

A peer-reviewed version of this preprint was published in PeerJ on 1 March 2017.

[View the peer-reviewed version](https://doi.org/10.7717/peerj.3061) (peerj.com/articles/3061), which is the preferred citable publication unless you specifically need to cite this preprint.

Papageorgiou L, Megalooikonomou V, Vlachakis D. 2017. Genetic and structural study of DNA-directed RNA polymerase II of *Trypanosoma brucei*, towards the designing of novel antiparasitic agents. PeerJ 5:e3061 <https://doi.org/10.7717/peerj.3061>

Genetic and structural study of DNA-directed RNA polymerase II of *Trypanosoma brucei*, towards the designing of novel antiparasitic agents.

Louis Papageorgiou^{1,2,3}, Vasileios Megalooikonomou³, Dimitrios Vlachakis^{Corresp. 2,3}

¹ Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, Athens, Greece

² Computational Biology & Medicine Group, Biomedical Research Foundation, Academy of Athens, Athens, Attica, Greece

³ Computer Engineering and Informatics Department, University of Patras, Patra, Greece

Corresponding Author: Dimitrios Vlachakis
Email address: dvlachakis@bioacademy.gr

Trypanosoma brucei brucei (TBB) belongs to the unicellular parasitic protozoa organisms, specifically to the *Trypanosoma* genus of the *Trypanosomatidae* class. A variety of different vertebrate species can be infected by TBB including humans and animals. Under particular conditions, the TBB can be hosted by wild and domestic animals; thereby an important reservoir of infection always remains available to transmit through the tsetse flies. Although the TBB parasite is one of the leading causes of death in the most underdeveloped countries, to date, there is neither vaccination available nor any drug against TBB infection. The subunit RPB1 of the TBB DNA-directed RNA polymerase II (DdRpII) constitutes an ideal target for the design of novel inhibitors, since its instrumental role is vital for the parasite's survival, proliferation, and transmission. A major goal of the described study is to provide insights for novel anti-TBB agents via a state of the art drug discovery approach of the TBB DdRpII RPB1. In an attempt to understand the function and action mechanisms of this parasite enzyme related to its molecular structure, an in-depth evolutionary study has been conducted in parallel to the *in silico* molecular designing of the 3D enzyme model, based on state of the art comparative modelling and molecular dynamics techniques. Based on the evolutionary studies results nine new invariant, first-time reported, highly conserved regions have been identified within the DdRpII family enzymes. Consequently, those patches have been examined both at the sequence and structural level and have been evaluated in regards to their pharmacological targeting appropriateness. Finally, the pharmacophore elucidation study enabled us to virtually *in silico* screen hundreds of compounds and evaluate their interaction capabilities with the enzyme. It was found that a series of Chlorine-rich set of compounds were the optimal inhibitors for the TBB DdRpII RPB1 enzyme. All-in-all, herein we present a series of new sites on the TBB DdRpII RPB1 of high pharmacological interest, alongside the

construction of the 3D model of the enzyme and the suggestion of a new *in silico* pharmacophore model for fast screening of potential inhibiting agents.

1 **Genetic and structural study of DNA-directed RNA polymerase II of Trypanosoma brucei,**
2 **towards the designing of novel antiparasitic agents.**

3
4 **Louis Papageorgiou^{1,2,3}, Vasileios Megalooikonomou² and Dimitrios Vlachakis^{1,2,#}**

5
6
7 ¹Computational Biology & Medicine Group, Biomedical Research Foundation, Academy of Athens, Soranou
8 Efessiou 4, Athens 11527, Greece

9 ²Computer Engineering and Informatics Department, University of Patras, Patras University Campus, 26504,
10 Greece

11 ³Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, University
12 Campus, Athens, 15784, Greece

13
14
15

16

17

18

19

20

21

22

23

24

25

26

27

28 #Correspondence to:

29 Dimitrios Vlachakis,

30 Bioinformatics & Medical Informatics Team, Biomedical Research Foundation, Academy of
31 Athens, SoranouEfessiou 4, Athens 11527, Greece

32 Tel: + 30 210 6597 647, Fax: +30 210 6597 545

33 E-mail: dvlachakis@bioacademy.gr

34 **Abstract**

35 *Trypanosoma brucei brucei* (TBB) belongs to the unicellular parasitic protozoa organisms,
36 specifically to the *Trypanosoma* genus of the *Trypanosomatidae* class. A variety of different
37 vertebrate species can be infected by TBB including humans and animals. Under particular
38 conditions, the TBB can be hosted by wild and domestic animals; thereby an important
39 reservoir of infection always remains available to transmit through the tsetse flies. Although the
40 TBB parasite is one of the leading causes of death in the most underdeveloped countries, to
41 date, there is neither vaccination available nor any drug against TBB infection. TBB DNA-
42 dependent RNA polymerase II (DdRpII subunit RPB1) is an ideal target for the design of novel
43 inhibitors against TBB. This enzyme plays a critical role in parasite's survival, proliferation, and
44 transmission. A major goal of the described study is to provide insights for novel anti-TBB
45 agents via a state of the art drug discovery approach of the TBB DdRpII RPB1. In an attempt to
46 understand the function and action mechanisms of this parasite enzyme related to its
47 molecular structure, an in-depth evolutionary study has been conducted in parallel to the *in*
48 *silico* molecular designing of the 3D enzyme model, based on state of the art comparative
49 modelling and molecular dynamics techniques. Based on the evolutionary studies results nine
50 new invariant, first-time reported, highly conserved regions have been identified within the
51 DdRpII family enzymes. Consequently, those patches have been examined both at the sequence
52 and structural level and have been evaluated in regards to their pharmacological targeting
53 appropriateness. Finally, a 3D pharmacophore model was constructed specifically for the TBB
54 DdRpII RPB1 enzyme. All-in-all, herein we present a series of new sites on the TBB DdRpII RPB1
55 of high pharmacological interest, alongside the construction of the 3D model of the enzyme and
56 the suggestion of a new *in silico* pharmacophore model for fast screening of potential inhibiting
57 agents.

58

59

60 Introduction

61 African trypanosome parasites cause human sleeping sickness and nagana in Africa, Asia, and
62 South America. More than 95% of reported cases are caused by two subspecies of
63 *Trypanosoma brucei brucei* (TBB), the *Trypanosoma brucei gambiense* (TBG) and the
64 *Trypanosoma brucei rhodesiense* (TBR) which is found in western and central Africa (Berriman
65 et al. 2005; World Health Organization 2015). The parasitic infection is transmitted by tsetse
66 flies, which breed in warm and humid areas. Tsetse flies are found living in 36 countries in sub-
67 Saharan Africa, thus putting 60 million people at risk. Currently, about 10,000 new cases each
68 year are reported by the World Health Organization (WHO). Moreover, it is believed that many
69 cases are undiagnosed and unreported. Sleeping sickness can be curable with medication, but
70 it may be fatal if it is left untreated. It is estimated that Human deaths caused by Sleeping
71 sickness are of about 48,000 annually. Bites by the tsetse fly erupt into a red sore on the skin
72 and in the following weeks, the person may have to deal with several symptoms including
73 fever, swollen lymph glands, aching muscles, headaches, and irritability. In advanced stages,
74 the TBB parasite attacks the central nervous system of the host, and in general causes some
75 disorders in personality, circadian rhythm, serenity, speech, and difficulties in walking. Despite
76 the significant treatment advances for patients with sleeping sickness, the parasite's
77 progression is often inevitable and needs more treatment options. Until today, drugs can only
78 be used in the early stages of the disease and without providing 100% reassurance for full
79 convalescence of the patient (Ridley 2002; Ross et al. 2007; Trouiller et al. 2002). The TBB parasite
80 starts its activity after each invasion through its proteins, specifically with its replication
81 enzymes including helicases and polymerases. Such enzymes are ideal targets for inhibitor
82 design since those proteins are crucial for the TBB parasite survival. Being already in possession
83 of the widely known sequence of the DNA-dependent RNA polymerase II (DdRpII) RPB1 (Chung
84 et al. 1993) which plays a significant role in the replication of the parasite, our primary goal is to
85 suppress its function towards replication itself when it infects a human. Although TBB has been
86 reported many times in the past, the three-dimensional structure of its essential enzymes like
87 DdRpII remains unknown so far (Malvy & Chappuis 2011).

88 Protein structure has been found to be three to ten times more conserved than
89 sequence (Illergard et al. 2009). Thus, when possible, it is preferable to study an enzyme's 3D
90 structure rather than its sequence. Knowledge of the tertiary structure can assist in the
91 understanding of relationships between structure and function (Berg et al. 2002). Herein, the
92 three-dimensional structure of DdRpII subunit RPB1 has been modelled, in an effort to predict
93 the 3D molecular structure that is linked to the function of this enzyme (Bayele 2009; Koch et
94 al. 2016). Two molecular models have been constructed using conventional molecular
95 modelling techniques and two different homolog 3D structures as templates. The established
96 molecular models of the DdRpII RPB1 enzyme of TBB exhibits all known structural motifs that
97 are unique to the DdRpII RPB1 enzymes.

98 Upon successful completion of the 3D structure prediction of the TBB DdRpII RPB1
99 protein, molecular dynamics simulations have been performed to structurally improve and
100 benchmark the quality of the 3D models. Moreover, the reliability and viability of the TBB
101 DdRpII RPB1 models were checked using several *in silico* scoring tools such as MOE and
102 Procheck. After the model validation process, a *de novo* structure-based drug design approach
103 has been performed based on two models, which led to the establishment of a 3D novel

104 pharmacophore model that is highly specific for the DdRpII RPB1 enzyme of TBB. The generated
105 pharmacophore model may be used in future experiments involving the high throughput virtual
106 screening of large compound databases towards the identification of novel anti-TBB agents
107 (Loukatou et al. 2014). The present work opens the field for the design of novel compounds
108 with improved biochemical and clinical characteristics in the future.

109
110
111

112 **Methods**

113

114 **Database sequence search**

115 The full-length protein sequences related to the DdRpII family were extracted from the NCBI
116 database. In total, 36 DdRpII protein sequences were downloaded from several species with
117 fully sequenced genomes (Supplementary data 1).

118

119 **Genetic and evolutionary analyses**

120 Multiple sequence alignment of the DdRpII protein family sequences were performed using two
121 different programs, MUSCLE (Edgar 2004) and CLUSTALW (Chenna et al. 2003; Thompson et al.
122 1994). In the next step, multiple sequence alignment was checked with ProtTest3 (Darriba et al.
123 2011) to estimate the appropriate model of sequence evolution. Phylogenetic analyses were
124 performed by two different ways, and two representative phylogenetic trees were constructed
125 for the DdRpII dataset (Vlachakis et al. 2014b). The first phylogenetic tree was constructed
126 using the MEGA software (Stecher et al. 2014) utilizing Bayesian and Maximum Likelihood
127 statistical methods as described in with 100 bootstrap replicates (Figure 1 and Supplementary
128 data 2). The second phylogenetic tree was constructed using the Jalview software (Waterhouse
129 et al. 2009) utilizing the neighbour joining statistical method in with 100 bootstrap replicates
130 (Supplementary Figure 1 and Supplementary data 3).

131

132 **Conserved motifs exploration**

133 The phylogenetic trees that derived from the phylogenetic analyses (Jalview and MEGA) were
134 separated in sub-trees, in order to extract the most highly related protein sequences of the TBB
135 DdRpII RPB1 family for the conserved motifs exploration (Figure 2). The full-length amino acid
136 sequences of the closely related proteins with the TPP DdRpII RPB1 protein were aligned using
137 the CLUSTALW (Thompson et al. 1994) statistical method. The evolutionary conserved
138 sequences motifs that were derived from the multiple sequence alignment were identified
139 through the consensus sequence and logo graph where generated using Jalview (Waterhouse
140 et al. 2009) (Figure 2).

141

142 **Molecular modelling**

143 All calculations and visual constructions were performed using the Molecular Operating
144 Environment (MOE) version 2013.08 software package developed by Chemical Computing
145 Group (Montreal, Canada) on a cloud-based multi core High Performance Computing (HPC)
146 cluster (Loukatou et al. 2014).

147

148 Identification of templates structures and sequence alignment

149 The amino acid sequence of the TBB DdRpII RPB1 was retrieved from the conceptual translation
150 of the trypanosomal RNA polymerase largest subunit genes at the NCBI database
151 (<http://www.ncbi.nlm.nih.gov/>) (UniProtKB/Swiss-Prot: P17545.1) (Das et al. 2006; Evers et al.
152 1989). The blastp algorithm (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>) was used to identify
153 homologous structures by searching in the Protein Data Bank (PDB). The multiple sequence
154 alignment was performed using MOE (Vilar et al. 2008).

155

156 Homology Modelling

157 The homology modelling of the Tbb DdRpII RPB1 was carried out using MOE. The selection of
158 template crystal structures for homology modelling was based on the primary sequence
159 identity and similarity (Figure 3, Supplementary Figures 2 and 3), and the crystal resolution
160 (Nayeem et al. 2006). The crystal structure of *Schizosaccharomyces pombe* DdRpII RPB1 (PDB:
161 3H0G) was used as template structure for the model A, while the crystal structure of *Bos taurus*
162 DdRpII RPB1 (PDB: 5FLM) was used for building model B. The MOE homology model method is
163 separated into four main steps. First, comes a primary fragment geometry specification. Second
164 the insertion and deletions task. The third step is the loop selection and the side-chain packing,
165 and the last step is the final model selection and refinement (Figures 4, 5 and Supplementary
166 data 4 and 5) (Papageorgiou et al. 2014; Vlachakis et al. 2013b). Subsequently, energy
167 minimization was done in MOE initially using the Amber99 (Wang et al. 2000) force-field as
168 implemented into the same package. The energy minimization process was applied up to a
169 gradient of 0.0001, in an effort to remove the geometrical strain (Vlachakis et al. 2013a).

170

171 Molecular electrostatic potential

172 Molecular electrostatic potential surfaces were calculated by solving the non-linear Poisson –
173 Boltzmann equation using finite difference method as implemented into the MOE and PyMol
174 Software (Seeliger & de Groot 2010; Vilar et al. 2008). The potential was calculated on solid
175 points per side. Protein contact potential is an automated representation where the false
176 red/blue charge-smoothed surface is shown on the protein (Figure 4). Amber99 charges and
177 atomic radii were used for this calculation.

178

179 Molecular dynamics

180 The Molecular Dynamics simulations of both TBB DdRpII RPB1 3D models A and B were
181 executed in a periodic cell, which was explicitly solvated with simple point charge (SPC) water.
182 The truncated octahedron box was chosen for solvating the models, with a set distance of 7 Å
183 clear of the protein. The molecular dynamic simulations were conducted at 300 K, 1 atm with a
184 set 2 femtosecond step size for a total of one hundred nanoseconds. For the purposes of this study
185 we opted for a NVT ensemble in a canonical environment (Vlachakis et al. 2014a). NVT stands
186 for Number of atoms, Volume, and Temperature that remain constant throughout the
187 calculation (Vlachakis 2009). The intricate zinc ions were included in the molecular dynamics
188 simulations as integral parts of the modelled biological system (Chakravorty & Merz 2014;
189 Temiz et al. 2010). However, due to the nature of the ions, we had to limit the allowed degrees
190 of freedom for those molecules. Thus, the potential of the zinc ions was constrained in the
191 three dimensional conformational space in the vicinity of the TBB DdRpII RPB1 3D models. The

192 ions were prepositioned in the 3D models of TBB DdRpII RPB1, after structural superposition to
193 the template x-ray structure. The models were structurally optimized and adjusted locally by
194 subsequent energy minimizations, in an effort to eliminate any molecular clashes and minimize
195 the constrain energy. A radius of 6Å around each ion was given full degrees of freedom during
196 the abovementioned structural optimizations. Provided that the TBB DdRpII RPB1 is a
197 nucleotide processing enzyme, whose structure coordinates a repertoire of ions (e.g. Zinc,
198 Mg⁺⁺), the AMBER99 forcefield was selected (Figure 6). The AMBER99 forcefield is fully
199 parameterized for our biological system as it implements ff10 parameters for amino acids and
200 nucleic acids as well as EHT for small molecules, such as ions/cations at the same time (Vilar et
201 al. 2008). AM1-BCC charges were applied since the molecular system included the ion
202 molecules. The results of the molecular dynamics simulations for both models were collected
203 into a database by MOE for further analysis. The full simulation trajectories and molecular
204 dynamics graphs for both models are presented in Figure 7 and Supplementary Figures 4-7.

205

206 **Model evaluation**

207 The produced models were initially evaluated within the MOE package by a residue packing
208 quality function, which depends on the number of buried non-polar side-chain groups and on
209 hydrogen bonding. Moreover, the suite PROCHECK (Laskowski et al. 1996) was employed to
210 further evaluate the quality of the produced models. Finally, MOE and its build in protein check
211 module was used to evaluate whether the models of DdRpII RPB1 domains are similar to known
212 protein structures of this family (Supplementary data 6, 7 and 8).

213

214 **Pharmacophore Elucidation**

215 A pharmacophoric feature characterizes a particular property and is not tied to a specific
216 chemical structure; indeed different chemical groups may share the same property and so be
217 represented by the same feature (Vlachakis et al. 2013a). It is thus a mistake to name as
218 pharmacophoric features chemical functionalities such as guanidines or sulfonamides or typical
219 structural skeletons such as flavones or steroids.

220 The term pharmacophore modeling refers to the generation of a pharmacophore
221 hypothesis for the binding interactions in a particular active site (Vlachakis et al. 2015). Several
222 different pharmacophore models for the same active site can be overlaid and reduced to their
223 shared features so that common interactions are retained. Such a consensus pharmacophore
224 can be considered as the largest common denominator shared by a set of active molecules.

225 In MOE, the computerized representation of a hypothesized pharmacophore is called a
226 pharmacophore query. A MOE pharmacophore query is a set of query features that are
227 typically created from ligand annotation points. Annotation points are markers in space that
228 show the location and type of biologically important atoms and groups, such as hydrogen
229 donors and acceptors, aromatic centers, projected positions of possible interaction partners or
230 R-groups, charged groups, and bioisosteres. The annotation points on a ligand are the potential
231 locations of the features that will constitute the pharmacophore query. Annotation points
232 relevant to the pharmacophore are converted into query features with the addition of an extra
233 parameter: a non-zero radius that encodes the permissible variation in the pharmacophore
234 query's geometry.

235 Once generated, a pharmacophore query can be used to screen virtual compound libraries for
236 novel ligands. Pharmacophore queries can also be used to filter conformer databases, e.g.
237 output from molecular docking runs, for biologically active conformations.

238

239 Results

240

241 Phylogenetic Analysis

242 In the present study, two phylogenetic analyses of DdRpII family proteins in all available
243 genomes, with putative full-length protein sequences were performed using two different
244 statistical methods from the Jalview and MEGA software. Based on findings, putative members
245 of the DdRpII family were identified in the *Animalia*, *Fungi*, *Plantae*, *Protista* and
246 *Chromalveolata* kingdom major eukaryotic taxonomic division, as well as viruses (Figure 1 and
247 Supplementary Figure 1). In our analyses, in agreement with previous reports (Smith et al.
248 1989), we found that DdRpII family is split into two main subunits the RPB1 and the RPB2. The
249 two subunits of the DdRpII family are clearly separated in the phylogenetic trees as two major
250 sub-trees were obtained for each one of them (Figure 1 and Supplementary Figure 1). The
251 monophyletic sub-tree of the RPB1 subunit contains the TBB DdRpII RPB1, as well as another 17
252 leaves, which are related to RPB1 subunit. Furthermore, in the phylogenetic trees, the TBB
253 DdRpII RPB1 forms a distinct monophyletic branch with the *Euplotes octocarinatus* DdRpII RPB1
254 and the *Plasmodium falciparum* DdRpII RPB1, which is basal to a clade that corresponds to
255 other parasites. The Newick format of the phylogenetic trees is provided (Supplementary Data
256 1 and 2).

257

258 Conserved motifs exploration

259 Multiple sequence alignment of the DdRpII subunit RPB1 protein sequences from a variety of
260 several species were included in the first sub-tree, highlights important conserved functional
261 domains as described previously by Janet L. Smith and Judith R. Levin (Smith et al. 1989). Good
262 conservation is evident throughout the whole length of the sequence, especially among species
263 that belong to the same taxonomic division (Figure 2).

264 In this study, an effort has been done to suggest motifs that were probably included in
265 the DdRpII of the subunit RPB1. Regions conserved across all species (eukaryotic and viruses)
266 are indicative of important functional domains of the DdRpII RPB1 enzyme. Finally, the
267 consensus sequence of the multiple sequence alignment highlights nine conserved motifs
268 which are conserved between all species. All of the conserved motifs identified here have not
269 been reported previously, and indisputably deserve further study (Figures 2 and 3). It is
270 remarkable that all 18 polymerases, from the phylogenetic sub-tree of the subunit RPB1, have
271 high identity score and remain undamaged during the evolution (Figures 1 and 2). The highly
272 conserved motifs in protein families are directly related to their active sites and functionality
273 (Koonin & Galperin 2003; Papageorgiou et al. 2016).

274

275 3D models A and B of the *Trypanosoma brucei brucei* DdRpII RPB1

276 Homologous solved 3D structures from the Protein Data Bank (PDB) have been identified from
277 the Protein Data Bank (PDB) using the NCBI/BLASTp algorithm. Based on BLASTp report many
278 3D structures were determined suitable as templates for the homology modelling including the

279 crystal structure of the *Schizosaccharomyces pombe* DdRpII RPB1 (PDB: 3H0G) (Spahr et al.
280 2009), the crystal structure of the *Saccharomyces cerevisiae* DdRpII RPB1 (PDB: 4A3C and 1I3Q)
281 (Cheung et al. 2011; Cramer et al. 2001), the electron microscopy structure *Bos taurus* DdRpII
282 RPB1 (PDB: 5FLM) and the electron microscopy structure of the *Human* DdRpII RPB1 (PDB:
283 3JOK) (Bernecky et al. 2011). The final choice of a template structure was not only based on the
284 percent sequence identity/similarity and the structure resolution, but also on the results of the
285 phylogenetic trees. Two models were prepared. Model A was based on the
286 *Schizosaccharomyces pombe* DdRpII RPB1 x-ray structure, while model B was based on the *Bos*
287 *taurus* DdRpII RPB1 x-ray structure (Figure 3). Although the *Human* DdRpII RPB1 could also be
288 used to build the *Trypanosoma brucei* DdRpII RPB1 3D model, it was avoided in an effort to
289 minimize potential toxicity issues during the drug design process. Nonetheless, the sequence of
290 the *Human* DdRpII and the corresponding sequence of the *Trypanosoma brucei* and *Bos taurus*
291 were aligned in an effort to identify sequence-based differences and/or similarities for the
292 modelling and drug design process (Supplementary Figure 2). A multiple sequence alignment
293 was constructed including the *Trypanosoma brucei brucei* DdRpII RPB1 (NCBI: P17545.1) (Das et
294 al. 2006), the *Trypanosoma brucei gambiense* DdRpII RPB1 (NCBI: XP_011773113.1) (Jackson et
295 al. 2010), the crystal structure of *Schizosaccharomyces pombe* DdRpII RPB1 (PDB: 3H0G A chain)
296 (Spahr et al. 2009), the crystal structure of *Saccharomyces cerevisiae* DdRpII RPB1 (PDB: 1I3Q A
297 chain) (Cramer et al. 2001) , *Bos taurus* DdRpII RPB1 (PDB: 5FLM) (Bernecky et al. 2016). and the
298 crystal structure of *Human* DdRpII RPB1 (PDB: 3JOK A chain) (Bernecky et al. 2011) towards to
299 identify all the suggested conserved motifs within the highlighted domains of the RPB1 and the
300 major sequences differences and similarities (Supplementary Figure 2).

301 The above-mentioned sequence alignments were used to identify all the nine canonical
302 and conserved motifs as expected (Figures 2 and 3). The model of TPP DdRpII was first
303 structurally superimposed and subsequently structurally compared to its template using the
304 MOE software (Figure 4). The TPP DdRpII model exhibited an alpha-carbon RMSD lower than
305 1.3 angstroms (Figure 5 and Supplementary Data 8). Furthermore, the model was evaluated in
306 regards to its geometry and its compatibility with the template structure using the build in
307 protein check module of MOE (Supplementary Data 8). These results, confirmed the structural
308 viability of the 3D *in silico* model.

309

310 **Comparison of the *Trypanosoma brucei brucei* DdRpII RPB1 model A and model B.**

311 It was decided to produce two models using the aforementioned template structures. Model A
312 was build based on the *Schizosaccharomyces pombe* DdRpII RPB1 (PDB: 3H0G) X-ray structure
313 and model B was based on the *Bos taurus* DdRpII RPB1 (PDB: 5FLM) structure. *Bos taurus*
314 DdRpII RPB1 is a new released electron microscopy structure with 3.4 Å resolution, homolog to
315 *Trypanosoma brucei brucei* DdRpII RPB1. The sequence alignment between the *Trypanosoma*
316 *brucei brucei* DdRpII RPB1 and the *Bos taurus* DdRpII RPB1 template revealed 40% Identity and
317 56% similarity, same scores with the *Schizosaccharomyces pombe* DdRpII crystal structure, but
318 the overall sequence alignment length was shorter than the *Schizosaccharomyces pombe*
319 DdRpII crystal structure about 100 amino acids (Supplementary Figure 3). Furthermore, in the
320 sequence alignment of the *Trypanosoma brucei* DdRpII RPB1 and *Bos taurus* DdRpII RPB1 all
321 nine conserved motifs were identified, as expected. The root mean square deviation (RMSD)
322 between model A and its template is 1.3 Å whereas the RMSD between model B and *Bos taurus*

323 template is 2.7 Å. Nevertheless, the overall RMSD between the two models and the two
324 templates isn't bigger than 2,7 Å. (Figure 5 and Supplementary Data 8). Overall, we used to
325 prepare in parallel a 3D model based on the *Bos taurus* structure as it bears better validation
326 statistics and its sequence similarity to the *Trypanosoma brucei brucei* is higher. However, after
327 performing another full coarse of MDs for model B, it was concluded that the added value of
328 model B, when compared to model A is not significant, as models A and B are quite similar
329 indeed (Figure 7 and Supplementary Figures 4-7).

330

331 Discussion

332

333 Description of the *Trypanosoma brucei brucei* DdRpII RPB1 models.

334 RNA Polymerase II is a multi-subunit enzyme that transcribes protein-coding genes in
335 eukaryotes (Sentenac 1985). Transcription in eukaryotes is dependent by three classes of
336 nuclear RNA polymerases I-III. The genes encoding the largest subunits of eukaryotic RNA
337 polymerases I, II and III have been isolated and are single copy genes, except *Trypanosoma* RNA
338 polymerase II which contain two alleles (Smith et al. 1989). Structural and sequence differences
339 between the two alleles are minor, but the C-terminal domain of those enzymes has a highly
340 unusual structure. TBB DdRpII RPB1 model is the first protein subunit of the ten subunits multi-
341 complex of RNA Polymerase II (Hahn 2004; Suh et al. 2013). The RPB1 subunit is very critical in
342 RNA polymerase formation and function. The RPB1 active site and the RPB2 hybrid-binding
343 region combine in a single fold that forms the active centre of the RpII (Figure 4). There are two
344 metal ions at the RNA polymerase II active site. It has been previously reported that a Mg metal
345 ion interacts with the three invariant aspartates of RPB1 (Cramer et al. 2001). The latter
346 aspartate residues, which were found in all RPB1 sequences were aligned and fitted in a motifs
347 exploration study. Consequently, those residues have now been marked as motif 4b in the TBB
348 DdRpII RPB1 3D models.

349 The swinging motion of the clamp dictates the degree of opening of the cleft in DdRpII
350 and permits the insertion of promoter DNA for the initiation of transcription (Suh et al. 2013).
351 Based on previous studies, it is established that, upon closure of a transcribing complex, the
352 RPB1 clamp serves as a multi-functional tool, sensing the DNA/RNA hybrid conformation and
353 splitting DNA and RNA strands at the upstream end of the transcription complex (Cramer et al.
354 2001). The clamp is formed by N- and C-terminal regions of RPB1 and a part of the C-terminal
355 region of RPB2 (Chen et al. 2007; Hahn 2004; Li et al. 2014). The clamp is primarily stabilized by
356 three Zn ions within the RPB1 subunit (also marked in the TPP DdRpII RPB1) which forms zinc –
357 finger conformations; two within the “clamp core” and one in the “clamp head”. Accordingly,
358 two Zinc-finger formations were identified and highlighted in the TBB DdRpII RPB1 model
359 (Figure 6). The first formation can be recognized between a Zn ion and four cysteine residues in
360 the suggested motif 1a, also known as CX(2)CXnCX2C/H (Das et al. 2006) (Figure 3). Mutations
361 in the first Zn-finger formation confer a lethal phenotype of RNA polymerase II (Donaldson &
362 Friesen 2000). The second Zinc –finger can be recognized in the next four cysteine residues
363 (Figures 3 and 6). In the proposed motif 1b, the first two cysteine residues were identified,
364 which constitute part of the second Zing finger formation. Finally, according to our molecular
365 dynamics simulations, the main role of the Rpb1 and Rpb2 subunits is to provide stability within
366 the overall structure formation of the RNA polymerase II molecule in the 3D space.

367

368 **3D Pharmacophore Elucidation**

369 3D Pharmacophore design techniques take into account both the three-dimensional structures
370 and binding modes of receptors and inhibitors towards identifying regions that are favorable or
371 not for a particular receptor-inhibitor interaction (Vlachakis & Kossida 2013). The description of
372 the receptor-inhibitor interaction pattern is determined through a correlation between the
373 specific properties of the inhibitors and their action on enzymatic activity (Balatsos et al. 2009;
374 Vlachakis et al. 2012). The pharmacophore for TBB DdRpII RPB1 (Figure 8) was based on
375 structural information from the enzyme's catalytic site including all steric and electronic
376 features that are necessary to ensure optimal non-covalent interactions. The pharmacophoric
377 features were investigated including positively or negatively ionized regions, hydrogen bond
378 donors and acceptors, aromatic regions and hydrophobic areas. Firstly, there should be one
379 electron-donating group in the proximity of the Ser1172 (colored green). The electron-donating
380 region indicates a particular property of the inhibitor and is not necessarily confined to a
381 specific chemical structure. Moreover, this interaction site may not strictly represent a
382 hydrogen bond, but water or ion mediated bridges since the distance from the catalytic amino
383 acids varies between 3-9 Å. An aromatic PAP (colored orange) was positioned in the proximity
384 of Phe1179, which established pi-stacking interactions. Two electron accepting PAPs (colored
385 red) were positioned in the proximity of the two Arginine residues (Arg1171 and Arg1203).
386 Finally, a set of two adjacent PAPs were positioned in the center of the active site, where the
387 Zn⁺⁺ is coordinated in the crystal structure. Those yellow-colored PAPs are indicative of S-S
388 bonds and bridges or even S-C interactions, following the Michael acceptor moiety pattern. The
389 surrounding Cysteines are Cys1173, Cys1155, Cys1152, and Cys1270. However, the most
390 important factor of the latter PAPs was the optimal positioning of these groups in the 3D
391 conformational space of the TBB DdRpII RPB1 active site, rather than the amount of
392 conjugation or interaction with the protein.

393

394 **Conclusion**

395 The *Trypanosoma brucei brucei* DdRpII RPB1 enzyme was evolutionary analyzed, and nine new
396 conserved motifs were identified. Using the X-ray crystal structure of the *Schizosaccharomyces*
397 *pombe* DdRpII RPB1, the 3D model of the *Trypanosoma brucei brucei* DdRpII RPB1 was designed
398 using homology modelling techniques. The model was *in silico* evaluated and displayed high
399 conservation of the functional domains previously reported in other DdRpII subunit RPB1
400 species. The *Trypanosoma brucei brucei* DdRpII RPB1 model structure provides a basis for
401 interpretation of available data and the design of new experiments towards the *Trypanosoma*
402 *brucei brucei* inhibition. We, therefore, propose the use of the *Trypanosoma brucei brucei*
403 DdRpII RPB1 model A as a pharmacological targeting platform for advanced, *in silico* drug
404 design experiments using the novel findings of this study, both in the sequence and structural
405 level. The 3D models and sequence datasets that derived from this study will be made available
406 to the public, in an effort to pave the way for fellow scientists of multidiscipline backgrounds to
407 word in a synergic way towards the designing of novel anti-malarial agents with improved
408 biochemical and clinical characteristics in the future.

409

410

411 **Abbreviations**

412	DdRpII	DNA-directed RNA polymerase II
413	TBB	Trypanosoma brucei brucei
414	TBG	Trypanosoma brucei gambiense
415	TBR	Trypanosoma brucei rhodesiense
416	MOE	Molecular Operating Environment

418 **Figures and Data legend**

419

420 **Figure 1: Phylogenetic reconstruction of *Trypanosoma brucei brucei* DdRpII RPB1 protein**
421 **sequences.** The tree was generated using the DdRpII family dataset (36 full length protein
422 sequences samples). The tree was constructed by Matlab Bioinformatics Toolbox utilizing
423 Neighbour – Joining statistical method for 100 bootstrap replicates and visualized using MEGA
424 cycle option. In the tree representation there are clearly separated in two monophyletic
425 branches the RNA polymerases II subunits RPB1 (colored green) and RPB2 (colored blue).
426 *Trypanosoma brucei* DdRpII RPB1 protein sequence was correctly classified and separated in
427 the monophyletic sub-tree of the RPB1 group (highlight with red dots).

428

429

430 **Figure 2: Representative conserved motifs for the DdRpII subunit RPB1.** The nine suggested
431 conserved motifs were extracted based on the multiple sequence alignment of the 18 protein
432 sequences were classified and clearly separated in the DdRpII subunit RPB1 monophyletic sub-
433 tree. The conserved motifs were identified through the consensus sequence and logo graph
434 where generated using Jalview software.

435

436 **Figure 3: Sequence alignment between the *Trypanosoma brucei brucei* DdRpII RPB1 and the**
437 **corresponding sequence of the crystal structure of the *Schizosaccharomyces pombe* DdRpII**
438 **RPB1. (A)** Alignment of DdRpII RPB1 from *Trypanosoma brucei* DdRpII RPB1 (Labeled as “TB”)
439 with *Schizosaccharomyces pombe* DdRpII RPB1 (Labeled as “SB”) was initially carried out with
440 BLASTp and then manually adjusted. The nine suggested conserved motifs (Motifs 1a, 1b, 2, 3a,
441 3b, 3c, 4a, 4b, 4c) based on figure 2, domains and domain-like regions of *Trypanosoma brucei*
442 DdRpII RPB1 represented in different colours. The amino acid residue numbers at the domain
443 boundaries are indicated. Important structural elements and prominent regions involved in
444 subunit interactions are also noted. Residues involved in the Zn and Mg coordination are
445 highlighted in blue. **(B)** Domains and domain-like regions of the DdRpII subunit Rpb1. The
446 amino acid residue numbers at the domain boundaries are indicated.

447

448 **Figure 4: Model of the *Trypanosoma brucei brucei* DdRpII RPB1. (A and B)** Ribbon
449 representation of the produced *Trypanosoma brucei brucei* DdRpII RPB1 model (colored
450 Orange) superposed with the corresponding *Schizosaccharomyces pombe* DdRpII RPB1 (in
451 purple). **(C and D)** The nine suggested conserved motifs and the domains and domain-like
452 regions of the *Trypanosoma brucei brucei* DdRpII RPB1. The motifs and RPB1 domains have
453 been color-coded according to the Figures 2 and 3, and are shown in CPK format (Usual space
454 filling). **(E and F)** Electrostatic surface potential for the *Trypanosoma brucei brucei* DdRpII RPB1.
455 Represented with blue is the area of negative charge. Red is the area of positive charge and
456 white is the un-charged region.

457

458 **Figure 5: Structural superposition of the TBB DdRpII RPB1 models A and B. (A and B)** Ribbon
459 representation of the produced *Trypanosoma brucei brucei* DdRpII RPB1 model A (colored
460 Orange) and model B (colored Blue) superposed with the corresponding *Schizosaccharomyces*

461 *pombe* DdRpII RPB1 (in Purple) and *Bos taurus* DdRpII RPB1 (in Grey). The four 3D structures are
462 highly conserved in their active sites with few differences in the outer layer with overall RMSD
463 2.775 Å. **(C)** Ribbon representation of the produced *Trypanosoma brucei brucei* DdRpII RPB1
464 model A (colored Orange) superposed with the corresponding *Schizosaccharomyces pombe*
465 DdRpII RPB1 (in purple). (RMSD = 1.242 Å). **(D)** Ribbon representation of the produced
466 *Trypanosoma brucei brucei* DdRpII RPB1 model B (colored Blue) superposed with the *Bos taurus*
467 DdRpII RPB1 (in Grey) respectively. (RMSD = 2.757 Å).

468

469 **Figure 6: Zinc-finger formations in the *Trypanosoma brucei brucei* DdRpII RPB1 model.** Ribbon
470 representation of the produced *Trypanosoma brucei brucei* DdRpII RPB1 model. In the
471 produced model were highlighted 3 main zing-finger domain formations (colored grey) were
472 contained in the clam core, clam head and active site region. Domains and domain-like regions
473 of the *Trypanosoma brucei brucei* DdRpII RPB1 have been color-coded according to conventions
474 of Figure 3.

475

476 **Figure 7: Molecular dynamics simulation charts for the *Trypanosoma brucei brucei* DdRpII**
477 **RPB1 models. (A)** The root mean square deviation (RMSD) of the model A during the time. **(B)**
478 The root mean square fluctuation (RMSF) of the model A during the time. **(C)** The root mean
479 square deviation (RMSD) of the model B during the time. **(D)** The root mean square fluctuation
480 (RMSF) of the model B during the time.

481

482 **Figure 8: The 3D pharmacophore model for the *Trypanosoma brucei brucei* DdRpII RPB1**
483 **model.** In total 5 distinct pharmacophoric features were identified. An aromatic region (colored
484 orange), an electron donating region (colored green), two electron accepting regions (colored
485 red) and a sulphur specific S-S interacting region (colored yellow).

486

487 **Supplementary Figure 1: Phylogenetic reconstruction of *Trypanosoma brucei brucei* DdRpII**
488 **RPB1 model DdRpII RPB1 protein sequences.** The tree was generated using the DdRpII family
489 dataset (36 fowl length protein sequences samples) and the Jalview software. Tree was
490 constructed using the average distance statistical method with PAM 250. In the tree
491 representation there are clearly shown the two RNA polymerases II subunits RPB1 and RPB2 as
492 two main monophyletic sub-trees. *Trypanosoma brucei* DdRpII RPB1 protein sequence was
493 correctly classified in the monophyletic sub-tree of the RPB1 group.

494

495 **Supplementary Figure 2: Multiple sequence alignment.** The alignment was performed using
496 the *Trypanosoma brucei brucei* DdRpII RPB1, the *Trypanosoma brucei gambiense* DdRpII RPB1,
497 the crystal structure of *Schizosaccharomyces pombe* DdRpII RPB, the crystal structure of
498 *Saccharomyces cerevisiae* DdRpII RPB1 and the electron microscopy structure of Human DdRpII
499 DdRpII RPB1. **(A)** All nine suggested conserved motifs and major domains of DdRpII RPB1 have
500 been marked (Motifs 1a, 1b, 2, 3a, 3b, 3c, 4a, 4b, 4c). Additionally, in the multiple sequence
501 alignment were presented the major differences. **(B)** Domains and domainlike regions of the
502 DdRpII subunit Rpb1. The amino acid residue numbers at the domain boundaries are indicated.

503

504 **Supplementary Figure 3: Multiple sequence alignment.** The alignment was performed using
505 the *Trypanosoma brucei brucei* DdRpII RPB1, the crystal structure of *Schizosaccharomyces*
506 *pombe* DdRpII RPB and the electron microscopy structure of *Bos taurus* DdRpII RPB1. All five
507 sub-domains (A-E) as referred in Pfam database have been marked with different colours.

508

509 **Supplementary Figure 4: Molecular dynamics simulation charts of the root mean square**
510 **deviation (RMSD) for the *Trypanosoma brucei brucei* DdRpII RPB1 sub domains of the model**
511 **A.** The energy (Kcal/mol) vs time (ns) plot of the 100ns simulation trajectory of the TBB DdRpII
512 RPB1 model A. Sub-domain regions of the *Trypanosoma brucei brucei* DdRpII RPB1 have been
513 separated according to conventions of Supplementary Figure 3. **(A)** Domain A RMSD. **(B)**
514 Domain B RMSD. **(C)** Domain C RMSD. **(D)** Domain D RMSD. **(E)** Domain E RMSD.

515

516 **Supplementary Figure 5: Molecular dynamics simulation charts of the root mean square**
517 **fluctuation (RMSF) for the *Trypanosoma brucei brucei* DdRpII RPB1 sub domains of the model**
518 **A.** Sub-domain regions of the *Trypanosoma brucei brucei* DdRpII RPB1 have been separated
519 according to conventions of Supplementary Figure 3. **(A)** Domain A RMSF. **(B)** Domain B RMSF.
520 **(C)** Domain C RMSF. **(D)** Domain D RMSF. **(E)** Domain E RMSF.

521

522

523 **Supplementary Figure 6: Molecular dynamics simulation charts of the root mean square**
524 **deviation (RMSD) for the *Trypanosoma brucei brucei* DdRpII RPB1 sub domains of the model**
525 **B.** The energy (Kcal/mol) vs time (ns) plot of the 100ns simulation trajectory of the TBB DdRpII
526 RPB1 model B. Sub-domain regions of the *Trypanosoma brucei brucei* DdRpII RPB1 have been
527 separated according to conventions of Supplementary Figure 3. **(A)** Domain A RMSD. **(B)**
528 Domain B RMSD. **(C)** Domain C RMSD. **(D)** Domain D RMSD. **(E)** Domain E RMSD.

529

530 **Supplementary Figure 7: Molecular dynamics simulation charts of the root mean square**
531 **fluctuation (RMSF) for the *Trypanosoma brucei brucei* DdRpII RPB1 sub domains of the model**
532 **B.** Sub-domain regions of the *Trypanosoma brucei brucei* DdRpII RPB1 have been separated
533 according to conventions of Supplementary Figure 3. **(A)** Domain A RMSF. **(B)** Domain B RMSF.
534 **(C)** Domain C RMSF. **(D)** Domain D RMSF. **(E)** Domain E RMSF.

535

536

537 **Supplementary Data 1: DdRpII related proteins dataset.**

538

539 **Supplementary Data 2: MEGA software phylogenetic tree in newick format.** The tree was
540 constructed the Neighbour – Joining statistical method for 100 bootstrap replicates and the 36
541 extracted samples of the DpRpII.

542

543 **Supplementary Data 3: Jalview software phylogenetic tree in newick format.** The tree was
544 constructed using the average distances statistical method and the 36 extracted samples of the
545 DpRpII.

546

547 **Supplementary Data 4: *Trypanosoma brucei brucei* DdRpII RPB1 model A in .pdb format.**

548

549 **Supplementary Data 5: *Trypanosoma brucei brucei* DdRPII RPB1 model B in .pdb format.**

550

551 **Supplementary Data 6: Protein structure report of the template.**

552

553 **Supplementary Data 7: Protein structure report of the model.**

554

555 **Supplementary Data 8: Protein structure report of the superposed models and templates.**

556

557

558 **References**

559

560 Balatsos NA, Vlachakis D, Maragozidis P, Manta S, Anastasakis D, Kyritsis A, Vlassi M, Komiotis D, and
561 Stathopoulos C. 2009. Competitive inhibition of human poly(A)-specific ribonuclease (PARN) by
562 synthetic fluoro-pyranosyl nucleosides. *Biochemistry* 48:6044-6051. 10.1021/bi900236k

563 Bayele HK. 2009. Trypanosoma brucei: a putative RNA polymerase II promoter. *Exp Parasitol* 123:313-
564 318. 10.1016/j.exppara.2009.08.007

565 S0014-4894(09)00233-1 [pii]

566 Berg JM, Tymoczko JL, and Stryer L. 2002. *Biochemistry*. New York: W.H. Freeman.

567 Bernecky C, Grob P, Ebmeier CC, Nogales E, and Taatjes DJ. 2011. Molecular architecture of the human
568 Mediator-RNA polymerase II-TFIIF assembly. *PLoS Biol* 9:e1000603.
569 10.1371/journal.pbio.1000603

570 Bernecky C, Herzog F, Baumeister W, Plitzko JM, and Cramer P. 2016. Structure of transcribing
571 mammalian RNA polymerase II. *Nature* 529:551-554. 10.1038/nature16482

572 nature16482 [pii]

573 Berriman M, Ghedin E, Hertz-Fowler C, Blandin G, Renauld H, Bartholomeu DC, Lennard NJ, Caler E,
574 Hamlin NE, Haas B, Bohme U, Hannick L, Aslett MA, Shallom J, Marcello L, Hou L, Wickstead B,
575 Alsmark UC, Arrowsmith C, Atkin RJ, Barron AJ, Bringaud F, Brooks K, Carrington M, Cherevach I,
576 Chillingworth TJ, Churcher C, Clark LN, Corton CH, Cronin A, Davies RM, Doggett J, Djikeng A,
577 Feldblyum T, Field MC, Fraser A, Goodhead I, Hance Z, Harper D, Harris BR, Hauser H, Hostetler
578 J, Ivens A, Jagels K, Johnson D, Johnson J, Jones K, Kerhornou AX, Koo H, Larke N, Landfear S,
579 Larkin C, Leech V, Line A, Lord A, Macleod A, Mooney PJ, Moule S, Martin DM, Morgan GW,
580 Mungall K, Norbertczak H, Ormond D, Pai G, Peacock CS, Peterson J, Quail MA, Rabbinowitsch E,
581 Rajandream MA, Reitter C, Salzberg SL, Sanders M, Schobel S, Sharp S, Simmonds M, Simpson
582 AJ, Tallon L, Turner CM, Tait A, Tivey AR, Van Aken S, Walker D, Wanless D, Wang S, White B,
583 White O, Whitehead S, Woodward J, Wortman J, Adams MD, Embley TM, Gull K, Ullu E, Barry JD,
584 Fairlamb AH, Opperdoes F, Barrell BG, Donelson JE, Hall N, Fraser CM, Melville SE, and El-Sayed
585 NM. 2005. The genome of the African trypanosome *Trypanosoma brucei*. *Science* 309:416-422.
586 10.1126/science.1112642

587 Chakravorty DK, and Merz KM, Jr. 2014. Studying allosteric regulation in metal sensor proteins using
588 computational methods. *Adv Protein Chem Struct Biol* 96:181-218.
589 10.1016/bs.apcsb.2014.06.009

590 S1876-1623(14)00010-8 [pii]

591 Chen HT, Warfield L, and Hahn S. 2007. The positions of TFIIF and TFIIE in the RNA polymerase II
592 transcription preinitiation complex. *Nat Struct Mol Biol* 14:696-703. nsmb1272 [pii]

593 10.1038/nsmb1272

594 Chenna R, Sugawara H, Koike T, Lopez R, Gibson TJ, Higgins DG, and Thompson JD. 2003. Multiple
595 sequence alignment with the Clustal series of programs. *Nucleic Acids Res* 31:3497-3500.

596 Cheung AC, Sainsbury S, and Cramer P. 2011. Structural basis of initial RNA polymerase II transcription.
597 *EMBO J* 30:4755-4763. 10.1038/emboj.2011.396

598 emboj2011396 [pii]

599 Chung HM, Lee MG, Dietrich P, Huang J, and Van der Ploeg LH. 1993. Disruption of largest subunit RNA
600 polymerase II genes in *Trypanosoma brucei*. *Mol Cell Biol* 13:3734-3743.

601 Cramer P, Bushnell DA, and Kornberg RD. 2001. Structural basis of transcription: RNA polymerase II at
602 2.8 angstrom resolution. *Science* 292:1863-1876. 10.1126/science.1059493

603 1059493 [pii]

604 Darriba D, Taboada GL, Doallo R, and Posada D. 2011. ProtTest 3: fast selection of best-fit models of
605 protein evolution. *Bioinformatics* 27:1164-1165. 10.1093/bioinformatics/btr088

- 606 btr088 [pii]
607 Das A, Li H, Liu T, and Bellofatto V. 2006. Biochemical characterization of Trypanosoma brucei RNA
608 polymerase II. *Mol Biochem Parasitol* 150:201-210. S0166-6851(06)00232-5 [pii]
609 10.1016/j.molbiopara.2006.08.002
610 Donaldson IM, and Friesen JD. 2000. Zinc stoichiometry of yeast RNA polymerase II and characterization
611 of mutations in the zinc-binding domain of the largest subunit. *J Biol Chem* 275:13780-13788.
612 275/18/13780 [pii]
613 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic
614 Acids Res* 32:1792-1797. 10.1093/nar/gkh340
615 32/5/1792 [pii]
616 Evers R, Hammer A, Kock J, Jess W, Borst P, Memet S, and Cornelissen AW. 1989. Trypanosoma brucei
617 contains two RNA polymerase II largest subunit genes with an altered C-terminal domain. *Cell*
618 56:585-597. 0092-8674(89)90581-3 [pii]
619 Hahn S. 2004. Structure and mechanism of the RNA polymerase II transcription machinery. *Nat Struct
620 Mol Biol* 11:394-403. 10.1038/nsmb763
621 nsmb763 [pii]
622 Illergard K, Ardell DH, and Elofsson A. 2009. Structure is three to ten times more conserved than
623 sequence--a study of structural response in protein cores. *Proteins* 77:499-508.
624 10.1002/prot.22458
625 Jackson AP, Sanders M, Berry A, McQuillan J, Aslett MA, Quail MA, Chukualim B, Capewell P, MacLeod A,
626 Melville SE, Gibson W, Barry JD, Berriman M, and Hertz-Fowler C. 2010. The genome sequence
627 of Trypanosoma brucei gambiense, causative agent of chronic human african trypanosomiasis.
628 *PLoS Negl Trop Dis* 4:e658. 10.1371/journal.pntd.0000658
629 Koch H, Raabe M, Urlaub H, Bindereif A, and Preusser C. 2016. The polyadenylation complex of
630 Trypanosoma brucei: Characterization of the functional poly(A) polymerase. *RNA Biol* 13:221-
631 231. 10.1080/15476286.2015.1130208
632 Koonin EV, and Galperin MY. 2003. *Sequence - evolution - function : computational approaches in
633 comparative genomics*. Boston: Kluwer Academic.
634 Laskowski RA, Rullmannn JA, MacArthur MW, Kaptein R, and Thornton JM. 1996. AQUA and PROCHECK-
635 NMR: programs for checking the quality of protein structures solved by NMR. *J Biomol NMR*
636 8:477-486.
637 Li W, Giles C, and Li S. 2014. Insights into how Spt5 functions in transcription elongation and repressing
638 transcription coupled DNA repair. *Nucleic Acids Res* 42:7069-7083. 10.1093/nar/gku333
639 gku333 [pii]
640 Loukatou S, Papageorgiou L, Fakourelis P, Filntisi A, Polychronidou E, Bassis I, Megalooikonomou V,
641 Makalowski W, Vlachakis D, and Kossida S. 2014. Molecular dynamics simulations through GPU
642 video games technologies. *J Mol Biochem* 3:64-71.
643 Malvy D, and Chappuis F. 2011. Sleeping sickness. *Clin Microbiol Infect* 17:986-995. 10.1111/j.1469-
644 0691.2011.03536.x
645 S1198-743X(14)61376-8 [pii]
646 Nayeem A, Sitkoff D, and Krystek S, Jr. 2006. A comparative study of available software for high-accuracy
647 homology modeling: from sequence alignments to structural models. *Protein Sci* 15:808-824.
648 10.1110/ps.051892906
649 Papageorgiou L, Loukatou S, Koumandou VL, Makalowski W, Megalooikonomou V, Vlachakis D, and
650 Kossida S. 2014. Structural models for the design of novel antiviral agents against Greek Goat
651 Encephalitis. *PeerJ* 2:e664. 10.7717/peerj.664

- 652 Papageorgiou L, Loukatou S, Sofia K, Maroulis D, and Vlachakis D. 2016. An updated evolutionary study
653 of Flaviviridae NS3 helicase and NS5 RNA-dependent RNA polymerase reveals novel invariable
654 motifs as potential pharmacological targets. *Mol Biosyst*. 10.1039/c5mb00706b
- 655 Ridley RG. 2002. Medical need, scientific opportunity and the drive for antimalarial drugs. *Nature*
656 415:686-693. 10.1038/415686a
657 415686a [pii]
- 658 Ross L, Lim ML, Liao Q, Wine B, Rodriguez AE, Weinberg W, and Shaefer M. 2007. Prevalence of
659 antiretroviral drug resistance and resistance-associated mutations in antiretroviral therapy-
660 naive HIV-infected individuals from 40 United States cities. *HIV Clin Trials* 8:1-8.
661 L61238775J512643 [pii]
662 10.1310/hct0801-1
- 663 Seeliger D, and de Groot BL. 2010. Ligand docking and binding site analysis with PyMOL and
664 Autodock/Vina. *J Comput Aided Mol Des* 24:417-422. 10.1007/s10822-010-9352-6
- 665 Sentenac A. 1985. Eukaryotic RNA polymerases. *CRC Crit Rev Biochem* 18:31-90.
- 666 Smith JL, Levin JR, Ingles CJ, and Agabian N. 1989. In trypanosomes the homolog of the largest subunit of
667 RNA polymerase II is encoded by two genes and has a highly unusual C-terminal domain
668 structure. *Cell* 56:815-827.
- 669 Spahr H, Calero G, Bushnell DA, and Kornberg RD. 2009. Schizosaccharomyces pombe RNA polymerase II
670 at 3.6-Å resolution. *Proc Natl Acad Sci U S A* 106:9185-9190. 10.1073/pnas.0903361106
- 671 Stecher G, Liu L, Sanderford M, Peterson D, Tamura K, and Kumar S. 2014. MEGA-MD: molecular
672 evolutionary genetics analysis software with mutational diagnosis of amino acid variation.
673 *Bioinformatics* 30:1305-1307. 10.1093/bioinformatics/btu018
674 btu018 [pii]
- 675 Suh H, Hazelbaker DZ, Soares LM, and Buratowski S. 2013. The C-terminal domain of Rpb1 functions on
676 other RNA polymerase II subunits. *Mol Cell* 51:850-858. 10.1016/j.molcel.2013.08.015
677 S1097-2765(13)00588-1 [pii]
- 678 Temiz AN, Benos PV, and Camacho CJ. 2010. Electrostatic hot spot on DNA-binding domains mediates
679 phosphate desolvation and the pre-organization of specificity determinant side chains. *Nucleic*
680 *Acids Res* 38:2134-2144. 10.1093/nar/gkp1132
681 gkp1132 [pii]
- 682 Thompson JD, Higgins DG, and Gibson TJ. 1994. CLUSTAL W: improving the sensitivity of progressive
683 multiple sequence alignment through sequence weighting, position-specific gap penalties and
684 weight matrix choice. *Nucleic Acids Res* 22:4673-4680.
- 685 Trouiller P, Olliaro P, Torreele E, Orbinski J, Laing R, and Ford N. 2002. Drug development for neglected
686 diseases: a deficient market and a public-health policy failure. *Lancet* 359:2188-2194. S0140-
687 6736(02)09096-7 [pii]
688 10.1016/S0140-6736(02)09096-7
- 689 Vilar S, Cozza G, and Moro S. 2008. Medicinal chemistry and the molecular operating environment
690 (MOE): application of QSAR and molecular docking to drug discovery. *Curr Top Med Chem*
691 8:1555-1572.
- 692 Vlachakis D. 2009. Theoretical study of the Usutu virus helicase 3D structure, by means of computer-
693 aided homology modelling. *Theor Biol Med Model* 6:9. 10.1186/1742-4682-6-9
694 1742-4682-6-9 [pii]
- 695 Vlachakis D, Bencurova E, Papangelopoulos N, and Kossida S. 2014a. Current state-of-the-art molecular
696 dynamics methods and applications. *Adv Protein Chem Struct Biol* 94:269-313. 10.1016/B978-0-
697 12-800168-4.00007-X
698 B978-0-12-800168-4.00007-X [pii]

- 699 Vlachakis D, Fakourelis P, Megalooikonomou V, Makris C, and Kossida S. 2015. DrugOn: a fully integrated
700 pharmacophore modeling and structure optimization toolkit. *PeerJ* 3:e725. 10.7717/peerj.725
701 725 [pii]
- 702 Vlachakis D, Kontopoulos DG, and Kossida S. 2013a. Space constrained homology modelling: the
703 paradigm of the RNA-dependent RNA polymerase of dengue (type II) virus. *Comput Math*
704 *Methods Med* 2013:108910. 10.1155/2013/108910
- 705 Vlachakis D, and Kossida S. 2013. Molecular modeling and pharmacophore elucidation study of the
706 Classical Swine Fever virus helicase as a promising pharmacological target. *PeerJ* 1:e85.
707 10.7717/peerj.85
- 708 Vlachakis D, Koumandou VL, and Kossida S. 2013b. A holistic evolutionary and structural study of
709 flaviviridae provides insights into the function and inhibition of HCV helicase. *PeerJ* 1:e74.
710 10.7717/peerj.74
- 711 Vlachakis D, Pavlopoulou A, Roubelakis MG, Feidakis C, Anagnou NP, and Kossida S. 2014b. 3D molecular
712 modeling and evolutionary study of the *Trypanosoma brucei* DNA Topoisomerase IB, as a new
713 emerging pharmacological target. *Genomics* 103:107-113. 10.1016/j.ygeno.2013.11.008
714 S0888-7543(13)00224-3 [pii]
- 715 Vlachakis D, Pavlopoulou A, Tsiliki G, Komiotis D, Stathopoulos C, Balatsos NA, and Kossida S. 2012. An
716 integrated in silico approach to design specific inhibitors targeting human poly(a)-specific
717 ribonuclease. *PLoS One* 7:e51113. 10.1371/journal.pone.0051113
718 PONE-D-12-20013 [pii]
- 719 Wang J, Cieplak P, and Kollman P. 2000. How well does a restrained electrostatic potential (resp) model
720 perform in calculating conformational energies of organic and biological molecules. *J Comp*
721 *Chem* 21:1049-1071.
- 722 Waterhouse AM, Procter JB, Martin DM, Clamp M, and Barton GJ. 2009. Jalview Version 2--a multiple
723 sequence alignment editor and analysis workbench. *Bioinformatics* 25:1189-1191.
724 10.1093/bioinformatics/btp033
725 btp033 [pii]
- 726 World Health Organization. 2015. Human African trypanosomiasis
727
728

Figure 1

Phylogenetic reconstruction of *Trypanosoma brucei brucei* DdRpII RPB1 protein sequences.

The tree was generated using the DdRpII family dataset (36 full length protein sequences samples). The tree was constructed by Matlab Bioinformatics Toolbox utilizing Neighbour - Joining statistical method for 100 bootstrap replicates and visualized using MEGA cycle option. In the tree representation there are clearly separated in two monophyletic branches the RNA polymerases II subunits RPB1 (colored green) and RPB2 (colored blue).

Trypanosoma brucei DdRpII RPB1 protein sequence was correctly classified and separated in the monophyletic sub-tree of the RPB1 group (highlight with red dots).

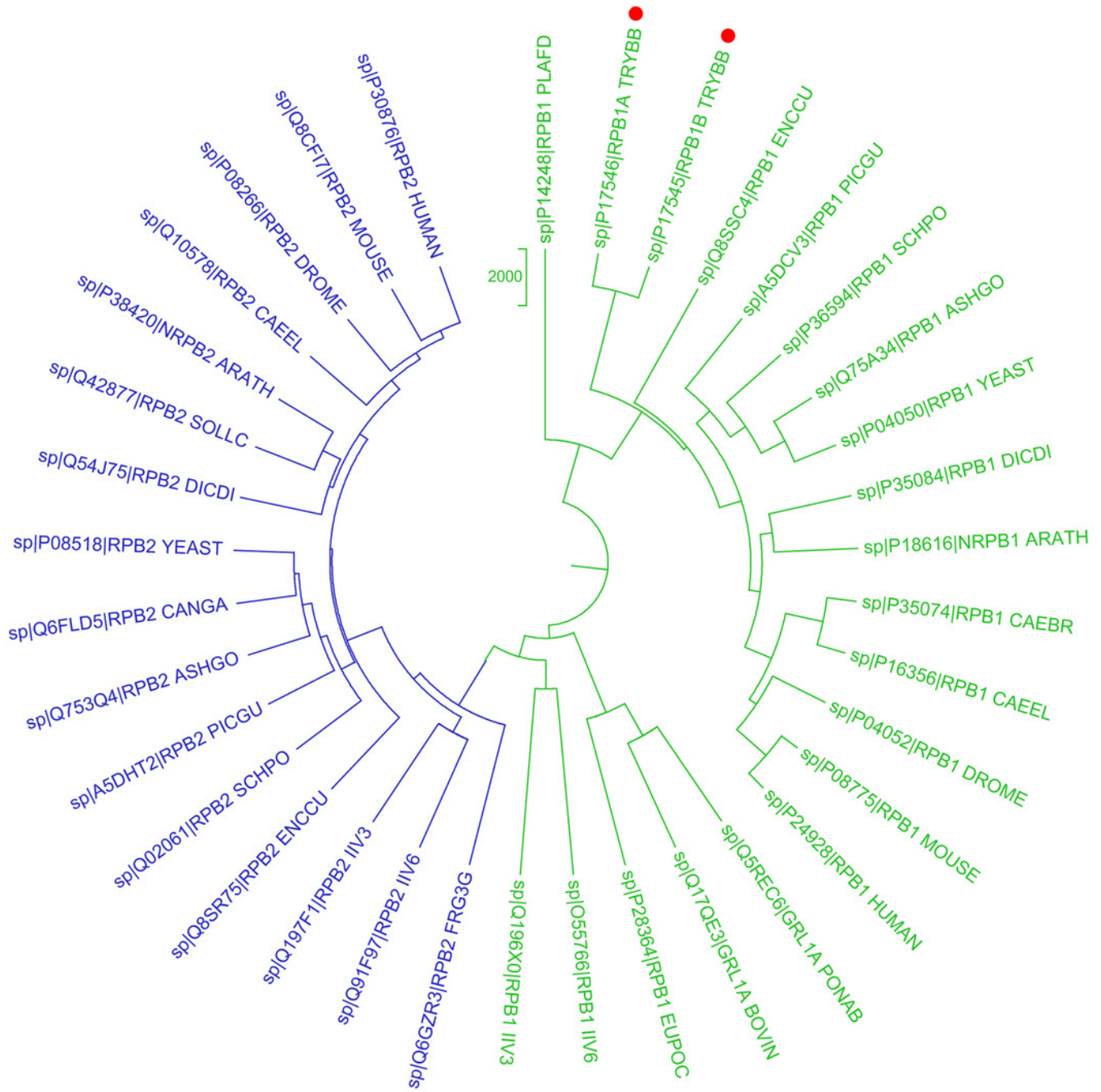
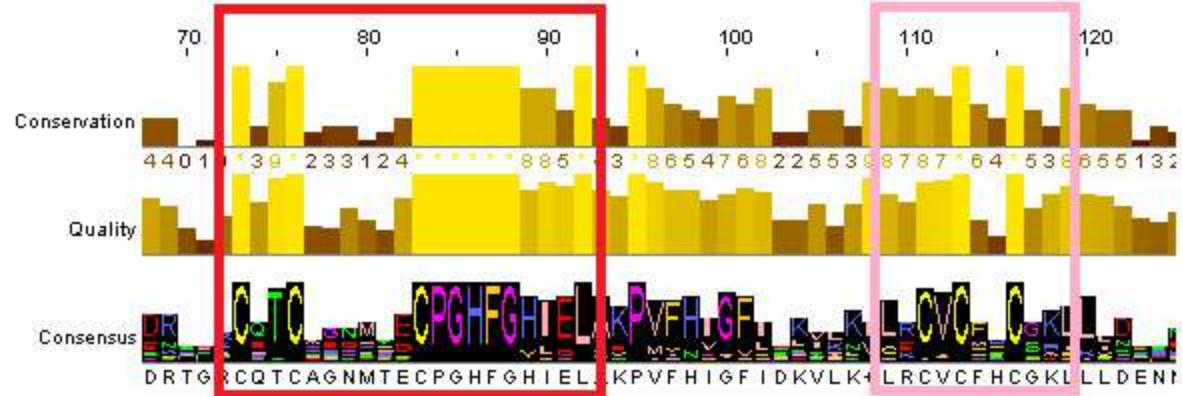


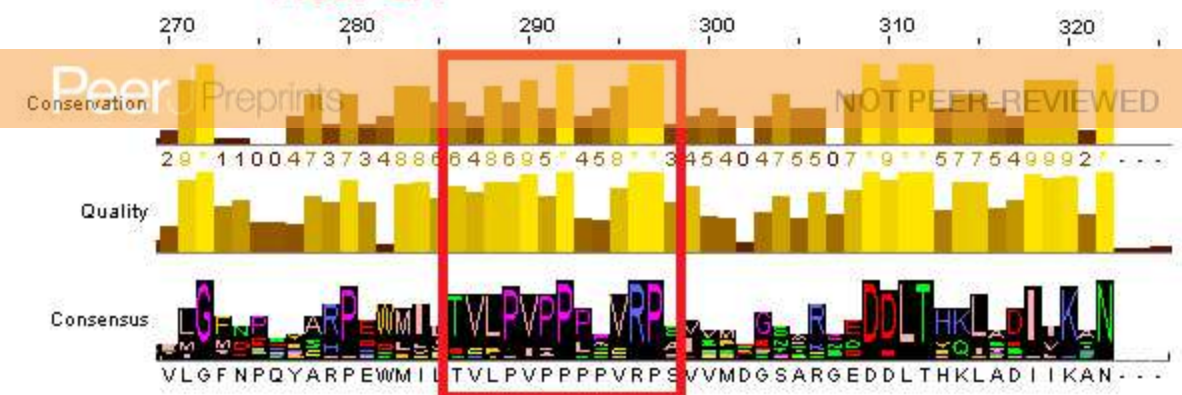
Figure 2 (on next page)**Representative conserved motifs for the DdRpII subunit RPB1.**

The nine suggested conserved motifs were extracted based on the multiple sequence alignment of the 18 protein sequences were classified and clearly separated in the DdRpII subunit RPB1 monophyletic sub-tree. The conserved motifs were identified through the consensus sequence and logo graph where generated using Jalview software.

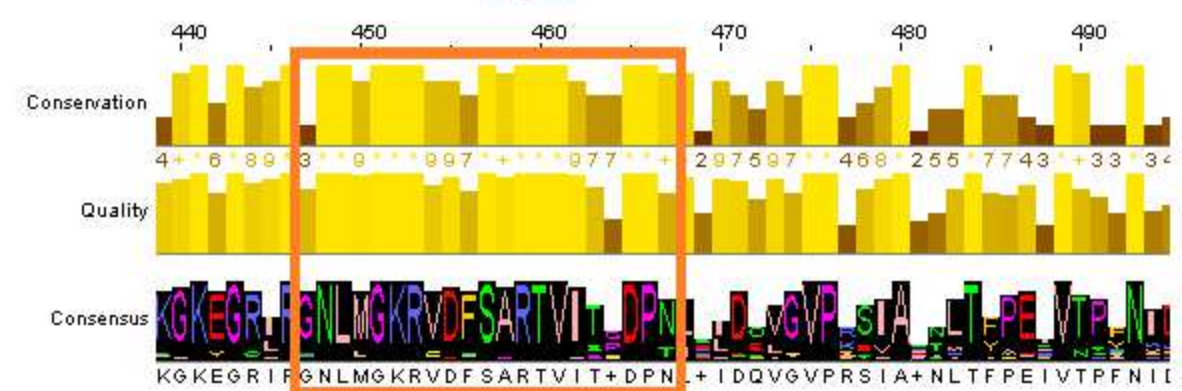


Motif 1A

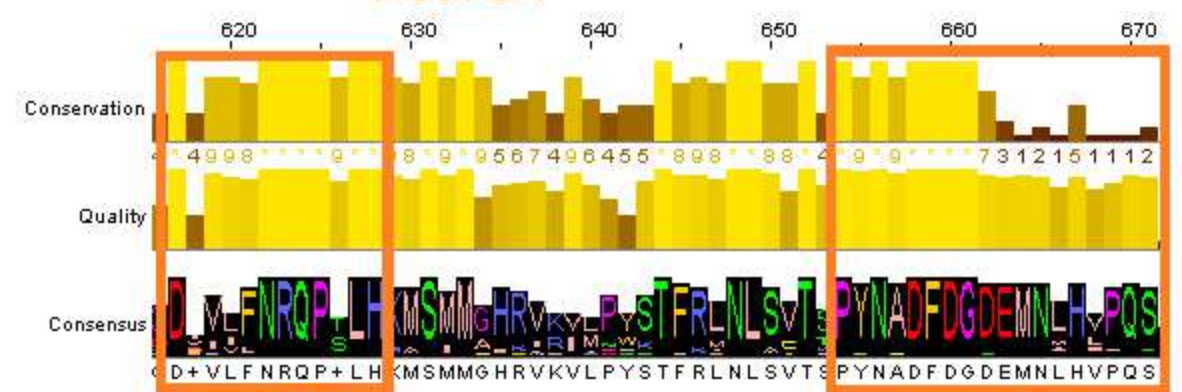
Motif 2



Motif 1B

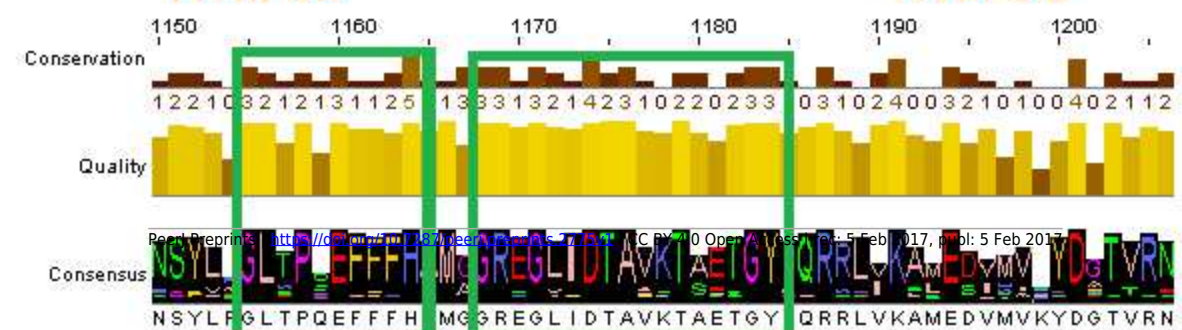


Motif 3A



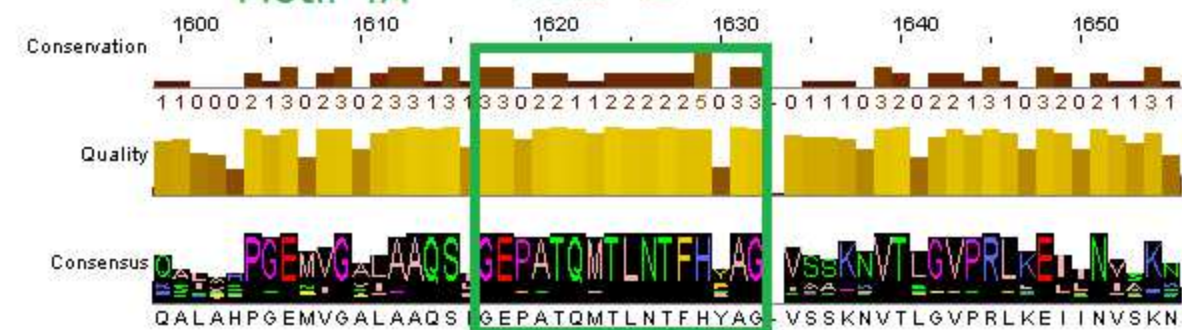
Motif 3B

Motif 3C



Motif 4A

Motif 4B



Motif 4C

Figure 3

Sequence alignment between the *Trypanosoma brucei brucei* DdRpII RPB1 and the corresponding sequence of the crystal structure of the *Schizosaccharomyces pombe* DdRpII RPB1.

(A) Alignment of DdRpII RPB1 from *Trypanosoma brucei* DdRpII RPB1 (Labeled as “TB”) with *Schizosaccharomyces pombe* DdRpII RPB1 (Labeled as “SB”) was initially carried out with BLASTp and then manually adjusted. The nine suggested conserved motifs (Motifs 1a, 1b, 2, 3a, 3b, 3c, 4a, 4b, 4c) based on figure 2, domains and domain-like regions of *Trypanosoma brucei* DdRpII RPB1 represented in different colours. The amino acid residue numbers at the domain boundaries are indicated. Important structural elements and prominent regions involved in subunit interactions are also noted. Residues involved in the Zn and Mg coordination are highlighted in blue. **(B)** Domains and domain-like regions of the DdRpII subunit Rpb1. The amino acid residue numbers at the domain boundaries are indicated.

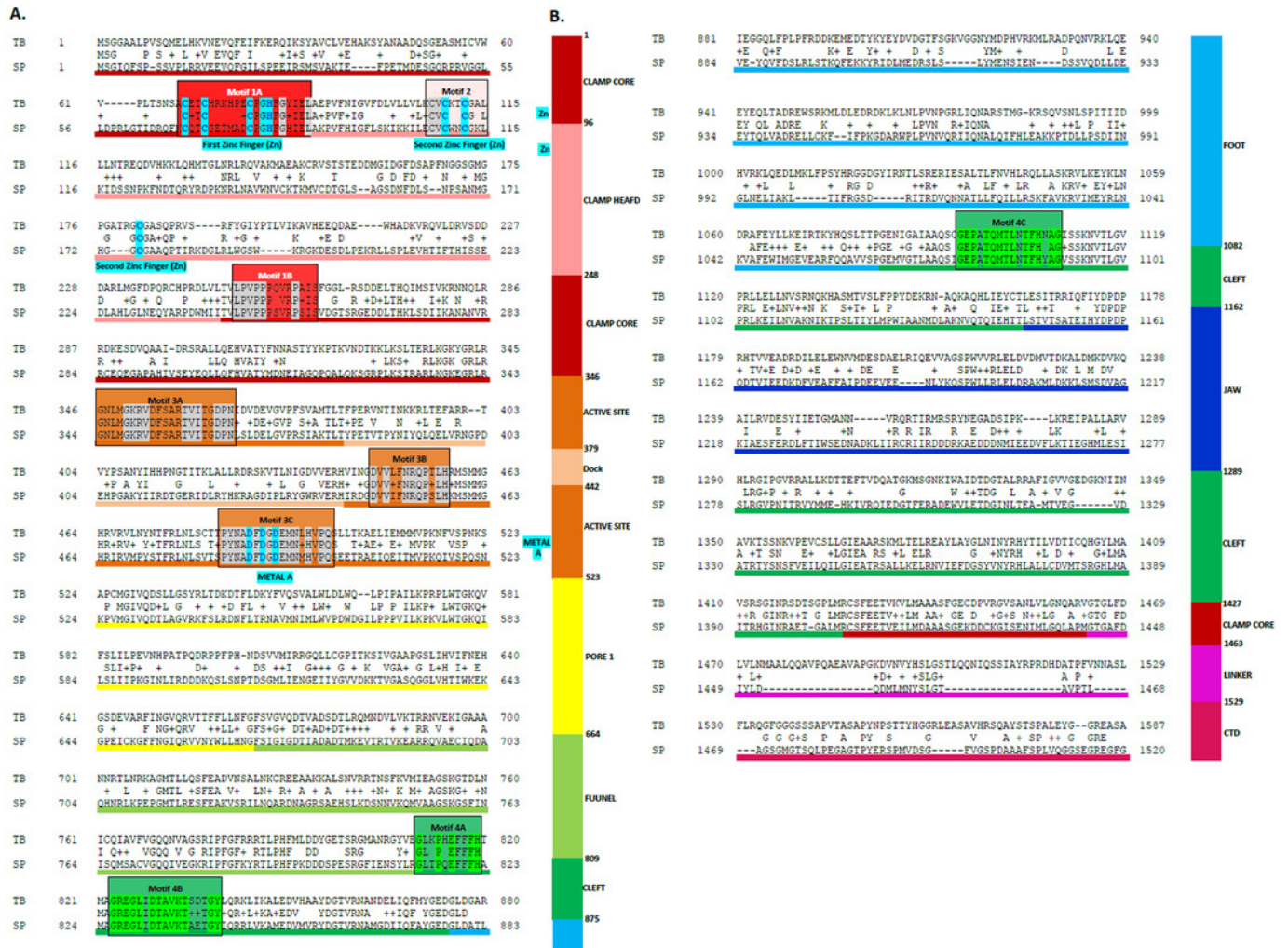
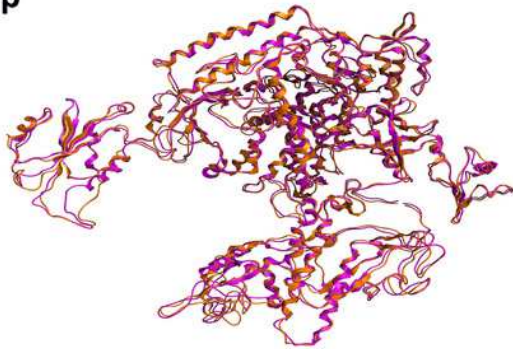


Figure 4

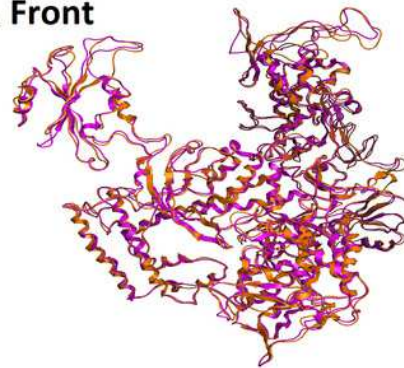
Model of the *Trypanosoma brucei brucei* DdRPII RPB1 .

(A and B) Ribbon representation of the produced *Trypanosoma brucei brucei* DdRPII RPB1 model (colored Orange) superposed with the corresponding *Schizosaccharomyces pombe* DdRPII RPB1 (in purple). **(C and D)** The nine suggested conserved motifs and the domains and domain-like regions of the *Trypanosoma brucei brucei* DdRPII RPB1 . The motifs and RPB1 domains have been color-coded according to the Figures 2 and 3, and are shown in CPK format (Usual space filling). **(E and F)** Electrostatic surface potential for the *Trypanosoma brucei brucei* DdRPII RPB1 . Represented with blue is the area of negative charge. Red is the area of positive charge and white is the un-charged region.

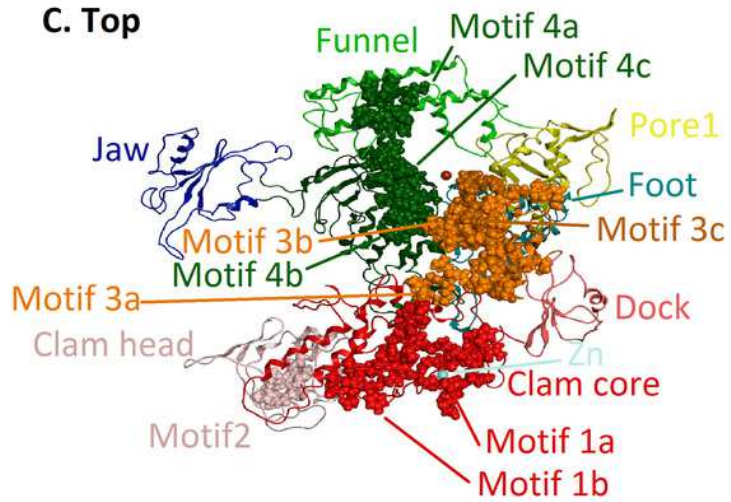
A. Top



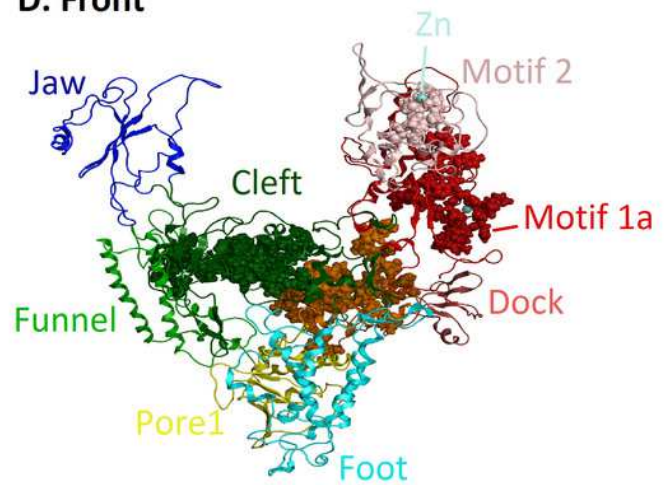
B. Front



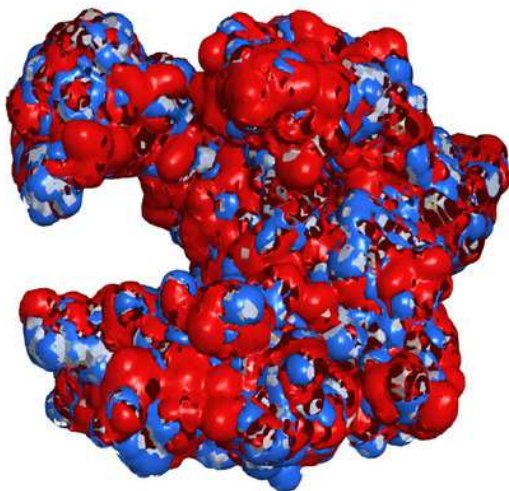
C. Top



D. Front



E. Top



F. Front

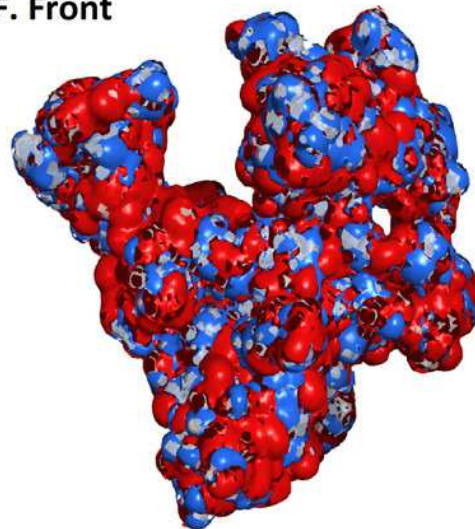


Figure 5

Structural superposition of the TBB DdRPII RPB1 models A and B

(A and B) Ribbon representation of the produced *Trypanosoma brucei brucei* DdRPII RPB1 model A (colored Orange) and model B (colored Blue) superposed with the corresponding *Schizosaccharomyces pombe* DdRplI RPB1 (in Purple) and *Bos taurus* DdRplI RPB1 (in Grey). The four 3D structures are highly conserved in their active sites with few differences in the outer layer with overall RMSD 2.775 Å. **(C)** Ribbon representation of the produced *Trypanosoma brucei brucei* DdRPII RPB1 model A (colored Orange) superposed with the corresponding *Schizosaccharomyces pombe* DdRplI RPB1 (in purple). (RMSD = 1.242 Å). **(D)** Ribbon representation of the produced *Trypanosoma brucei brucei* DdRPII RPB1 model B (colored Blue) superposed with the *Bos taurus* DdRplI RPB1 (in Grey) respectively. (RMSD = 2.757 Å).

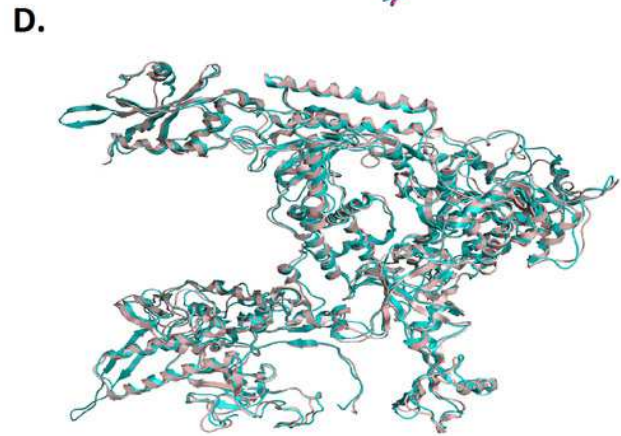
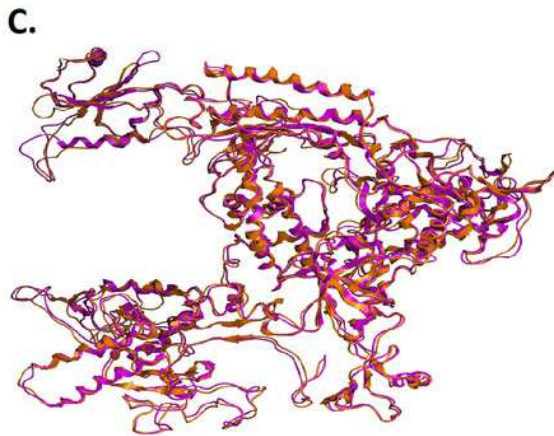
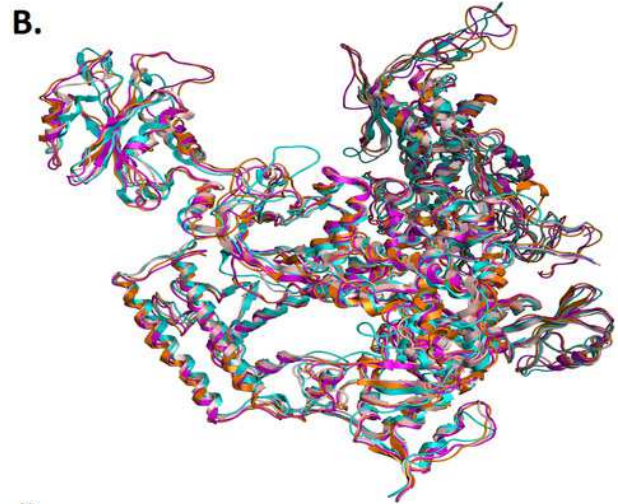
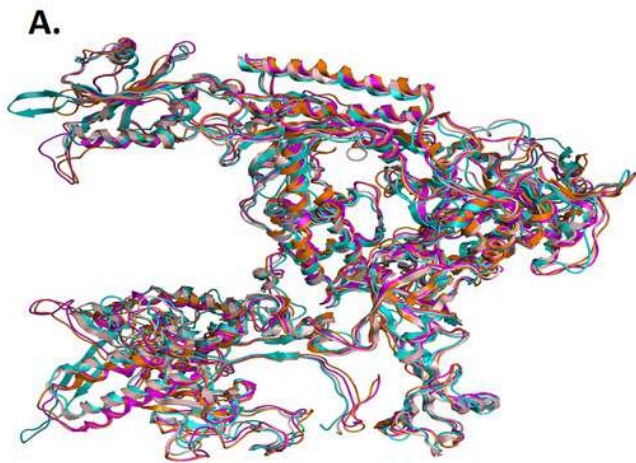


Figure 6

Zinc-finger formations in the *Trypanosoma brucei brucei* DdRPII RPB1 model.

Ribbon representation of the produced *Trypanosoma brucei brucei* DdRPII RPB1 model. In the produced model were highlighted 3 main zing-finger domain formations (colored grey) were contained in the clam core, clam head and active site region. Domains and domain-like regions of the *Trypanosoma brucei brucei* DdRPII RPB1 have been color-coded according to conventions of Figures 3.

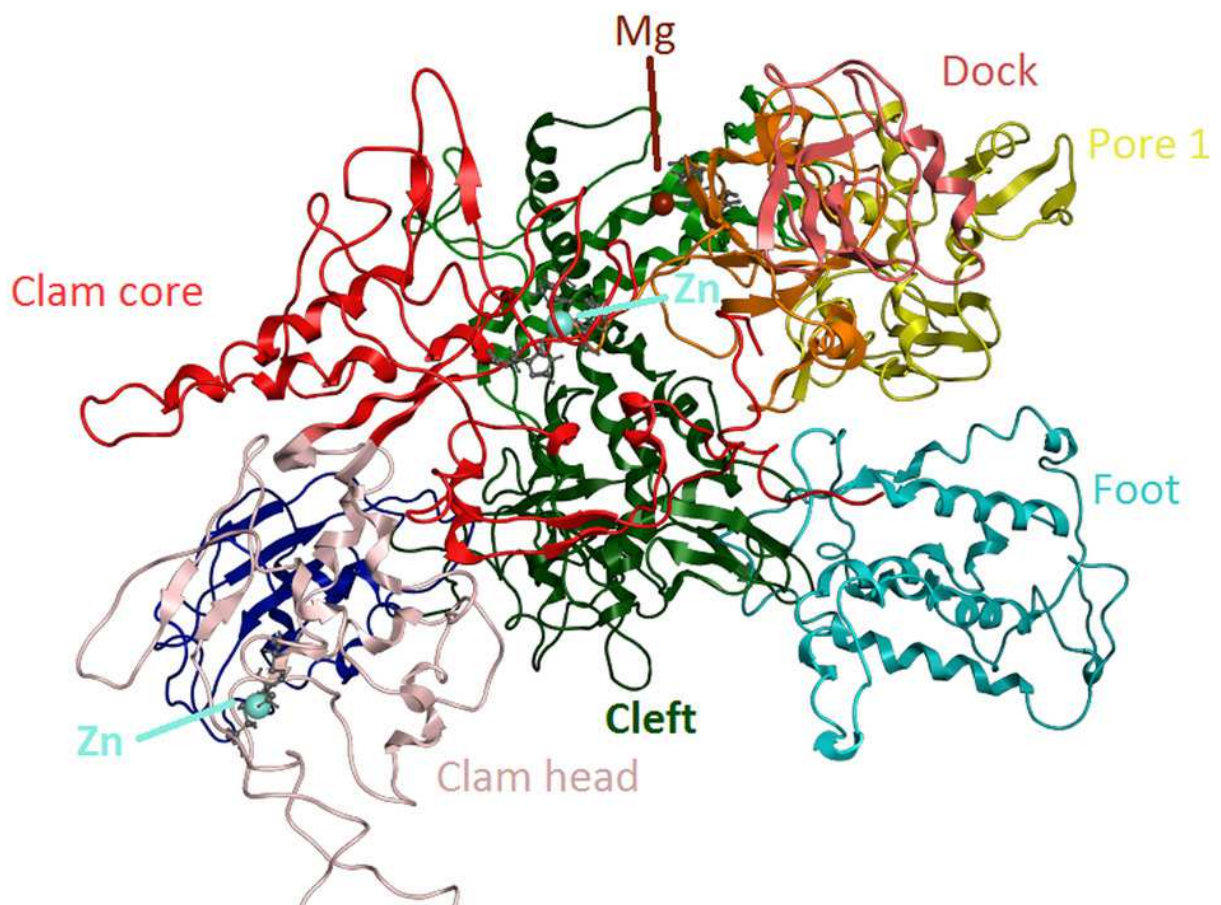


Figure 7

Molecular dynamics simulation charts for the *Trypanosoma brucei brucei* DdRpII RPB1 models.

(A) The root mean square deviation (RMSD) of the model A during the time. **(B)** The root mean square fluctuation (RMSF) of the model A during the time. **(C)** The root mean square deviation (RMSD) of the model B during the time. **(D)** The root mean square fluctuation (RMSF) of the model B during the time.

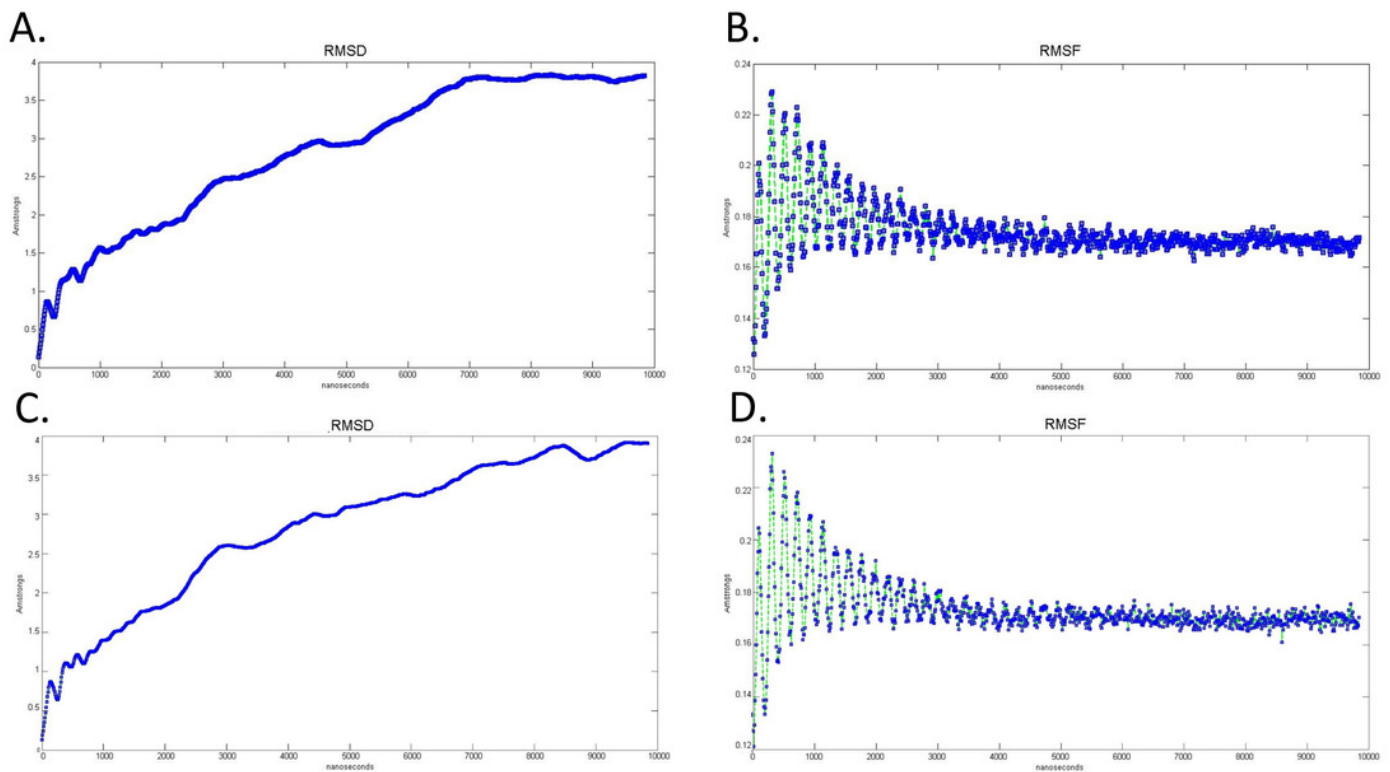


Figure 8

The 3D pharmacophore model for the *Trypanosomabrucei brucei* DdRP11 RPB1 model.

In total 5 distinct pharmacophoric features were identified. An aromatic region (colored orange), an electron donating region (colored green), two electron accepting regions (colored red) and a sulphur specific S-S interacting region (colored yellow).

