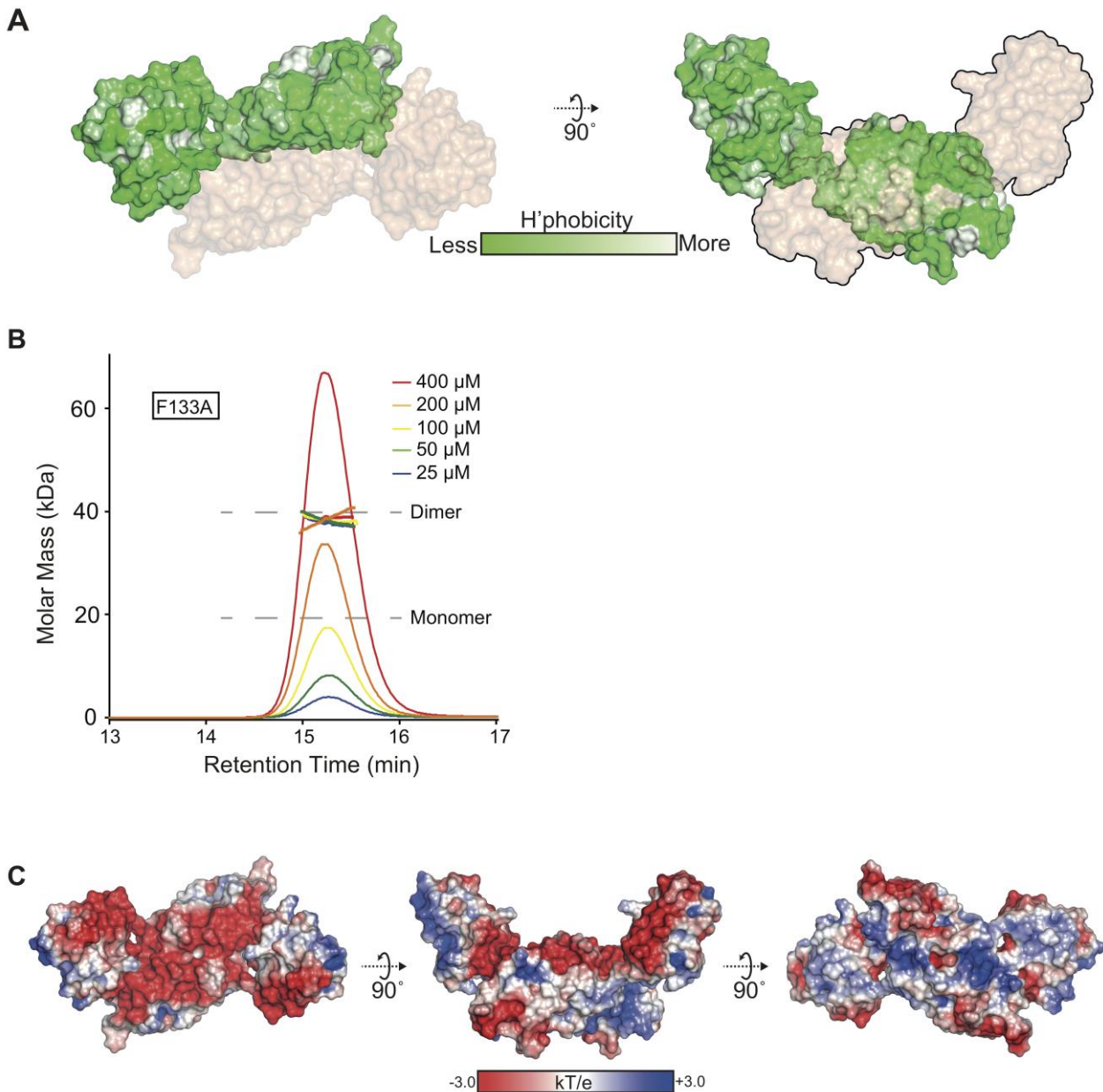
<sup>l</sup>The weight averaged molecular weight from Global SE analysis. The error is the range of molecular weights observed for each multi-speed sample when fitted individually.

<sup>m</sup>The range of the rms deviations observed for each multi-speed sample when fitted individually.

<sup>n</sup>The reduced chi-squared for the global fitting of all multispeed data.

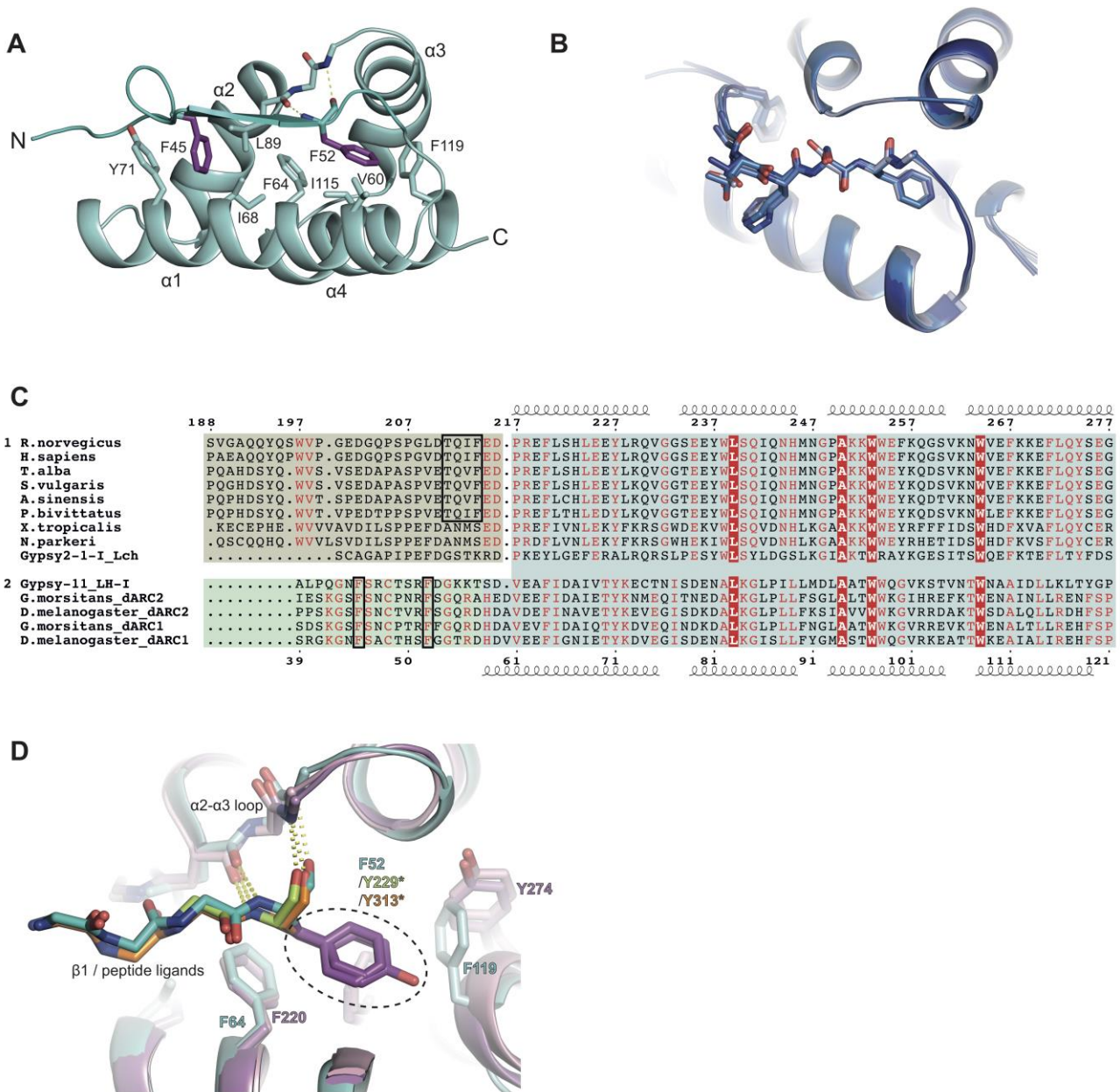


**Fig. S1. Crystal structures of dARC1 CA.** (A) Structural alignment of orthorhombic (P2<sub>1</sub>2<sub>1</sub>2<sub>1</sub>) and hexagonal (P6<sub>1</sub>22) crystal forms of dARC1 CA that contain an almost identical dimer. The ASUs are shown in cartoon representation, aligned using the dimeric CTDs (RMSD of 0.247 Å across 133 C $\alpha$  pairs). NTDs are coloured in blue shades, CTDs in red shades, according to the legend. (B) Alignment of all four NTDs only, average RMSD 0.34  $\pm$  0.09 Å over 66 $\pm$ 2 Cas. (C) Alignment of all four CTDs only, average rmsd 0.32  $\pm$  0.12 Å over 74 $\pm$ 4 Cas. (D) Alignment of all four chains through the CTD, highlighting that structural differences between chains are due to slight flexibility of the NTD-CTD linker region.



**Fig. S2. dARC1 CA dimer.** (A) Surface representation of dARC1 CA-CTD dimer displaying the distribution of surface hydrophobicity/hydrophilicity calculated using the pymol script ([https://pymolwiki.org/index.php/Color\\_h](https://pymolwiki.org/index.php/Color_h)). Greater Hydrophilicity is represented by darker green shading, the lighter and non-coloured regions represent the most hydrophobic areas of the molecule. The orientation in the left- and right-hand panels is the same as in **Fig. 1A** with monomer A displaying surface hydrophobicity, and monomer B shown as a wheat surface. (B) SEC-MALLS analysis of dARC1 CA (F133A). The sample loading concentrations were 400  $\mu$ M (8 mg/mL) (red), 200  $\mu$ M (4 mg/mL)

(orange), 100  $\mu\text{M}$  (2 mg/mL) (yellow), 50  $\mu\text{M}$  (1 mg/mL) (green) and 25  $\mu\text{M}$  (0.5 mg/mL) (blue). The differential refractive index (dRI) is plotted against column retention time and the molar mass, determined at 1-second intervals throughout the elution of each peak, is plotted as points. The dARC1 CA monomer and dimer molecular mass is indicated with the grey dashed lines. **(C)** Electrostatic surface potential of the dARC1 CA dimer, calculated with APBS. The model was modified to replace SeMet with native Met residues. The view in the left-hand and centre panels is the same as in **A**. The right-hand panel is a view from the C-terminal “underside” of the dimer.

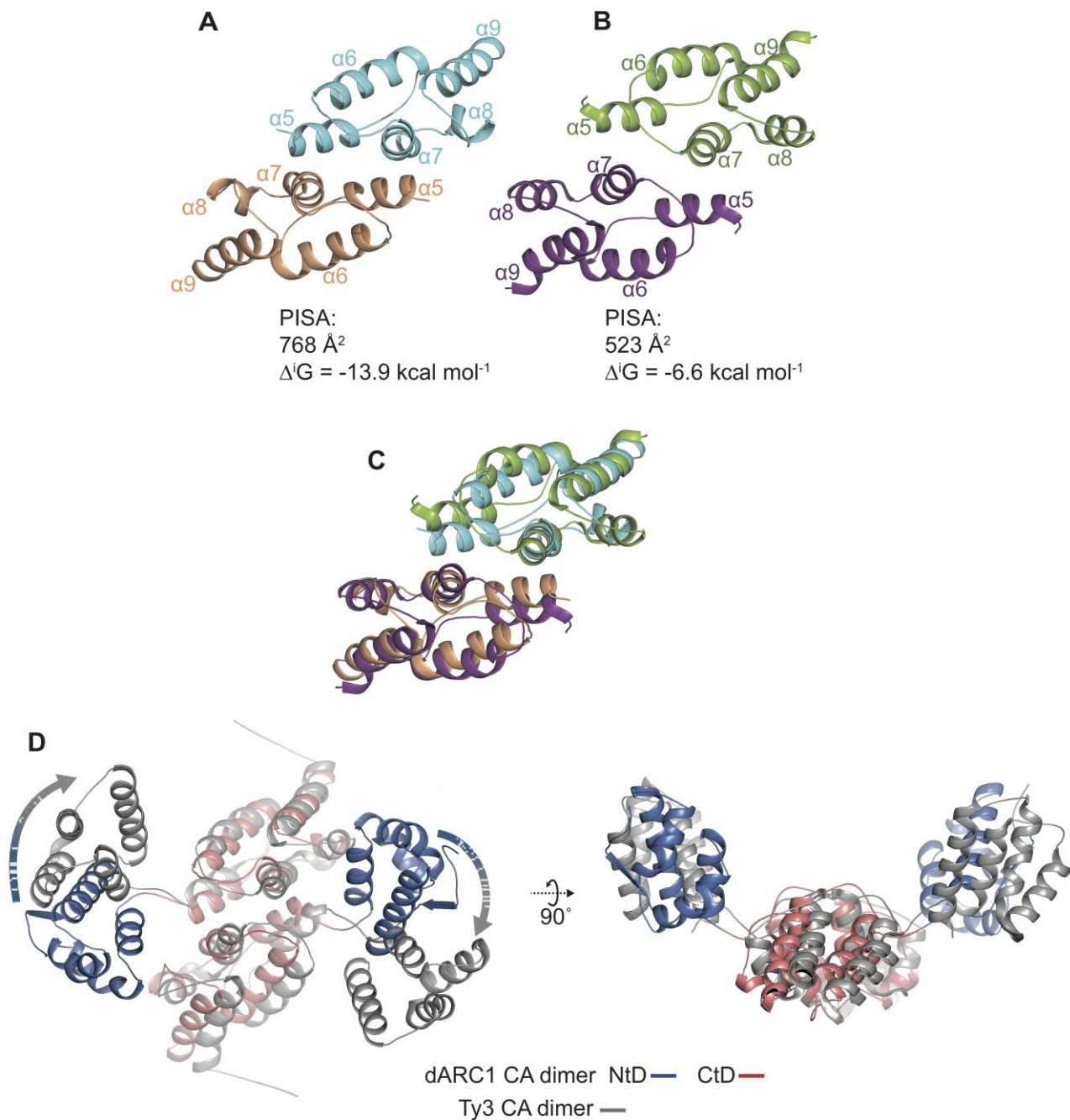


**Fig. S3. Comparison of dARC1 and mam-ARC CA-NTDs. (A)** The backbone of the dARC1 CA-NTD is shown in cartoon representation in cyan. Residues that contribute to the hydrophobic interface between the core domain, and the native N-terminal strand are shown in stick representation, coloured by atom type, with the conserved aromatic residues F45 and F52 that are buried in the interface coloured purple. The main chain hydrogen bonding interactions between the F52 backbone amide and carbonyl with the carbonyl of L89 and the amide of Y91 are shown as dashed lines. **(B)** Overlay of the four dARC1 CA NTDs from the two crystal structures, shown in different shades of blue as in



**fig S1B**, highlighting the identical conformation of the native N-terminal strand (sticks). **(C)** Multiple sequence alignment of NTDs of ARC and dARC CA sequences. Group 1 contains tetrapod ARC (tARC) sequences and the closely related *L. ch* Gypsy2 retrotransposon. Above, rARC secondary structure, numbers according to the rARC (*R. norvegicus*) sequence). Group 2 contains dARC CA sequences and the closely related *L. h* Gypsy11 retrotransposon. Below, dARC1 secondary structure, numbers according to the dARC1 (*D. melanogaster*) sequence). Red box, white text, invariant residues shared between groups. Red text, residues conserved within a group. Blue highlight, alpha-helical core which is conserved between the two groups. Yellow highlight, native N-terminal strand which is conserved within group 1. TQIF motif is boxed in black. Green highlight, native N-terminal strand which is conserved within group 2. Conserved, strand-burying aromatic residues are boxed in black. **(D)** Close-up view of the dARC1 CA-NTD (cyan), aligned with the rARC CA NTD-CaMK2B (light pink-orange) and rARC CA NtD-TARPy2 (dark pink-lime) structures, showing identical mechanisms of strand burial in the core domain. The native strand (dARC1) or rARC-bound peptides (TARPy2 and CaMK2B) as well as the NTD  $\alpha$ 2- $\alpha$ 3 loop are shown in stick representation, the equivalent strand-anchoring aromatic residue (circled) in each case is coloured purple.





**Fig. S4. Comparison of dARC1 and Ty3 CA. (A and B)** Cartoon representations of CA-CTD dimers. **(A)** dARC1, cyan and wheat, **(B)** Ty3 (PDB: 6r23), magenta and pale green. CTD helices are labelled  $\alpha 5$  to  $\alpha 9$  according to the dARC1 crystal structure. The buried surface area ( $\text{\AA}^2$ ) and free energy of interaction ( $\Delta^i G$ ) of each interface, calculated in PDBePISA is displayed below each structure. **(C)** Structural alignment of dARC1 CA-CTD and Ty3 CA-CTD. The structures are aligned with respect to the dARC1 dimer (RMSD = 3.3 Å, 143 C $\alpha$ ), protein backbones are coloured as in **A and B**. **(D)** View of the dARC1 and

Ty3 dimers showing the positioning of NTDs with respect to the CTD dimer. The structures are aligned on the CTD dimer as in **C**. Ty3 is coloured in grey, dARC1 CA-NTD and CA-CTD are coloured blue and red respectively. The dARC1 CA-NTDs have to undergo a slight rotational operation to match the conformation seen in the Ty3 dimer (represented by the shaded arrows). Despite this difference, both dimers have the same “glacial trough” arrangement and the dARC1 crystallographic dimer has an NTD-CTD orientation very close to that observed Ty3 icosahedral particles.