

# A comparative analysis of the foamy and ortho virus capsid structures reveals an ancient domain duplication

William R. Taylor<sup>a</sup>, Jonathan P. Stoye<sup>b</sup>, Ian A. Taylor<sup>c</sup>

*Francis Crick Institute, 1 Midland Rd., London NW1 1AT, UK*

<sup>a</sup>*Computational Cell and Molecular Biology,*

<sup>b</sup>*Retrovirus-Host Interactions,*

<sup>c</sup>*Molecular Structure Laboratories,*

---

## Abstract

to be written

*Keywords:* Virus capsid structure, foamy virus evolution, protein structure comparison

---

## 1. Introduction

1 Taxonomically, the *Orthoretrovirinae* (orthoretroviruses) and *Spumaretro-*  
2 *virinae*<sup>1</sup> (spumaviruses) make up the two subfamilies of *Retroviridae*. They  
3 share many similarities, including overall genome structures with gag, pol  
4 and env genes encoding proteins for replication and life cycles involving re-  
5 verse transcription and integration into the chromosomes of infected cells.  
6 However, there are also a number of differences distinguishing these viral  
7 subfamilies, including finer details of genome organisation, the absence of a  
8 Gag-Pol fusion protein in spumaviruses and the timing of reverse transcrip-  
9 tion.  
10

11 Gag is the major structural protein of both Ortho and Foamy viruses and  
12 also displays both important differences and similarities. Ortho and Foamy  
13 viral Gag are required for particle assembly, budding from the cell, reverse  
14 transcription and delivery of the viral nucleic acid into the newly infected

---

<sup>1</sup>This class is also commonly referred to as the Foamy viruses (after the morphological effect they have on infected cells) and will be referred by this name frequently below, with the term orthoretroviruses also contracted to "Ortho viruses".

cell. However, there are a number of striking differences including how the Gag precursor is targeted to the cell membrane, the absence of a Major Homology Region and Cys-His box in Foamy viruses and very different patterns of processing during viral maturation. In all Ortho viruses, Gag is proteolytically cleaved to form distinct, well-studied proteins, matrix (MA), capsid (CA) and nucleocapsid (NC), found in mature virions but in spumaviruses Gag processing does not occur.

The recent solution of the Foamy Gag protein structure has shed new light on this relationship by revealing that the capsid structures of both viral classes share a common protein fold, with the implication that their gag proteins may be evolutionarily related [1]. An intriguing aspect of this relationship was an ambiguity in the degree of relatedness between the two domains of the gag proteins, with the Spumaretroviral Gag domains appearing almost equally similar to both the amino- and carboxy-terminal domains of the orthoretroviruses. In this paper, we investigate the nature of this relationship in greater detail and discuss its evolutionary implications.

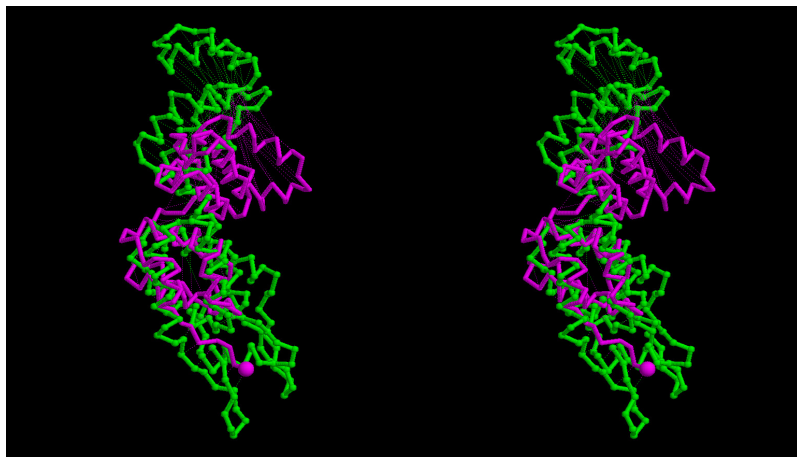
## 2. Results

### 2.1. Full-length comparison

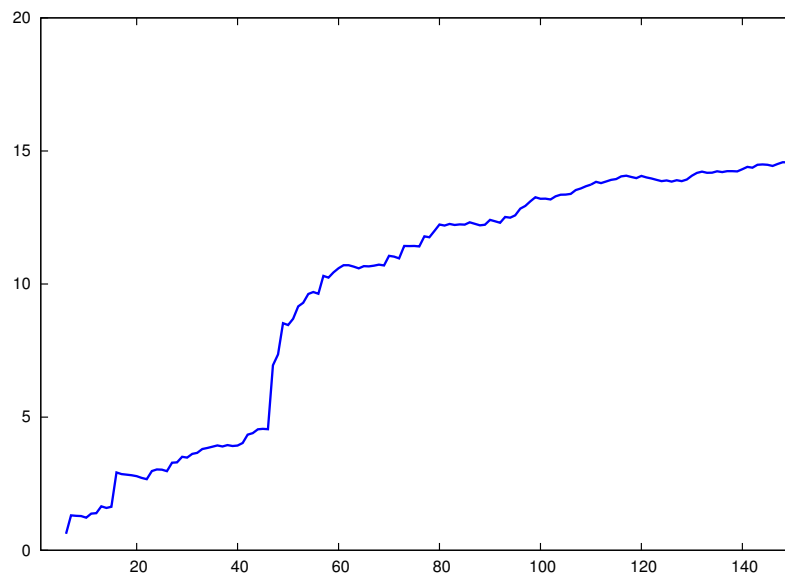
To investigate the structural relationship between the capsid structure of the ortho viruses (HIV, MLV, etc.), and the new structure of the foamy virus capsid [1] (PDB codes: 5m1g, 5m1h), the foamy virus structure was compared to one of the few full double domain ortho virus structures, the HIV capsid with PDB code: 3nte, using the flexible superposition program SAP [2]. Even though this program has a tolerant approach to relative domain shifts, the comparison produced a high RMSD value of 14Å over the 100 best superposed positions. The amino (N) terminal domain positions roughly corresponded but shifts in the relative orientation of the carboxy (C) terminal domain resulted in large deviations between equivalent helices. The superposed structures are shown in Figure 1(a) and the domain divergence can be seen clearly as a jump in the cumulative RMSD plot (Figure 1(b)).

### 2.2. DALI searches

Although this initial superposition (Figure 1) did not appear encouraging, the foamy virus structure was scanned across the Protein DataBank (PDB), using the DALI program [3] to search for any similarities.



(a)



(b)

Figure 1: **Full ortho/foamy virus capsid superposition.** The superposed structures are shown in part (a) as a stereo pair, coloured as green = ortho virus (HIV, PDB code: **3nte-A**) and magenta = foamy virus capsid. (The amino terminus is marked by a small sphere). Part (b) shows the cumulative RMSD plot for this superposition which plots the RMSD value (Y-axis) for increasingly larger sets of residues as ranked by their **SAP** similarity score (X-axis). The sharp rise in this trace marks the transition into subsets that include positions from the displaced domain.

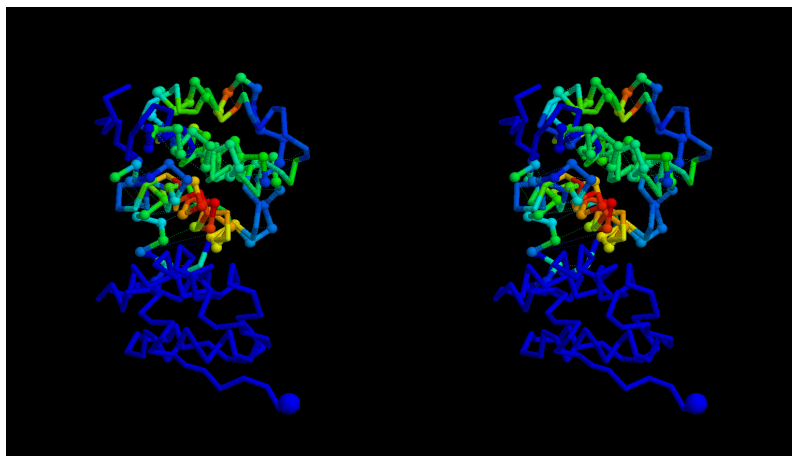
No:	Chain	Z	rmsd	lali	nres	%id	PDB	Description
1:	4x3x-A	5.0	3.1	66	82	11	PDB	MOLECULE: ACTIVITY-REGULATED CYTOSKELETON-ASSOC
2	3g29-A	3.7	2.7	60	77	8	PDB	MOLECULE: GAG POLYPROTEIN;
3	3g0v-A	3.7	2.9	62	76	8	PDB	MOLECULE: GAG POLYPROTEIN;
4:	2v50-D	3.6	2.2	41	998	7	PDB	MOLECULE: MULTIDRUG RESISTANCE PROTEIN MEXB;
5:	3j39-i	3.6	2.5	40	113	3	PDB	MOLECULE: 60S RIBOSOMAL PROTEIN L10A-2;
6	4ph2-A	3.6	3.2	69	127	7	PDB	MOLECULE: BLV CAPSID - N-TERMINAL DOMAIN;
7:	11qp-E	3.6	3.8	69	326	7	PDB	MOLECULE: RFCS;
8:	4gco-A	3.6	3.7	55	120	11	PDB	MOLECULE: PROTEIN STI-1;
9	3g29-B	3.6	2.8	62	77	8	PDB	MOLECULE: GAG POLYPROTEIN;
10	3g1i-B	3.6	2.9	62	75	8	PDB	MOLECULE: GAG POLYPROTEIN;
11	3g21-A	3.6	2.8	60	77	8	PDB	MOLECULE: GAG POLYPROTEIN;
12:	2a0u-A	3.5	3.1	68	374	4	PDB	MOLECULE: INITIATION FACTOR 2B;
13:	1j7q-A	3.5	2.9	60	86	5	PDB	MOLECULE: CALCIUM VECTOR PROTEIN;
14:	2a0u-B	3.5	8.1	80	367	4	PDB	MOLECULE: INITIATION FACTOR 2B;
15:	11qp-A	3.5	3.7	70	326	7	PDB	MOLECULE: RFCS;
16	4ph0-C	3.5	4.6	101	199	8	PDB	MOLECULE: BLV CAPSID;
17	4ph0-D	3.5	4.2	101	198	8	PDB	MOLECULE: BLV CAPSID;
18	4ph2-B	3.5	3.3	69	127	7	PDB	MOLECULE: BLV CAPSID - N-TERMINAL DOMAIN;
19:	1sxj-B	3.4	3.5	65	316	3	PDB	MOLECULE: ACTIVATOR 1 95 KDA SUBUNIT;
20:	2afd-A	3.4	2.7	59	88	14	PDB	MOLECULE: PROTEIN ASL1650;

Figure 2: **Top structural similarities** found by the DALI program in the 90% non-redundant PDB (PDB-90) using the full length foamy virus capsid as a query (145 residues). The columns are: the ranked number of the hit (No.), marked by a '|' for a capsid protein, otherwise ':'; the PDB entry identifier (Chain, with the chain designation after the dash); the DALI Z-score (Z) (significance estimate); the root-mean-square-deviation (rmsd) over aligned  $\alpha$ -carbon positions; the number of aligned positions (lali); the number of residues in the matched structure (nres); the percentage sequence identity of the match (%id) followed by a description of the molecule. It can be seen from the number of matched positions (lali) that most matches are partial, covering typically less than half the query structure.

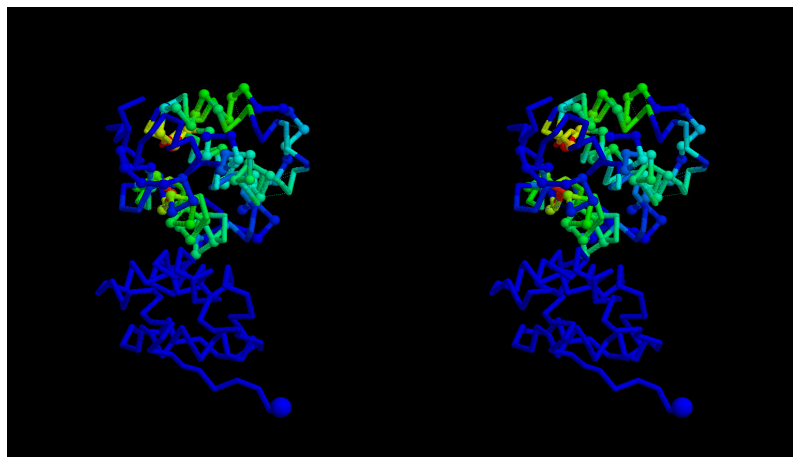
### 2.2.1. Full chain scan

A scan of the full-length foamy structure using the DALI server<sup>2</sup> over the 90% non-redundant protein structure databank identified a wide selection of retroviral capsid structures. In the ranked list of structure hits, capsids were identified from position 2 to position 550. The top hits are shown in Figure 2 (See Supplementary material for a summary of the full 550 with Z-scores over 2). Many capsids are found in the top 20 hits and although the top scoring hit is not obviously a capsid protein, it is thought to have originated from the Ty3/Gypsy retrotransposon family gag gene [4]. However, almost all of these are partial hits, covering little more than half the query structure. The structural alignment of the top two hits is shown in Figure 3 coloured to emphasise the matched regions.

<sup>2</sup>[http://ekhidna.biocenter.helsinki.fi/dali\\_server](http://ekhidna.biocenter.helsinki.fi/dali_server), see Methods section for details.



(a) 4x3x-A



(b) 3g29-A

Figure 3: **Top hits superposed.** The top two DALI hits to the full foamy virus capsid are shown as a  $\alpha$ -carbon backbone (stereo pair) coloured using the residue similarity score calculated by SAP. (red = strong similarity, blue = none). The amino terminus of the foamy structure is marked by a large ball and the other structure is distinguished by small balls on its  $\alpha$ -carbon atoms. (a) a cytoskeleton associated protein (fragment) of the arc/arg3.1 gene (PDB code: 4x3x-A), (which is thought to have originated from a Ty3/Gypsy retrotransposon family capsid) and (b) the structure of the capsid C-terminal domain of the Rous sarcoma virus (PDB code: 3g29-A).

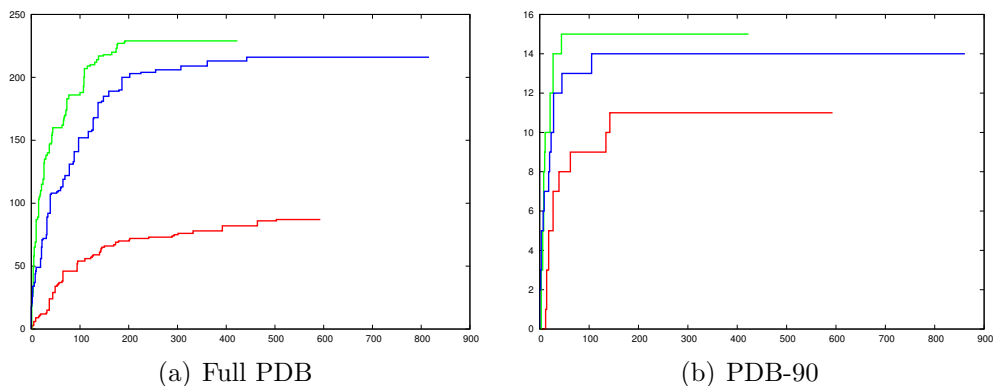


Figure 4: **PDB capsid structure matches.** The number of capsid structures identified by the DALI program in (a) the full PDB and (b) the 90% non-redundant PDB (PDB-90) is shown for queries using the full foamy capsid structure (red), the carboxy terminal domain (green) and the amino terminal domain (blue). The number of capsid hits (Y-axis) is plotted against the order of all hits ranked by Z-score down to a value of 2. A curve approaching the top left corner indicates greater specificity and the extent of a curve to the right indicates the total number of hits.

The result of the DALI search indicated that the Foamy virus structure shares some similarity with the capsid structure of the ortho-viruses. However, the matches consist only of a small number of helices and appears barely more convincing than other matches to proteins that seem very unlikely to have any meaningful connection to a viral capsid. The preponderance of capsid matches throughout the list of hits might seem to add some support to the relationship but may simply be a reflection of the number of capsid structures in the structure databank.

Adding confusion to the ortho/foamy relationship is the additional observation that the distribution of matches to the ortho-virus structures between the amino (N) and carboxy (C) terminal domains are mixed. For example; taking the top 10 matches, the N-terminal domain of the Foamy structure aligns with 6 C-terminal domains and 4 N-terminal domains of the ortho viruses and the best match with the corresponding Foamy C-terminal domain aligns with an ortho N-terminal domain.

### 2.2.2. Domain scans

To clarify the domain match specificity, the two domains of the Foamy virus (1–88 and 89–180, as defined automatically [5]) were scanned separately

79 using the DALI program. The individual domains were much more specific  
80 at matching known capsid structures<sup>3</sup>, both in the full PDB and PDB-90  
81 collections as can be seen from the plots in Figure 4.

82 The results of these scans strengthened the identification of the rela-  
83 tionship to the ortho capsids and supported the swapped specificity for the  
84 N-terminal match of the Foamy structure with the C-terminal match of the  
85 ortho virus and *vica versa*, with all top 12 hits of each domain matching  
86 their opposed counterpart. The structure-based sequence alignments of each  
87 domain based on this equivalence are shown in Figure 5.

88 Although domain transposition is not impossible in viral genomes, it is  
89 sufficiently unexpected to warrant deeper investigation, especially as it is  
90 hard to imagine how an ancestral capsid protein could tolerate such a large  
91 rearrangement and still pack to form a competent shell. We therefore under-  
92 took a more thorough evaluation using alternative methods to assess the  
93 statistical significance of these structural similarities.

### 94 2.3. Structural alignment significance

#### 95 2.3.1. Reversed-structure searches

96 For each comparison, the DALI program calculates an empirical Z-score,  
97 combining an estimation of significance with protein length normalisation.  
98 The program reports all matches over  $Z=2$ , however, when the proteins are  
99 small and especially when the structures being compared are both predom-  
100 inantly alpha-helical in nature, then matches over this cutoff include many  
101 functionally unrelated hits where the similarity has arisen through the for-  
102 tuitous alignment of a few helices.

103 Therefore, to calculate a stricter cutoff on score, we created a decoy probe  
104 by reversing the alpha-carbon backbone then reconstructing the full atomic  
105 structure, using a simple algorithm to regenerate a full backbone<sup>4</sup>). Figure 6  
106 plots the ranked DALI Z-scores for the separate (native) foamy domains. As  
107 would be expected, the larger C-terminal domain has hits with a higher sig-  
108 nificance than the smaller N-terminal domain: the former covers the range  
109  $Z=2.5$  to  $Z=5$  over the true hits (magenta dots) whereas the latter tracks a

---

<sup>3</sup>True/false hits were defined by protein descriptions with the words "CAPSID",  
"GAG" or "P24".

<sup>4</sup>Note that reversing the  $\alpha$ -carbon backbone does not change the chirality of the  $\alpha$ -  
helices but as DALI requires a full atomic backbone, this must be restored on the reversed  
chain.

Nter	PIGTVPIQHIRSVTGEPPRNPREIPIWLGKRNAPIDGVFPVTTPLRCRIINAILGGNIGLSLTPGDCLTWDSAVATLFI	RTHGTFP
	: : :   : :   : : : : : : : : :   : :	
3g1gA	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3tirA	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3g1iA	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3g29A	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3g0vA	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3g29B	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3g1iB	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3g26A	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
3dtjC	-----SILD--IRQGP--EPFRDYVDRFYKTLR--VKNW--MTATLLVQANPD-TILKGPGA--TLEEMTA	-CQGV--
3dtjB	-----SILD--IRQGP--EPFRDYVDRFYKTLR--VKNW--MTATLLVQANPD-TILKGPGA--TLEEMTA	-CQGV--
3dtjA	-----SILD--IRQGP--EPFRDYVDRFYKTLR--VKNW--MTATLLVQANPD-TILKGPGA--TLEEMTA	-CQGV--
3g21A	-----PWAD--IMQGPS--SFVDFANRLIKAVEGSDL-ARAPVIIDCFRQKSQPQQLI--PSTL-TTPGEI	IKYVLD
Cter	MHQLGNVIKGIQVDEGVATAYTLGMMLSGQNYQLVSGIIRGYLPQAVVTALQRLDQEDNQTRAETFIQHLNAVYEILGLNARGQSIRL	
	: :   : : : : :   : : : : :   : : : : :   : :   : :	
116nA	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----PTSI
3j34U	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
4u0bF	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
4u0bG	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
3h4eB	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
2jprA	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
1afvB	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
4u0bE	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
4u0bK	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
4u0bH	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
2gonA	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS
1afvA	SPRTLNAWVKVVEEKA-IPMFSALSC--GATPDLNLTMLNTVGGHQAAMQMLKETINEEA--EIKRWIILGLNKIVRMY	-----SPTS

Figure 5: **Top domain similarity alignments.** The sequence alignments are shown for the top 12 capsid domain matches found by the DALI program using the foamy virus capsid N and C domains separately as a query over the full PDB. The sequence of the N-terminal domain (N-ter) is shown at the top of the first alignment block and the sequences of the C-terminal domain (C-ter) at the top of the second block. The sequences of the ortho-viruses aligned below these all come from the "swapped" relationship of C and N terminal domains, respectively. These alignments, which are determined by structure not sequence, exhibit no specific similarity beyond what would be expected from aligning similar secondary structures from similar sized domains. (Amino acid identities are marked by a bar and similarities by a colon).



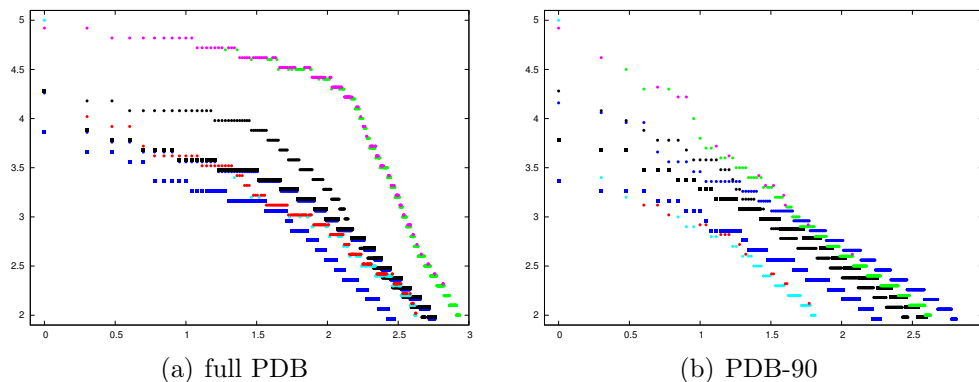


Figure 6: **Ranked DALI scores with decoys.** The DALI Z-scores (Y-axis) are plotted against the  $\log_{10}$  of their ranked position in the list of hits (X-axis) with the amino-terminal domain (N) as T=red, F=cyan dots and the carboxy-terminal domain (C) as T=magenta and F=green dots, where T is a true capsid hit and F is a false hit to a non-capsid protein. Four sets of decoys are compared to these, consisting of the reversed foamy capsid domains in black and the reversed HIV capsid domains in dark-blue (with a circle = N and a square = C domains in both). The DALI score for each set of hits has been slightly displaced to prevent coincident dots from being obscured. (This happens because of the integral number of residues and the DALI score being specified to only one decimal place).

110 similar profile running one Z-value unit lower (2–4 over true red dots). Plot-  
 111 ting the Z-scores against the log of their rank produces almost linear traces  
 112 for the hits from the PDB-90, making it easy to compare N-domain (red/cyan  
 113 dots) with C-domain (magenta/green dots) (for T/F hits) in Figure 6.

114 The equivalent scans with the reversed domain structures, using both  
 115 the foamy and ortho (HIV) structures (neither of which should have any  
 116 particular relationship to the capsid or any other natural protein) also found  
 117 hits with high Z-scores (black and blue points in Figure 6, respectively).  
 118 When compared with the native domains (Figure 6), these decoys had a  
 119 profile that tracked mostly above the N-terminal native domain but below  
 120 the C-terminal domain. However, with the latter domain, this was only  
 121 distinct in the hits to the full PDB whereas with the PDB-90, the native  
 122 domain was only clearly better over the top 10 matches, half of which were  
 123 to non-capsid structures.

124 The results with the simple reversed decoy using DALI suggested that  
 125 the match of the foamy virus domains to the ortho virus capsid N-terminal  
 126 domain may be due to chance and that the match to the C-terminal domain

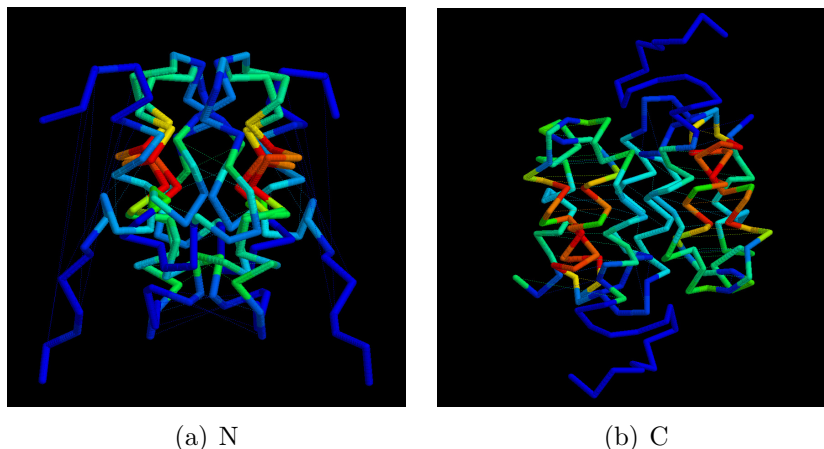


Figure 7: **Native/decoy similarity.** When superposed using the program SAP, both N-terminal (left) and C-terminal (right) domains have some degree of similarity to their reversed decoy 'doppleganger', which is more marked for the N domain. The superposed structures are coloured by the SAP residue-level score as red = high similarity, blue = low. The N domain has roughly 60 equivalent  $\alpha$ -carbon positions compared to only 24 in the larger C domain.

127 looks meaningful if based on the hits to the full PDB but may be only  
 128 marginal based on the PDB-90 hits.

129 However, both the N and C terminal domains pocess a degree of internal  
 130 symmetry which gives rise to a partial match with their reversed 'dopple-  
 131 ganger' decoys. The N-terminal domain superposed on its decoy had an  
 132 RMSD of 5.4/60 ( $\text{\AA}/\alpha$ -carbons) and 5.5/24 for the C-terminal domain. The  
 133 higher symmetry of the smaller domain may be sufficient to explain its poor  
 134 level of specificity seen in Figure 6 and to try and resolve this ambiguity, a  
 135 more diverse set of decoys were generated based on cyclic permutation and  
 136 segment swapping combined with chain reversal [6].

### 137 2.3.2. Customised decoy comparisons

138 To improve the statistical analysis of the foamy/ortho capsid similarity,  
 139 we employed a method based on the generation of a population of customised  
 140 'decoy' models to provide a background distribution of unrelated protein  
 141 scores [6]. This method retains the advantage of the simple reversed struc-  
 142 tures where every comparison that constitutes the random pool is between  
 143 two models of the same size and secondary structure composition as the pair  
 144 of native structures being compared. For this study we collected 12 capsid

145 N-terminal domains and 7 C-terminal domains, each of which were compared  
146 with the foamy N-terminal domain and the foamy C-terminal domain. (The  
147 structures are identified in Table 1 with full details in the Methods section).

148 For each domain pair to be compared, decoys were created using cyclic  
149 permutation and segment swapping with chain reversal to generate a family  
150 of customised decoys for each comparison [6]. All pairs of forward/reversed  
151 decoys were then compared, with each pair being drawn from a pool of mod-  
152 els generated from the two native structures. This ensures that the native  
153 domains (which may have different lengths) are always evaluated against a  
154 decoy pair with the same length combination. (See Methods section for de-  
155 tails). All the decoy comparisons, of which there are typically 150–300 for  
156 each comparison, can then be compared to the native pair on a plot of RMSD  
157 against the number of matched residues ( $\alpha$ -carbon atoms). An example is  
158 shown in Figure 8 for the comparison of the HIV1 structure (PDB codes:  
159 1ak4 (N) and 1a43 (C)) domains against the foamy virus gag domains.

### 160 2.3.3. Statistical analysis of the decoy comparisons

161 The quality of the comparisons in Figure 8 can be quantified as a combi-  
162 nation of their RMSD ( $R$ ) and the number of matched (superposed) positions  
163 ( $N$ ). However, as explained in the Methods section, for statistical analysis,  
164 it is easier to combine this pair of numbers as a single number, called the  
165  $a$ -value ( $\text{Equ}^n. 1$ ), which is the scaling factor that causes a theoretical curve  
166 to pass through the point ( $R, N$ ).

167 When expressed by a single  $a$ -value all the data points in a comparison,  
168 such as Figure 8(c), can be plotted as a frequency histogram and examined  
169 to see if they approximate a Normal distribution. The distributions were  
170 found to be a good fit to unskewed Gaussians and so were treated as normal  
171 distributions (rather than extreme value distributions that have also been  
172 considered previously as a model for random structure comparison scores  
173 [7, 6]). The frequency data from the comparison of the orthoN domain from  
174 HIV1 and the foamyC domain (Figure 8(c)) is shown in Figure 9(a) along  
175 with a Normal distribution that has the same mean ( $\mu$ ) and standard devia-  
176 tion ( $\sigma$ ) as the data. On this plot, the value of  $a$  ( $\text{Equ}^n. 1$ ) for the comparison  
177 of the native pair of domains is also plotted (blue triangle) and from its po-  
178 sition, a Z-score can be calculated.

179 In this way, the significance of all combinations of the native ortho and  
180 foamy domain superpositions were calculated, using the background distri-  
181 bution of 'customised' decoy comparisons based on each individual native

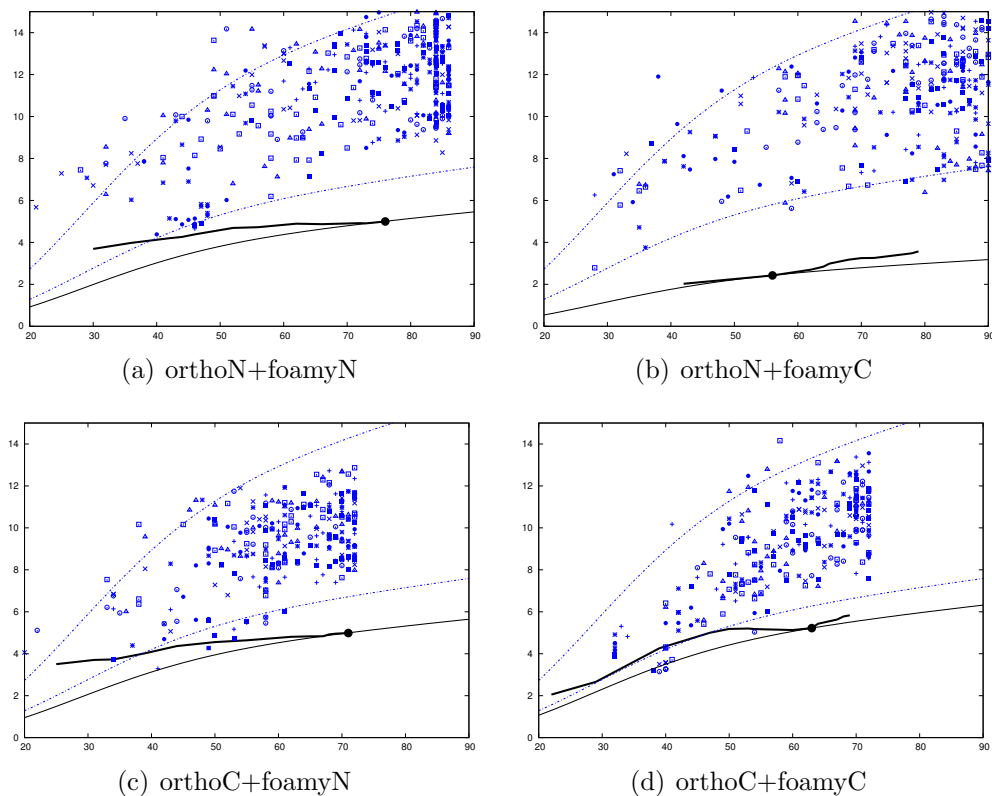


Figure 8: **ortho/foamy domains compared with customised decoys.** Each amino (N), carboxy (C) domain combination between the ortho retrovirus capsid structure (HIV1) and the foamy virus capsid structure is plotted as a line for increasingly large subsets of matched positions against their RMSD (Y-axis), as in Figure 2. The point on this line marks the lowest  $a$ -value (Equ<sup>n</sup>. 1), however, to be consistent with the decoy data, the full alignment length was used. The decoy comparison data (blue) is plotted in a variety of symbols with each representing a different combination of decoy construction. The dashed blue lines (which are the same in all plots) mark the approximate 10<sup>th</sup> percentile boundaries of the decoy generated distributions, with  $a = 1.7$  (upper) and  $a = 0.8$  (lower). (See Methods section for details).

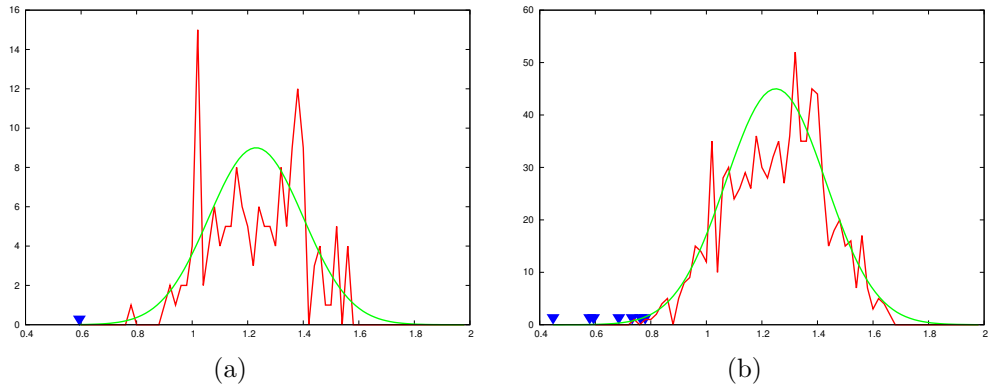


Figure 9: **ortho-C and foamy-N domain comparison statistics.** The  $a$ -value (normalised RMSD) for the comparison of the ortho-C and foamy-N decoy domains (Figure 8(b)) are plotted as a frequency distribution (red) along with a bell-shaped Normal distribution curve (green) with matching mean ( $\mu$ ) and spread ( $\sigma$ ). Part (a) shows the distribution for the HIV1 C-terminal domain ( $\mu = 1.23$ ) and spread ( $\sigma = 0.17$ ) with the position of the native structure comparison plotted as a blue (inverted) triangle. Its position lies 0.64 units below the mean giving a Z-score of  $0.64/0.17 = 3.76$ . Part (b) shows the combined data from seven representative viruses (in Table 1). These data comprise two distributions, that of the combined decoys and also the much smaller distribution of native scores (blue triangles). This allows a T-test to be made on the significance of their separation.

182 pair. The resulting Z-scores ( $\sigma$  units) are collected in Table 1. The degree  
 183 of similarity between the domains ranged from less than  $1\sigma$  to over  $5\sigma$ , with  
 184 the latter (highly significant) result being obtained for both a swapped (NC)  
 185 and forward (CC) combination. However, of the top 12 scores, only three  
 186 now came from swapped pairings.

187 *Asymmetry statistics:*

188 To quantify the degree of bias for domains of like-type (NN, CC) to be  
 189 more similar than those of mixed-type (NC, CN), the observed ranking of  
 190 like and mixed pairs, based on their Z-value (Table 1), was compared to  
 191 that expected by chance. The positions of all pairs in the list were shuffled  
 192 a million times and the asymmetry of each arrangement was quantified as  
 193 the number of like-pairs in the top half and also by their second moment:  
 194  $\sqrt{((\sum r_i^2)/N)}$ , where  $r$  is the rank of the like-pair  $i$  in a list of  $N$  pairs. The  
 195 chance of obtaining a distribution with more like-pairs being ranked higher  
 196 can be calculated by summing the area of the tail of each empirical dis-  
 197 tribution that lies beyond the observed value. However, these values were  
 198 calculated over all pairs and neglects the principle that emphasis should be  
 199 given to the more significant similarities. Rather than rely on a single signifi-  
 200 cance cutoff (like  $3\sigma$ ) or an arbitrary cutoff (like the "3-out-of-12" mentioned  
 201 above), we calculated statistics for all such cutoffs (Figure 10(a)).

202 The majority of values in Figure 10(a) lie below the 0.05 probability level  
 203 for the larger sample sizes, with those for the top-half bias statistic (blue line)  
 204 being more significant than the moment-based statistic (red line). While  
 205 confirming the visual trend towards a bias of higher scoring like-type domain  
 206 similarities, the analysis summarised in Figure 10(a) is complicated by having  
 207 unequal numbers of amino and carboxy domain comparisons and also by  
 208 including some closely related structures. To produce a more balanced data-  
 209 set, one of each pair of the two most similar carboxy domain structures was  
 210 discarded leaving five structures and for each of these, their matching amino  
 211 terminal domain was also retained, leaving: BLV-1, HIV-1, HML2, HTLV-1  
 212 and RSV. Despite having a smaller set of comparisons (5N + 5C domains  
 213 giving 20 rather than 38 Z-scores), the results for this reduced set indicated  
 214 an equally clear bias towards a preferred like-domain equivalence,  
 215 especially as measured by their occurrence in the upper half of the ranked  
 216 list, with several having a probability below the 0.05 level and a few below  
 217 the 0.005 level (Figure 10(b)).

218 *T-test statistic:*

<i>a</i>	ortho-N					
	foamy-N			foamy-C		
virus	pool	<i>a</i> -value	Z-score	pool	<i>a</i> -value	Z-score
BLV6	300	0.552	<b>4.073</b>	244	0.542	3.692
BLV	251	0.550	<b>4.494</b>	184	0.400	3.669
HIV6	312	0.551	3.781	220	0.405	3.579
HIV1	312	0.573	3.703	213	0.402	3.692
HML2	264	0.777	2.166	196	0.438	<b>4.594</b>
HTLV	400	0.592	<b>4.030</b>	328	0.457	<b>4.013</b>
JSRV	225	1.063	0.896	190	0.601	3.237
MLV	326	0.751	3.044	188	0.508	3.151
MPMV	269	0.565	<b>3.902</b>	185	0.523	2.918
PSIV	285	0.621	3.731	235	0.369	<b>5.019</b>
RELIK	234	0.639	3.688	237	0.700	3.297
RSV	204	0.543	3.123	239	0.526	3.542

<i>b</i>	ortho-C					
	foamy-N			foamy-C		
virus	pool	<i>a</i> -value	Z-score	pool	<i>a</i> -value	Z-score
BLV6	144	0.763	3.019	212	0.709	<b>4.046</b>
BLV	154	0.578	3.400	204	0.556	<b>4.047</b>
HIV1	157	0.593	3.760	174	0.705	3.362
HIV6	179	0.780	3.175	177	0.640	<b>4.380</b>
HML2	185	0.732	3.027	184	0.676	<b>3.900</b>
HTLV	156	0.685	3.847	163	0.694	2.807
RSV	155	0.448	3.754	235	0.403	<b>5.009</b>

Table 1: **Ortho and foamy domain comparison Z-score statistics.** For each amino (N) and carboxy (C) domain pair between an ortho virus structure and the foamy virus capsid structure, a **Z-score** is calculated based on the **a-value** (Equ<sup>n</sup>. 1) derived from the comparison RMSD and length, relative to the **pool** of background decoy comparisons. The ortho **virus** identity is indicated by the code to the left, full details of which can be found in the Methods section. The top 12 Z-scores are high-lighted in bold, only three of which support a swapped domain match.

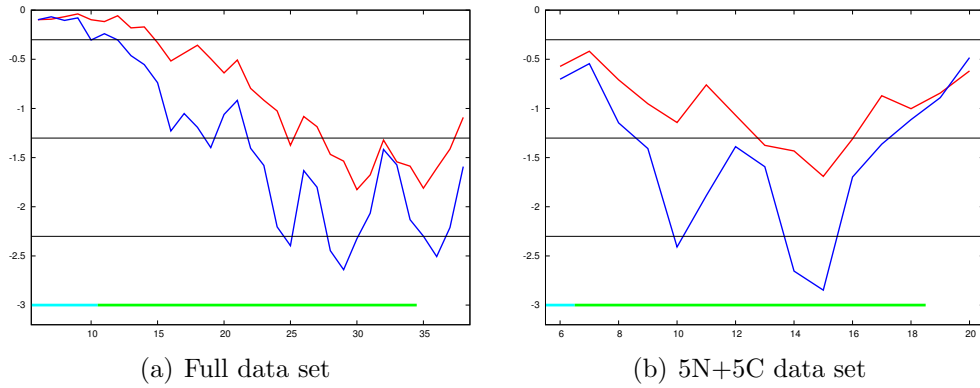


Figure 10: **Asymmetry statistics for like/mixed domain pairs.** Given the ranked list of domain pairings, the chance for more domain pairs of like-type to be found higher than the observed order was evaluated from empirical distributions measured by two statistics: the second moment of the rank value (red) and the number of like-type pairs in the top half (blue). These statistics were calculated for all subsets from the 6 top pairs up to the full set of comparisons (X-axis) and for each, the chance of a better score is plotted as the  $\log_{10}$  of the probability (Y-axis). The horizontal lines mark the 0.5, 0.05 and 0.005 levels. The line at the 0.001 level is coloured by the Z-score for each pair as: green = over 3 and cyan = over 4 sigma. Part (a) shows the probabilities calculated from the full set of 7 carboxy and 12 amino domains and part (b) shows the same values calculated on a more balanced set of 5 non-redundant carboxy domains and their matching amino domains.



	orthoN	orthoC
foamyN	Avg: 6.67e-01 < 1.32e+00 Tprob = 4.62e-21 **  StD: 1.61e-01 = 2.12e-01 Fprob = 1.84e-01	Avg: 6.51e-01 < 1.25e+00 Tprob = 2.35e-16 **  StD: 1.17e-01 = 1.89e-01 Fprob = 1.12e-01
foamyC	Avg: 4.92e-01 < 1.29e+00 Tprob = 4.09e-10 **  StD: 1.02e-01 < 2.21e-01 Fprob = 7.37e-03 **	Avg: 6.22e-01 < 1.30e+00 Tprob = 3.81e-23 **  StD: 1.12e-01 = 1.77e-01 Fprob = 1.20e-01

Table 2: **ortho and foamy capsid domain comparison T-test significance.** For each combination of domains between the ortho and foamy viruses, the probability is given that the two means from each distribution (Avg values) were sampled from the same distribution. (i.e., that the native and decoy comparisons are not distinct). All domain pairings are extremely significant. An F-test was used to test if the standard deviations (Std) of each sample were distinct and if not, the a T-test was made on the assumption of equal standard deviations.

219 An alternative to the above analysis, which still remains marginally sig-  
220 nificant, is to pool the raw comparison data for all the domain comparisons  
221 and their background distributions giving now not just a single value com-  
222 pared to a distribution but two distributions (Figure 9(b)). For these data,  
223 a significance was calculated using Student’s T-test, the values of which are  
224 given in Table 2.

225 From these results, it can be seen that all the four possible pairings are  
226 highly significant with probabilities ranging from  $10^{-10}$  to over  $10^{-20}$ . It is  
227 also clear that the two swapped pairings (NC and CN) have higher proba-  
228 bilities than the forward pairings (NN and CC). Combining the probabilities  
229 ( $P$ ) as:  $\Delta P = \log_{10}(P_{NN}P_{CC}) - \log_{10}(P_{NC}P_{CN})$ , gives a value of 17.7 (42.7 -  
230 25.0) which means that the swapped pairing is almost 18 orders of magnitude  
231 less likely than the forward pairing. Calculating the same statistic on the  
232 reduced 5N+5C domain data set gave a similar result but with a difference  
233 reduced 1000-fold to 15 orders of magnitude.

234 The unexpected swapped pairing, which was indicated originally by the  
235 DALI results, now seems less likely. The preferred, and biologically more  
236 reasonable, result is that the ortho virus domain are related to the foamy virus

domains as a result of genetic divergence from a common, double domain ancestor.

#### 2.4. Internal duplication

The transposed pairings of N/C and C/N (ortho/foamy) domains still retain a high structural significance and this suggests that the two domains are derived from a common ancestral structure, probably as the result of a prior gene-duplication event that has been retained more clearly in the less embellished foamy virus structures. Comparing the two foamy domains gives a Z-score of 2.077 sigma which, although of marginal significance, supports this model. (Figure 11(a, b)).

Such a relationship between the foamy domains implies an equivalent relationship in the ortho viruses and a similar comparison in structures of their N and C domains finds matches with Z-scores ranging from 2 to 4. As with the comparison of the ortho and foamy structures, these can be pooled to allow a joint T-test to be applied. This gave a probability of  $10^{-8}$  that the true N/C domain comparisons were drawn from the decoy distribution, adding strong support to the hypothesis of an ancient gene duplication occurring before the split of the ortho and foamy virus families. (Figure 11(c, d), blue triangles). Supporting this relationship, earlier studies also suggested an internal duplication in the ortho viruses but were based largely on very distant sequence similarity [8].

This test was applied only to the comparison of domains between viruses with known structures for both domains, however, it is not unreasonable to compare amino and carboxy domains across all viruses. The longer loops in the ortho virus domains gives greater scope of structural variation and a wide range of variation was seen ranging from RMSD values under 4 to over 12. When normalised for length ( $a$ -value from Equ<sup>n</sup>. 1) and partial matches under 60 positions excluded, a distinct cluster remains between  $a = 0.5 \dots 0.8$  (4...6Å RMSD) but still with a long tail to higher values. Despite this tail, the T-test on the distributions is highly significant at  $2.7 \times 10^{-17}$ .

One of the better N/C ortho similarities is shown in Figure 12(a), along with the N/C ortho domain superposition in Figure 12(b).

#### 2.5. Fold-space representation

To summarise the structural relationships among the ortho and foamy domains, the matrix of pairwise comparisons was projected into a three-

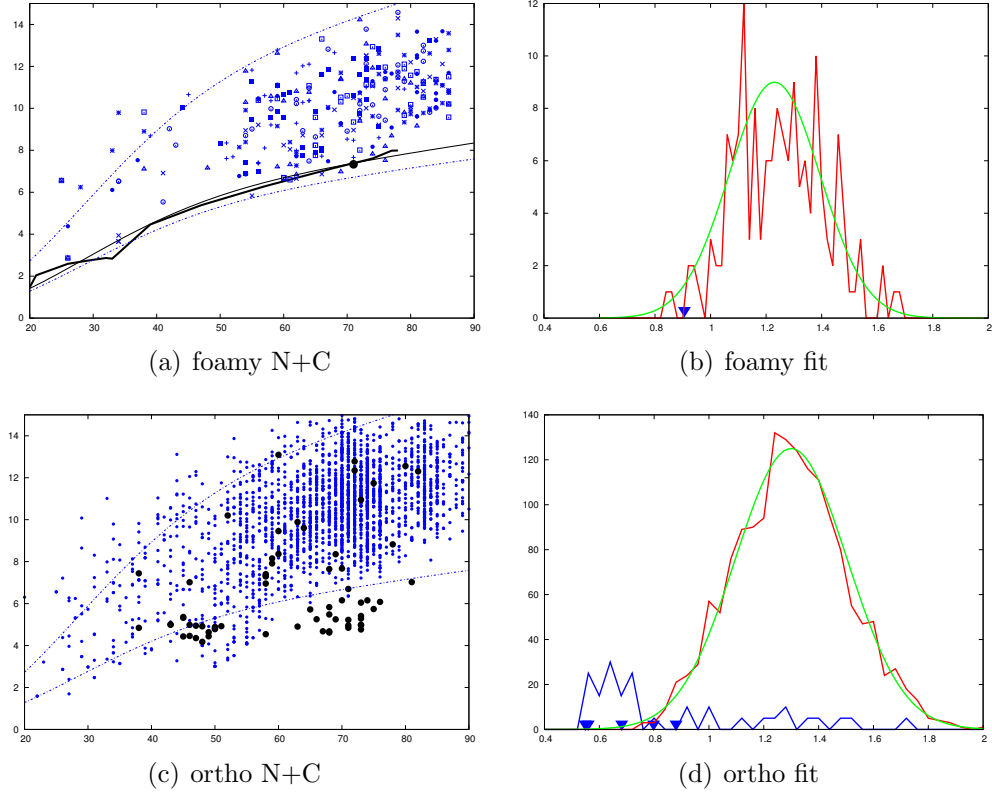
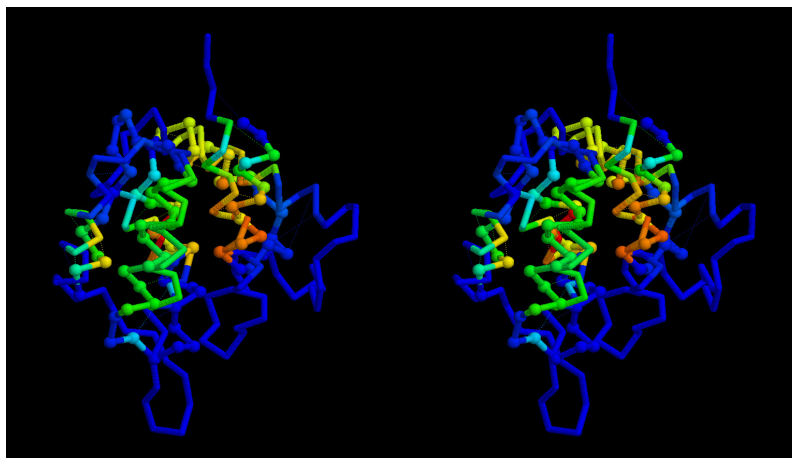
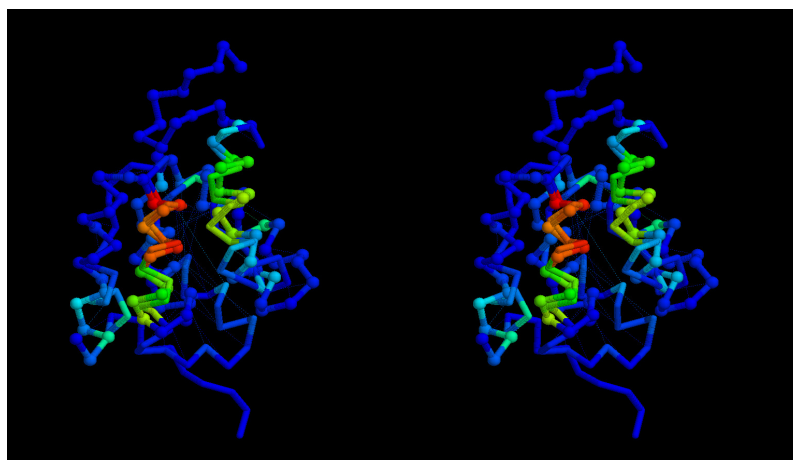


Figure 11: **N and C domains compared with customised decoys.** *a)* The N and C domains of the foamy virus (black) compared to decoys (blue) with *(b)* the derived frequency plot with the native comparison marked by a blue triangle. (See legend to Figure 8 for details). *c)* The N and C domains of the ortho virus combinations (black) with *(d)* the derived frequency plot showing the native comparison for pairs from the same virus (blue triangles) with the distribution of all native pairs shown as a scattered frequency plot (blue line). (See Methods section for details).



(a) ortho



(b) foamy

Figure 12: **Amino and carboxy domains superposed.** *a* ortho virus domains and *b* foamy virus domains are shown as a stereo pair with their  $\alpha$ -carbon backbones coloured by the residue similarity score calculated by SAP. (red = strong similarity, blue = none). The amino terminal domain is distinguished by small balls on its  $\alpha$ -carbon positions and the amino terminus lies to the top in both panels.

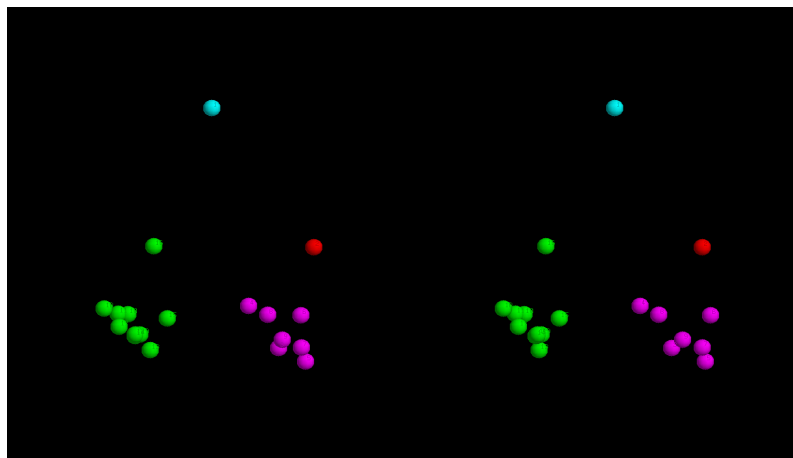


Figure 13: **Fold-space representation of all domains.** All the viral domains considered in the paper were projected into a 3D fold-space representing the relationship of their SAP weighted RMSD values. The domains are coloured as: foamyN = cyan, foamyC = red, orthoN = green and ortho C = magenta.

dimensional fold-space. (See methods for details). This produces a best visual representation of the RMSD values between domains.

As can be seen from Figure 13, the N and C domains of the ortho viruses form distinct clusters with the foamy C domain lying closer to the ortho C-domain cluster. The foamy N-domain, however, maintains a fairly equal distance from both ortho domain clusters but lies closer to its C-terminal partner.

### 3. Conclusions

#### 3.1. Structure comparison

##### 3.1.1. Pairwise significance

The comparison of small domains that are largely composed of  $\alpha$ -helices presents a challenging problem in how to interpret the significance of the RMSD values. As the individual helical secondary structure elements (SSEs) constitute a sizeable fraction of the domain, it takes only the chance alignment of a few helices to result in a low RMSD over a large proportion of the structure, giving an apparently meaningful result.

The use of the customised decoy sets attempts to avoid this problem by recreating a large number of possible folds that were generated using

290 the same (reconnected) SSEs. Moreover, to avoid any chance recreation  
291 of native fragments, each comparison always involved the comparison of a  
292 native (forward) chain direction with a reversed chain. Using these models, a  
293 background distribution of decoy/decoy comparisons allowed us to calculate  
294 Z-scores for each native/native comparison with the advantage that every  
295 comparison in the background distribution involved two models with the  
296 same length, density and secondary composition as the native pair. These  
297 values indicated a clearly significant relationship between the foamy and  
298 ortho structures.

### 299 *3.1.2. Direct or transposed domain order?*

300 However, the Z-scores did not point to a clear resolution of whether the  
301 domains should have a direct correspondance (NN and CC match) or a trans-  
302 posed relationship (NC and CN) with significant individual matches found  
303 across all pairings. Testing for a bias towards more significant like-domain  
304 pairings (NN,CC) in the list of similarities ranked by Z-score confirmed the  
305 visual bias towards a direct correspondance but only at a marginal level of  
306 significance (around 0.05). By contrast, the application of a T-test on the  
307 combined raw comparison data returned a very clear distinction between the  
308 direct and the transposed relationships, clearly favouring the more natural  
309 forward order.

310 Although the ‘astronomic’ probabilities calculated by the T-test seem  
311 very convincing, they must be viewed in the light of the much lower proba-  
312 bilities calculated from the asymmetry statistics. Both calculations involve  
313 assumptions and are limited by the small number of known structures so  
314 neither can be taken as definitive. It would seem likely that the “true” level  
315 of significance may lie somewhere between the two results but as both point  
316 in the direction of the NN and CC domain order, there is no reason to adopt  
317 the more unexpected transposed domain order.

### 318 *3.2. Evolutionary implications*

319 On the basis of these structural comparisons, and a variety of functional  
320 assays described elsewhere, we can conclude that the central domain of the  
321 spumavirus Gag gene encodes a polypeptide sequence related to that of  
322 the corresponding region of orthoretroviruses, CA. It therefore seems rea-  
323 sonable to suppose that the last common ancestor of orthoretroviruses and  
324 spumaviruses possessed such a sequence. Moreover this region appears to be

325 made up from two related subdomains suggesting a gene duplication event  
326 in a common precursor.

327 In our initial search of the foamy virus capsid using the DALI program,  
328 we made the curious observation that the strongest similarity of the foamy  
329 virus capsid was with the ARC protein (Activity-Regulated Cytoskeleton-  
330 associated protein) that is active in neural synaptic growth and activity (and  
331 several other developmental associated functions). The ARC protein has  
332 widespread and clear (non-gag) sequence homologues as far back as insects  
333 and probably deeper, giving it a very ancient origin somewhere close to the  
334 metazoan root [8]. If ARC is considered to be a relic of an ancient Ty3/Gypsy  
335 retrotransposon [4], preserved as a 'living fossil' in the genomes of metazoa,  
336 this relationship would suggest an equally ancient origin for the foamy virus.

337 Alternatively, the foamy viruses may have co-opted an ARC protein to  
338 facilitate budding and their escape from the cell. As it is believed that the  
339 Ty3/Gypsy family of intracellular retrotransposons gave rise to retroviruses  
340 [9], it will therefore be of considerable interest to determine whether such  
341 elements possess CA proteins with a two-domain structure. Finally, it is  
342 worth noting that the Gag protein of Ty3 is significantly shorter than that  
343 of the retroviruses and it is possible that the N-terminal domains of the  
344 orthoretroviruses and spumaviruses were co-opted at different times to facil-  
345 itate budding from the cell surface. If so, the very different structures of this  
346 region in the orthoretroviruses and spumaviruses might suggest independent  
347 acquisition events.

## 348 4. Methods

### 349 4.1. Structural data

350 The foamy virus structures were obtained from the Protein Structure  
351 Databank (PDB code:5M1G) [1].

352 The ortho virus structures used, with their shorthand code in bold and  
353 PDB code in teletype, were:

- 354 • **BLV**: bovine leukemia virus (deltaretrovirus) **4PH1** (N-ter.dom) and  
355 **4PH2** (C-ter.dom) [10],
- 356 • **BLV6**: bovine leukemia virus (hexameric) **4PH0** (both dom.s) [10],
- 357 • **HIV1**: human immunodeficiency virus 1 (lentivirus) **1AK4** (N-ter.dom)  
358 **[11]** and **1A43** (C-ter.dom) [12],

- 359 • **HIV6**: human immunodeficiency virus 1 3H47 (both dom.s) [13],
- 360 • **HML2**: human endogenous retrovirus type-K (betaretrovirus) [14],
- 361 • **HTLV**: human T-cell leukemia virus (deltaretrovirus) 1QRJ (both dom.s)
- 362 [15],
- 363 • **JSRV**: jaagsiekte sheep Retrovirus (betaretrovirus) 2V4X (N-ter.dom)
- 364 [16],
- 365 • **MLV**: murine leukemia virus (gammaretrovirus) 1U7K (N-ter.dom) [17],
- 366 • **MPMV**: Mason-Pfizer monkey virus (betaretrovirus) 2KGF (N-ter.dom)
- 367 [18],
- 368 • **PSIV**: prosimian immunodefficiency virus (ancient lentivirus) 2XGV (N-
- 369 ter.dom) [19],
- 370 • **RELIK**: rabbit endogenous lentivirus type-K (ancient lentivirus) 2XGU
- 371 (N-ter.dom) [19],
- 372 • **RSV**: Rous sarcoma virus (alpharetrovirus) 3G1I (both dom.s) [20].

## 373 4.2. Structure comparison

### 374 4.2.1. DALI

375 The DALI method for searching the PDB with a structural query [3] was  
 376 accessed via the server at: [http://ekhidna.biocenter.helsinki.fi/dali\\_server](http://ekhidna.biocenter.helsinki.fi/dali_server).  
 377 The DALI method reports the significance of each match with an estimated Z-  
 378 score which is the raw comparison score, normalised by the combined length  
 379 of the proteins. Z-scores down to a value of 2 are reported by the program.

380 The list of DALI hits (ranked by Z-score) were assessed by how many high-  
 381 scoring capsid structures had been identified. These true/false (T/F) hits  
 382 were defined simply by protein descriptions that contained the words "CAP-  
 383 SID", "GAG" or "P24". This may have misclassified a few (low scoring) hits  
 384 to the matrix protein and missed some hits where the primary description  
 385 refers to a cyclophilin structure solved in complex with the capsid.

386 DALI reports structural hits in both the full PDB and a reduced collection  
 387 of structures that have no pair of proteins with over 90% sequence identity,  
 388 referred to as the 90% non-redundant or PDB-90 collection. It was found,  
 389 however, that some hits, seen in the full PDB were not found in the PDB-90,



390 for example in Figure 6, all of the top 31 hits of the N-domain against the  
391 full PDB are missing in the PDB-90 hits. The most likely explanation is  
392 that the PDB-90 secection has not been updated at the same time as the full  
393 collection. For this reason, hits to both databases were monitored.

#### 394 4.2.2. SAP

395 The **SAP** method for structure comparison [2] was run as local copy which  
396 can be accessed at: **mathbio**. As part of determining the alignment between  
397 two structures, the **SAP** program calculates a similarity score for each pair of  
398 matched positions which is how similar the rest of the structure looks from  
399 the viewing-frame of the superposed residues. This value can be used both to  
400 weight the importance of positions when calculating the (rigid-body) RMSD  
401 superposition and to colour positions in the superposed structures. [21]. (As  
402 in Figure 3).

403 If the matched positions are ranked by this value, then RMSD values can  
404 be calculated over increasingly larger subsets to high-light the extent of a  
405 well matched core before the contribution of variable loops, or domain shifts,  
406 leads to higher RMSD values. (As in Figure 1(b)).

#### 407 4.3. Decoy structure construction

##### 408 4.3.1. Reversed structure decoys

409 Simple structural decoys were generated from native PDB structures by  
410 reversing the order of the  $\alpha$ -carbon atoms in the PDB file using the Unix  
411 command line:

```
412 cat native.pdb | grep ' CA ' | sort -nr -k2 > reverse.pdb
```

413 The reversal of a protein chain does not alter the chirality of the alpha helix  
414 and these decoys can be used directly in **SAP**. However, **DALI** requires all  
415 main-chain atoms and these must be regenerated for the reversed decoys.  
416 This was done using the simple **ca2main** program which can also be found  
417 at: **mathbio**.

##### 418 4.3.2. Customised decoys

419 Customised structural decoys were generated for each comparison using  
420 each of the pair of structures being compared to create two pools of decoys  
421 then comparing all decoys in the first pool against all decoys from the second  
422 but with their chain reversed as described in the previous section.

423 The decoys were created as described by Taylor [6], starting by cyclising  
 424 the chain then introducing new termini in each surface loop to create cyclic  
 425 permutations. In addition, when three loop regions lie in close proximity,  
 426 their ends are also swapped. That is: if a chain runs from amino (N) to  
 427 carboxy (C) termini through three adjacent loop regions **a-b**, **c-d** and **e-f**  
 428 as: N,**a-b,c-d,e-f**,C then the swapped chain runs: N,**a-d,e-b,c-f**,C with  
 429 each switch being made at the least disruptive point. This swap does not  
 430 create any reversed segments which would otherwise form regions of local  
 431 matching when the whole chain is reversed.

432 In a pair of structures, if each have four surface loops where breaks can be  
 433 made, then including the native termini, this gives five cyclic permutations  
 434 and if two groups of loops can be reconnected then a total of 15 distinct decoys  
 435 can be made from each native starting structure. As these can be compared  
 436 pairwise, a pool of 225 decoy derived data points is generated that constitutes  
 437 the random background against which the native/native comparison can be  
 438 assessed.

439 For example, in Figure 8, the 36 data points marked by a solid circle come  
 440 from the comparison of six cyclic permutations of a native ortho domain  
 441 compared with six permutations of a reversed foamy domain that includes a  
 442 single loop reconnection.

443 Every pair drawn from this pool will have the same lengths as the two  
 444 native structures as well as the same secondary structure composition, surface  
 445 exposure and inertial properties but each decoy will have a different chain  
 446 fold.

#### 447 4.4. Statistical tests

##### 448 4.4.1. RMSD length normalisation

449 The quality of structure comparisons can be characterised by a combina-  
 450 tion of their RMSD value and the number of matched (superposed) positions.  
 451 How to combine these values has been the subject of much discussion over  
 452 the years and central to this is the expected random RMSD value for two  
 453 proteins of a given length [22, 23, 24]. However, when reviewed [6], all these  
 454 measures were approximations of a simple square-root function of the protein  
 455 length (as originally proposed by McLachlan on theoretical grounds [22]) but  
 456 with an added term to depress the RMSD values obtained with small units  
 457 or structure that are dominated by secondary structure elements (and super-  
 458 secondary structure motifs) giving a lower than expected RMSD value. The  
 459 formula that best captures this is:  $R = \sqrt{N(1 - \exp(-N^2/s^2))}$ , where,  $R$  is

460 the expected random RMSD for  $N$  matched positions and  $s$  is the damping  
 461 factor in the inverted Gaussian term (equivalent to the standard deviation  
 462 in the Normal distribution).

463 Any point that lies on this line can be considered "exactly" random with  
 464 those above it being "more" random and those below it "less" random. This  
 465 can be quantified as a single number which is the value of a scaling factor  
 466 ( $a$ ), which when applied to the curve, makes it pass through any given point.  
 467 If a comparison has an RMSD of  $R$  over  $N$  positions, then  $R = a\sqrt{N}(1 -$   
 468  $\exp(-N^2/s^2))$  and when

$$a = R/(\sqrt{N}(1 - \exp(-N^2/s^2))), \quad (1)$$

469 the line will pass through the data point. This reduces the pair of values  
 470 ( $R, N$ ) to a single value  $a$  that is a simpler quantity for statistical analysis.

471 The best value for  $s$  is slightly dependent on the nature of the proteins  
 472 being compared. For artificial (random-walk) models with no secondary struc-  
 473 ture, no modification will be needed but the proteins considered here have  
 474 segments of packed alpha helices that can be locally similar over two to three  
 475 helices. To correct for this, a value of  $s = 30$  was used (or  $1/s^2 = 0.11$ )  
 476 which is higher than the value of  $1/s^2 = 0.03$  used previously. That this is  
 477 a reasonable fit to the data can be seen in the way the dashed blue lines  
 478 in Figure 8 track the upper and lower boundary of the decoy comparison  
 479 results.

480 When  $a = 1$ , the point lies on the random line and when  $a = 0$ , the RMSD  
 481 is zero, so values of  $a$  that approach this lower bound will be of interest when  
 482 evaluating similarity.

#### 483 4.4.2. Frequency plots

484 The  $a$ -values obtained using Equ<sup>n</sup>. 1 were plotted as frequency histograms  
 485 using using only data points that had a length of  $N \pm 10$ , where  $N$  is the max-  
 486 imum number of matched positions. Previously, a cumulative plot of RMSD  
 487 was used to select an optimal value for  $N$  (giving the minimum  $a$ -value). This  
 488 can be important if the full set of matched positions is dominated by a high  
 489 deviations from variable loop regions. However, in the current application,  
 490 the small length of the foamy virus loops meant that this was not an impor-  
 491 tant aspect and the full number of matched positions was taken. Otherwise,  
 492 the same correction would have to be applied to all decoy comparisons to  
 493 maintain a fair comparison. (See Figure 8, where the black dot marks the  
 494 minimum  $a$ -value length).

495 The mean and standard deviation of the  $\alpha$ -values in the  $N \pm 10$  region  
496 were calculated and the corresponding Normal distribution used to calcu-  
497 late Z-scores for the associated native comparison. (See Figure 9(a), for an  
498 example).

#### 499 4.4.3. *T-tests*

500 Data from separate native/native comparisons, with their customised de-  
501 coy data, were combined giving not only a much larger background decoy  
502 derived population of scores but also a smaller distribution of native com-  
503 parison scores that can be tested to calculate the probability that they were  
504 drawn from the same population as the decoy data. To do this, a T-test was  
505 used which takes the size, mean, and standard deviation of each distribution  
506 and calculates a probability. The implementaion of this test was taken from  
507 the Numerical Reicpies collection [25] which implements one of two variants  
508 of the test depending on whether the distributions have statistically distinct  
509 standard deviations. (Routines `ttest()` and `tutest()`). The choice of rou-  
510 tine is based on a preapplication of an F-test on the standard-deviations.  
511 (Using the routine `ftest()`).

512 The values quoted in the Results section are for a two-tailed T-test, how-  
513 ever, as it is expected that the native comparisons should always be more  
514 similar than comparisons between random models, then a one-tailed T-test  
515 would be valid, which gives half the probability. As the values in the Tables  
516 are so significant and only the relative relationships are of interest, then the  
517 choice is unimportant.

#### 518 4.5. *Fold-space clustering*

519 The results of the pairwise similarity within a set of structures can be  
520 visualised by treating the RMSD values as Euclidean distances<sup>5</sup> and reducing  
521 their dimensionality to sufficiently few dimensions to be visualised: usually  
522 2D or, better 3D, to visualise the space with less distortion. Rather than use a  
523 simple multi-dimensional scaling (MDS) method ([26]), the more complicated  
524 method of multi-dimensional projection was used ([27], see [28] for a simpler  
525 exposition).

526 This method reduces the dimensionality of the projection in gradual  
527 stages with each step employing triangle-inequality balancing and hyper-

---

<sup>5</sup>In theory, pairwise RMSD values are guaranteed to constitute a consistent Euclidean metric, but only in  $N-1$  dimensions (where  $N$  is the number of structures compared).

528 dimensional real-space refinement. In the real-space refinement stages, a  
 529 weight can be applied to pairwise distances. (This cannot be done in direct  
 530 MDS projection, which can only assign a mass to each point). Weights were  
 531 assigned to distances as a function of their inverse RMSD, up to a maximum  
 532 value of 1.

533 The method is robust and has been widely applied to rough models ([29])  
 534 and predicted inter-residue distances that constitute highly non-metric data  
 535 sets ([30]).

## 536 References

- 537 [1] N. Ball, G. Nicastro, M. Dutta, D. Pollard, D. Goldstone, M. Sanz-  
 538 Ramos, A. Ramos, E. Mllers, K. Stirnnagel, N. Stanke, D. Lindemann,  
 539 J. Stoye, W. Taylor, P. Rosenthal, I. Taylor, Structure of a spumaretro-  
 540 virus gag central domain reveals an ancient retroviral capsid, PLoS  
 541 Path. ? (2016) ? Submitted.
- 542 [2] W. R. Taylor, Protein structure alignment using iterated double dy-  
 543 namic programming, Prot. Sci 8 (1999) 654–665.
- 544 [3] L. Holm, C. Sander, Protein-structure comparison by alignment of dis-  
 545 tance matrices, J. Molec. Biol. 233 (1993) 123–138.
- 546 [4] W. Zhang, J. Wu, M. Ward, S. Yang, Y. Chuang, M. Xiao, R. Li,  
 547 D. Leahy, P. Worley, Structural basis of arc binding to synaptic proteins:  
 548 implications for cognitive disease., Neuron 86 (2015) 490–500.
- 549 [5] W. R. Taylor, Protein structure domain identification, Prot. Engng. 12  
 550 (1999) 203–216.
- 551 [6] W. R. Taylor, Decoy models for protein structure score normalisation,  
 552 J. Molec. Biol. 357 (2006) 676–699.
- 553 [7] M. Levitt, M. Gerstein, A unified statistical framework for sequence  
 554 comparison and structure comparison, Proc Natl Acad Sci USA 95  
 555 (1998) 5913–5920.
- 556 [8] M. Campillos, T. Doerks, P. Shah, P. Bork, Computational characteri-  
 557 zation of multiple gag-like human proteins, Trends in Genetics 22 (2006)  
 558 285–589.

- 559 [9] C. Llorens, M. Fares, A. Moya, Relationships of gag-pol diversity be-  
560 tween Ty3/Gypsy and retroviridae LTR retroelements and the three  
561 kings hypothesis., *BMC Evol. Biol.* 8 (2008) e276.
- 562 [10] G. Obal, F. Trajtenberg, F. Carrion, L. Tome, N. Larrieux, X. Zhang,  
563 O. Pritsch, A. Buschiazio, Conformational plasticity of a native retrovi-  
564 ral capsid revealed by X-ray crystallography., *Science* 349 (2015) 95–98.  
565 DOI: 10.1126/science.aaa5182.
- 566 [11] T. Gamble, F. Vajdos, S. Yoo, D. Worthylake, M. Houseweart,  
567 W. Sundquist, C. Hill, Crystal structure of human cyclophilin A bound  
568 to the amino-terminal domain of HIV-1 capsid., *Cell* 87 (1996) 1285–  
569 1294.
- 570 [12] D. Worthylake, H. Wang, S. Yoo, W. Sundquist, C. Hill, Structures of  
571 the HIV-1 capsid protein dimerization domain at 2.6Å resolution., *Acta*  
572 *Crystallogr., Sect. D* 55 (1999) 85–92. DOI: 10.1107/S0907444998007689.
- 573 [13] O. Pornillos, B. Ganser-Pornillos, B. Kelly, Y. Hua, F. Whitby, C. Stout,  
574 W. Sundquist, C. Hill, M. Yeager, X-ray structures of the hexameric  
575 building block of the HIV capsid., *Cell* 137 (2009) 1282–1292. DOI:  
576 10.1016/j.cell.2009.04.063.
- 577 [14] G. Mortuza, M. Dodding, D. Goldstone, L. Haire, J. Stoye, I. Taylor,  
578 Structure of B-tropic MLV capsid N-terminal domain., *J. Mol. Biol.* 376  
579 (2008) 1493–1508.
- 580 [15] S. Khorasanizadeh, R. Campos-Olivas, C. Clark, M. Summers,  
581 Sequence-specific <sup>1</sup>H, <sup>13</sup>C and <sup>15</sup>N chemical shift assignment and sec-  
582 ondary structure of the HTLV-I capsid protein., *J. Biomol. NMR* 14  
583 (1999) 199–200.
- 584 [16] G. B. Mortuza, D. C. Goldstone, C. Pashley, L. F. Haire, M. Palmarini,  
585 W. R. Taylor, J. P. Stoye, I. A. Taylor, Structure of the capsid amino  
586 terminal domain from the betaretrovirus, Jaagsiekte sheep retrovirus.,  
587 *J. Molec. Biol.* 386 (2009) 1179–1192.
- 588 [17] G. B. Mortuza, L. F. Haire, A. Stevens, S. J. Smerdon, J. P. Stoye,  
589 I. A. Taylor, High-resolution structure of a retroviral capsid hexameric  
590 amino-terminal domain., *Nature* 431 (2004) 481–485.

- 591 [18] P. Macek, J. Chmelik, I. Krizova, P. Kaderavek, P. Padrta, L. Zidek,  
592 M. Wildova, R. Hadravova, R. Chaloupkova, I. Pichova, T. Ruml,  
593 M. Rumlova, V. Sklenar, NMR structure of the N-terminal domain  
594 of capsid protein from the mason-pfizer monkey virus, *J. Mol. Biol.* 392  
595 (2009) 100–114. DOI: 10.1016/j.jmb.2009.06.029.
- 596 [19] D. C. Goldstone, M. W. Yap, L. E. Robertson, L. F. Haire, W. R. Taylor,  
597 A. Katzourakis, J. P. Stoye, I. A. Taylor, Structural and functional anal-  
598 ysis of prehistoric lentiviruses uncovers an ancient molecular interface.,  
599 *Cell Host Microbe* 8 (2010) 248–259.
- 600 [20] G. Bailey, J. Hyun, A. Mitra, R. Kingston, Proton-linked dimerization  
601 of a retroviral capsid protein initiates capsid assembly, *Structure* 17  
602 (2009) 737–748. DOI: 10.1016/j.str.2009.03.010.
- 603 [21] F. Rippmann, W. R. Taylor, Visualization of structural similarity in  
604 proteins, *J. Molec. Graph.* 9 (1991) 3–16.
- 605 [22] A. D. McLachlan, How alike are the shapes of two random chains?,  
606 *Biopolymers* 23 (1984) 1325–1331.
- 607 [23] F. E. Cohen, M. J. E. Sternberg, On the prediction of protein structure:  
608 the significance of the root-mean-square deviation, *J. Molec. Biol.* 138  
609 (1980) 321–333.
- 610 [24] V. N. Maiorov, G. M. Crippen, Significance of root-mean-square devia-  
611 tion in comparing three-dimensional structures of globular proteins, *J.*  
612 *Mol. Biol.* 235 (1994) 625–634.
- 613 [25] W. H. Press, B. P. Flannery, S. A. Teukolsky, W. T. Vetterling, Numer-  
614 ical Recipes: The Art of Scientific Computing, Cambridge Univ. Press  
615 (Cambridge, UK), 1986.
- 616 [26] N. P. Brown, C. A. Orengo, W. R. Taylor, A protein structure compar-  
617 ison methodology, *Computers Chem.* 20 (1996) 359–380.
- 618 [27] A. Aszódi, W. R. Taylor, Hierarchical inertial projection: a fast distance  
619 matrix embedding algorithm., *Computers Chem.* 21 (1997) 13–23.
- 620 [28] W. R. Taylor, A. C. W. May, N. P. Brown, A. Aszódi, Protein structure:  
621 Geometry, topology and classification, *Rep. Prog. Phys.* 64 (2001) 517–  
622 590.

623 [29] W. R. Taylor, V. Chelliah, S. M. Hollup, J. T. MacDonald, I. Jonassen,  
624 Probing the “dark matter” of protein fold-space, *Structure* 17 (2009)  
625 1244–1252.

626 [30] A. Aszódi, W. R. Taylor, Folding polypeptide  $\alpha$ -carbon backbones by  
627 distance geometry methods, *Biopolymers* 34 (1994) 489–506.

628 *Acknowledgements:*. The work was supported by the Francis Crick Institute  
629 under awards: FC001179 (WRT), FC001162 (JPS) and FC001178 (IAT). The  
630 Crick receives its core funding from Cancer Research UK, the UK Medical  
631 Research Council, and the Wellcome Trust.