# Structural basis of SNAPc-dependent snRNA transcription initiation by RNA polymerase II

**Srinivasan Rengachari[1], Sandra Schilbach[1], Thangavelu Kaliyappan[2], Jerome Gouge[2], Kristina Zumer[1], Juliane Schwarz[3,4], Henning Urlaub[3,4], Christian Dienemann[1], Alessandro Vannini[2,5*] and Patrick Cramer[1*]**

*[1] Department of Molecular Biology, Max Planck Institute for Multidisciplinary Sciences, Am Fassberg 11, 37077 Göttingen, Germany.*

*[2] Division of Structural Biology, The Institute of Cancer Research, London, SW7 3RP, UK*

*[3] Max Planck Institute for Multidisciplinary Sciences, Bioanalytical Mass Spectrometry, Göttingen, Germany*

*[4] University Medical Center Göttingen, Institute of Clinical Chemistry, Bioanalytics Group, Göttingen, Germany*

*[5] Human Technopole, 20157, Milan, Italy*

*Correspondence: alessandro.vannini@fht.org; patrick.cramer@mpinat.mpg.de

**RNA Polymerase II (Pol II) carries out transcription of both protein-coding and non-coding genes. Whereas Pol II initiation at protein-coding genes has been studied in detail, Pol II initiation at non-coding genes such as small nuclear RNA (snRNA) genes is not understood at the structural level. Here we study Pol II initiation at snRNA gene promoters and show that the snRNA-activating protein complex (SNAPc) enables DNA opening and transcription initiation independent of TFIIE and TFIIH *in vitro*. We then resolve cryo-EM structures of the SNAPc-containing Pol II preinitiation complex (PIC) assembled on U1 and U5 snRNA promoters. The core of SNAPc binds two turns of DNA and recognizes the snRNA promoter-specific proximal sequence element (PSE) located upstream of the TATA box-binding protein TBP. Two extensions of SNAPc called wing-1 and wing-2 bind TFIIA and TFIIB, respectively, explaining how SNAPc directs Pol II to snRNA promoters. Comparison of structures of closed and open promoter complexes elucidates TFIIH-independent DNA opening. These results provide the structural basis of Pol II initiation at non-coding RNA gene promoters.**

Transcription by RNA polymerase II (Pol II) has been structurally well studied for protein-coding genes that produce messenger RNA (mRNA)[1-4]. Pol II however also carries out transcription of non-coding small nuclear RNAs (snRNAs) that are an integral part of the pre-mRNA splicing machinery[5]. Pol II transcribes four of the five snRNAs, namely U1, U2, U4 and U5 snRNAs, whereas Pol III transcribes U6 snRNA[6]. In contrast to the Pol III-dependent snRNA promoter, Pol II-dependent snRNA promoters lack a TATA box motif[7]. To produce snRNAs, Pol II uses many of its accessory factors that are used for mRNA synthesis, but additionally requires specific factors for transcription initiation and elongation[8].

Transcription initiation of snRNA genes relies on a specific factor called snRNA-activating protein complex (SNAPc). SNAPc binds a DNA motif in the upstream region of snRNA promoters, the so-called proximal sequence element, or PSE[9]. Human SNAPc contains

44   five subunits – SNAPC1, SNAPC2, SNAPC3, SNAPC4 and SNAPC5. The subunits SNAPC1,
45   SNAPC3 and SNAPC4 form the core of SNAPc[10], of which SNAPC3 and SNAPC4 posess
46   DNA-binding function[11,12]. The core subunits of SNAPc are conserved and have been
47   characterized in Drosophila melanogaster and Trypanosoma brucei, where they are sufficient
48   for activating snRNA transcription[13,14]. SNAPC2 and SNAPC5 however contribute to the
49   stability and activity of SNAPc[10,15,16].

50   The initiation of SNAPc-regulated Pol II snRNA transcription was reported to rely on
51   the general transcription factors (GTFs) TBP, TFIIA, TFIIB, TFIIE and TFIIF[17,18] . The role of
52   TFIIH in Pol II snRNA transcription remains unclear[17], although TFIIH is known to be required
53   for DNA opening at promoters of protein-coding genes[19]. The structure of SNAPc and its
54   molecular interactions with the Pol II pre-initiation complex (PIC) are also unknown. As a
55   consequence, the structural basis and the mechanism of snRNA transcription initiation remains
56   to be uncovered. Here we report structures of SNAPc-containing Pol II PICs bound to U1 and
57   U5 snRNA promoters. Our results show how SNAPc is structured, how it recognizes the PSE,
58   and how it positions a core Pol II PIC on snRNA promoters for DNA opening and transcription
59   initiation. More generally, this work adds to our understanding of the evolution of the three
60   eukaryotic transcription systems.

61

62   **RESULTS**
63   **Preparation of functional SNAPc**
64   We prepared two variants of recombinant human SNAPc, namely SNAPc-FL containing all
65   full-length subunits and SNAPc-core[10], containing SNAPC1, SNAPC3, SNAPC4 (residues 1-
66   516) and SNAPC5 (**Figure 1a**) (Methods). Both purified SNAPc variants were able to bind U1
67   and U5 snRNA promoter DNA (RNU1 and RNU5), both in the absence and in the presence of
68   TBP and TFIIB in an electrophoretic mobility shift assay (EMSA) (**Figure 1b**). EMSA also
69   showed that the SNAPc variants could facilitate binding of TBP to snRNA promoters that lack
70   a TATA box (**Figure 1b, 2a**), consistent with previous studies[20].

71   To test whether recombinant SNAPc could mediate Pol II transcription initiation from
72   snRNA gene promoters, we used an *in vitro* transcription assay. The assay showed that Pol II
73   could initiate transcription from a U1 promoter in the presence of TBP, TFIIA, TFIIB and TFIIF
74   and was stimulated ~4.5 fold and ~2.5 fold by the addition of SNAPc-FL or SNAPc-core,
75   respectively (**Figure 1d, e**). Addition of TFIIE reduced this increase in transcription activity to
76   ~1.5 fold and ~1.2 fold for SNAPc-FL and SNAPc-core, respectively (**Figure 1e**), suggesting
77   that TFIIE is not required for SNAPc-dependent transcription initiation and rather inhibitory in
78   our biochemical system. Further addition of TFIIH did override the stimulatory effect of
79   SNAPc and led to formation of non-specific transcripts at multiple sites (**Figure 1d**). In
80   conclusion, our recombinant SNAPc variants stimulate Pol II transcription initiation from
81   snRNA gene promoters in the absence of TFIIE and TFIIH *in vitro*.

82

83   **Cryo-EM analysis of SNAPc-containing PICs**
84   Based on these observations we reconstituted a functional SNAPc-containing Pol II PIC on a
85   U1 promoter DNA (Methods). We incubated SNAPc-core and S. scrofa Pol II (99.9% identical
86   to human Pol II) with human TBP, TFIIA, TFIIB, TFIIE and TFIIF, and subjected the resulting
87   complex to sucrose-gradient ultracentrifugation. Peak fractions contained apparent
88   stochiometric amounts of Pol II, SNAPc-core and the general transcription factors, indicating

89  formation of a stable 24-subunit SNAPc-containing PIC (**Figure 1c**). The sample was
90  crosslinked[21] and subjected to cryo-EM analysis (Methods). Initial trials showed that the PIC
91  containing the SNAPc-FL variant was less stable (**Figure 1c**), whereas the PIC containing
92  SNAPc-core was stable and suited for cryo-EM analysis, leading to a high-resolution single
93  particle dataset (**Extended Data Table 1**).

94       Reconstructions from 3D classification of this dataset showed two distinct particle
95  classes of the SNAPc-containing Pol II PIC (**Extended Data Figure 1**). Further 3D
96  classification and refinement identified these two states as the closed promoter complex (CC)
97  and the open promoter complex (OC) states of the PIC. We obtained structures of the CC and
98  OC states at an overall resolution of 3.4 Å and 3.0 Å, respectively (**Extended Data Figure 1,**
99  **3**). None of our maps revealed density for TFIIE, consistent with our in vitro transcription
100 assays that showed TFIIE was not required for initiation (**Figure 1d, e**). Densities for SNAPc
101 and upstream DNA containing the PSE were improved by focussed 3D classification and
102 masked refinements. The local resolution for this region was 3.5 Å for the OC state (**Extended**
103 **Data Figure 1, 3**).

104      In an effort to obtain a high-resolution structure of SNAPc, we additionally
105 reconstituted a SNAPc-containing Pol II PIC on a DNA that was based on the U5 promoter
106 sequence (Methods). The resulting cryo-EM dataset enabled refinement of the SNAPc-
107 containing PIC in the CC state at an overall resolution of 3.0 Å and with the local map of the
108 upstream region extending to 3.2 Å (**Extended Data Figure 2, 3**). The local map enabled
109 building of an atomic model for SNAPc and PSE-containing upstream DNA (Methods). We
110 then combined the resulting model with the known high-resolution structures of mammalian
111 Pol II PIC in CC and OC states[22]. After manual adjustment, refined structures of the SNAPc-
112 containing PIC in the CC and OC states containing the U1 and U5 promoters showed good
113 stereochemistry resulting in a total of three structures (**Extended Data Table 1**).

114

115 **Overall structure of SNAPc-containing PIC**
116 The overall structure of the SNAPc-containing Pol II PIC shows that SNAPc binds the promoter
117 DNA upstream of TBP (**Figure 2**). SNAPc recognizes the PSE motif and interacts with TFIIA
118 and TFIIB. Despite these multiple interactions, the presence of SNAPc does not alter the
119 canonical core PIC structure in any substantial way[22]. TBP binds to the AGGCTG sequence at
120 register –30 to –25 bp (**Figure 2a**) of the TATA-less U1 promoter and bends the DNA by 90°
121 similar to what is observed in a TBP-TATA DNA complex[23,24] (**Extended Data Figure 4a**). In
122 the following, we will first describe the SNAPc structure and SNAPc-DNA interactions based
123 on the U5-containing CC structure that is resolved at the highest resolution. We will then
124 describe promoter opening based on the CC and OC structures of the U1-containing PIC.

125

126 **SNAPc structure contains two protruding wings**
127 The high-resolution structure of the SNAPc core bound to the U5 promoter shows how the
128 subunits SNAPC1, SNAPC3 and SNAPC4 fold and interact (**Figure 3**). SNAPC1 possesses an
129 N-terminal VHS/ENTH-like domain[25] that forms a mainly helical structure (**Extended Data**
130 **Figure 4b**). SNAPC3 is saddle-shaped with a central 'ubiquitin-like domain' (ULD) and
131 additional α-helices and β-strands (**Extended Data Figure 4c**). Consistent with biochemical
132 studies[26], SNAPC3 contains two zinc fingers (ZF-1 and ZF-2). ZF-1 is a C2H2 type zinc finger
133 with residues Cys221, His313, Cys317 and His319 coordinating a $Zn^{2+}$ ion (**Extended Data**

134   **Figure 5f**). ZF-2 is a C4 type zinc finger with residues Cys354, Cys357, Cys380, Cys383
135   coordinating another $Zn^{2+}$ ion (**Extended Data Figure 5g**). SNAPC4 contains four complete
136   repeats (R1-R4) and a half repeat (Rh) of the Myb domain[12], of which we observe Rh, R1 and
137   R2 (residues 274-398) (**Figure 3b**, **Extended Data Figure 4d**). R1 and R2 contain three helices
138   forming canonical helix-turn-helix folds. The SNAPc core is stabilized by intricate interactions
139   of SNAPC3 with both SNAPC1 and SNAPC4. The N-terminal region of SNAPC3 interacts
140   mainly with SNAPC1, burying a surface area of ~1640 $\text{Å}^2$. The C-terminal region of SNAPC3
141   binds SNAPC4 and buries ~3010 $\text{Å}^2$. A total of four subunit interfaces are formed based on
142   hydrophobic interactions, salt bridges and polar contacts (**Figure 3c-f, Extended Data Figure**
143   **7**).

144   SNAPc also contains two protrusions that we refer to as 'wing-1' and 'wing-2'. The
145   wing-1 of SNAPc consists of a pair of helices that precede the Rh region of SNAPC4 (residues
146   184-256). The wing-2 of SNAPc is a four-helix bundle that is formed by two helices of
147   SNAPC1 (residues 162-234) and one helix each of SNAPC4 (residues 81-125) and SNAPC5
148   (residues 1-51) (**Extended Data Figure 4e**). Although the resolution in wing-2 is limited due
149   to mobility, AlphaFold2 prediction[27] and prior biochemical studies[16] led to a reliable model of
150   wing-2 that we confirmed by crosslinking mass-spectrometry (**Extended Data Figure 5k, 6**).
151   In conclusion, these efforts provided the structure of SNAPc, which contains a three-subunit
152   core and two protruding wings extending from the core.

153

154   **SNAPc core recognizes the snRNA promoter**
155   Our U5-containing CC structure also reveals details of how SNAPc binds the PSE motif in
156   promoter DNA (**Figure 4a**). The SNAPc core binds to the PSE motif through its subunits
157   SNAPC3 and SNAPC4 (**Extended Data Figure 8a**), consistent with biochemical data[11,28].
158   SNAPc contacts promoter DNA 8 bp upstream of the proximal edge of the TBP-binding site
159   (**Figure 4b, c**). The register of the modelled snRNA promoter is defined by the nucleotide on
160   the non-template (NT) strand at the upstream edge of TBP binding site starting at –30,
161   ascending in the 5' to 3' direction. SNAPC3 and SNAPC4 both bind this region through
162   contacts with the DNA backbone and bases on both strands of the PSE (**Extended Data Figure**
163   **8a**). DNA binding occurs both to the major and minor grooves. SNAPC3 inserts its helix α8
164   into the major groove and forms multiple contacts with DNA. K199 forms salt bridges with the
165   backbone phosphates of nucleotide G9 on the template strand. The residue K194 of the same
166   helix forms ionic interactions with the O6 atom of the nucleotide bases G –42 and G –43 of the
167   NT strand. Further downstream, H198 establishes hydrophobic contacts with the nucleotide
168   base T –45 on the template strand (**Figure 4b**).

169   Since most of these protein-DNA contacts are to the DNA backbone, the question arises
170   how SNAPc can recognize the PSE. Our structure suggests that recognition is at least partially
171   achieved by indirect readout. In particular, the DNA major groove is locally distorted at the
172   PSE and differs from canonical B-DNA at registers –51 to –41 (**Extended Data Figure 8b**).
173   At the position where SNAP3 helix α8 is inserted into the major groove, the duplex geometry
174   resembles A-form DNA[29] (**Extended Data Figure 8c**). This deviation is also reflected in the
175   minor grooves upstream and downstream of this site (**Extended Data Figure 8a, d**).

176   SNAPc also binds the minor groove of DNA with subunits SNAPC3 and SNAPC4.
177   Q152 of SNAPC3 a forms hydrogen bond with the nucleotide base T –48 of NT strand while
178   SNAPC4 residue Y372 interacts hydrophobically with the C3 atom of backbone sugar of the

4

179    nucleotide base A –50 of the template strand. Arginine residues R148 and R151 of SNAPC3
180    and R373 of SNAPC4 form salt bridges with the DNA backbone (**Figure 4c**). Our structure
181    also shows that the SNAPC4 Myb repeat R2 binds DNA via its helix α15 that contacts the
182    anterior major groove, and early biochemical studies indicated that the Myb repeats R3 and R4
183    are involved in DNA binding[12]. I388 establishes hydrophobic interactions with the nucleotide
184    base A –50 and the C5 atom of nucleotide C –51 on the template strand. The neighbouring
185    Y389 residue forms hydrogen bonds with the N7 atom of A –55  and hydrophobic interaction
186    with T -54 of the NT strand (**Figure 4c**). The residues K347, R373 and R390 of SNAPC4
187    interact with the DNA backbone. Although biochemical studies had identified SNAPC3 and
188    SNAPC4 as poor DNA binders when investigated in isolation[10,11], our results suggest that
189    formation of the SNAPc complex with its intricate interactions between these two subunits
190    enables tight binding of the PSE which explains how SNAPc recognizes the snRNA promoter.
191
192    **The wings of SNAPc bind TFIIA and TFIIB**
193    SNAPc also interacts with TFIIA and TFIIB that flank TBP in the PIC (**Figures 5, Extended**
194    **Data Figure 8a**). Whereas wing-1 of SNAPc binds to TFIIA, wing-2 binds TFIIB (**Extended**
195    **Data Figure 8a**). SNAPc interaction with TFIIA and TFIIB involves three interfaces that we
196    call A, B and C. In interface A, the wing-1 of SNAPC4 (helices α4, α5) slides under the four-
197    helix bundle of TFIIA like a wedge, stabilizing the flexible bundle region (**Figure 5a**). SNAPC4
198    additionally interacts with the β-barrel of TFIIA to form interface B (**Figure 5a**). Interfaces A
199    and B are formed by a combination of hydrophobic interactions, salt bridges and polar contacts.
200    Incidentally, the TFIIA bundle has also been shown to interact with TAF4 and TAF12 in lobe
201    B of the multisubunit TFIID complex that, like SNAPc, is important for promoter recognition[30].
202    Interface C is formed between wing-2 and the C-terminal cyclin fold of the TFIIB core (**Figure**
203    **5b**). The wing-2 helices from SNAPC1 and SNAPC5 form contacts with the terminal α-helix
204    of the TFIIB core. Interface C stabilizes the TFIIB core, which was suggested to play a key role
205    in the activation of snRNA transcription initiation[7]. Together, SNAPc wing-1 and wing-2 bind
206    TFIIA and TFIIB, respectively, to position the core PIC with respect to SNAPc and the PSE
207    promoter element.
208
209    **Promoter DNA opening**
210    Comparison of our CC and OC structures bound to the U1 promoter provides insights into the
211    mechanism of TFIIE- and TFIIH-independent DNA opening (**Figure 6**). Overall, closed and
212    open U1 promoter DNA follow trajectories within the Pol II cleft that are comparable to those
213    observed for protein-coding promoter DNA in PIC structure[22]. Also, as observed in PIC
214    structures lacking SNAPc[2,22], the OC state is associated with a closed Pol II clamp and an
215    ordered B-reader and B-linker elements in TFIIB (**Figure 6b**). However, DNA opening can
216    also be achieved spontaneously in the absence of TFIIE and TFIIH at some protein-coding
217    genes in yeast[31], and such spontaneous opening depends on the DNA duplex stability around
218    the transcription start site (TSS)[32]. Studies in yeast Pol II have further shown that an AT-rich
219    sequence increases the propensity of spontaneous promoter opening during transcription
220    intiation[31]. Similarly, we find that promoter sequences of snRNA-encoding genes are AT-rich
221    in the initially melted region (IMR) spanning positions –8 to +2 around the TSS (position +1)
222    (**Extended Data Figure 8e**). We propose that the AT-rich nature of the IMR enables

223    spontaneous DNA opening of the U1 promoter upon PIC binding. In summary, these results
224    suggest that DNA opening of snRNA gene promoters and the spontaneously melted protein-
225    coding genes rely on easily melting regions around the TSS and use similar mechanisms.
226
227    **Definition of the transcription start site**
228    We observe 19 nucleotides of the DNA template strand spanning from the TBP-binding site to
229    the upstream edge of the DNA bubble (at position -12). The templating nucleotide in open
230    promoter DNA reaches the active site of Pol II ~30 nucleotides downstream of the upstream
231    edge of the TBP-binding site (**Figures 6a, b**). The DNA strands forming the open DNA bubble
232    are mobile, leading to a weakly resolved map. Subsequently, 12 nucleotides further
233    downstream, we observe T +1 of the template strand immediately downstream of the catalytic
234    $Mg^{2+}$ ion at the active site. This posits residue G –1 as the template for RNA synthesis. The CA
235    dinucleotide is the signature of the Initiator sequence (Inr)[33] and is located at register –1 and +1
236    of the non-template strand. This observation suggests that the TSS position is defined by a fixed
237    distance from the site of TBP binding, as is known for protein-coding human genes that have
238    their TSS within a window of 28-33 bp downstream of the TATA box[34]. Since we also observe
239    a fixed position of SNAPc with respect to TBP, the TSS is apparently set by a fixed distance
240    from the PSE in snRNA promoters.
241         These observations suggest that Pol II transcription would initiate from a TSS that is
242    rather precise *in vivo*. To investigate this, we identified the main TSSs and determined their
243    'TSS precision scores' from a reanalysis of 5'- capped RNA sequencing data[35] for both mRNA
244    and snRNA encoding genes with a constitutive first or a single exon (Methods). A maximum
245    precision score of 1 means that all transcripts initiate at the main TSS (±2 bp). Indeed we find
246    that Pol II snRNA transcription generally initiates in this narrow, 5-bp window with a high
247    median precision score of 0.86, as exemplified by the *RNVU1-15* promoter (**Figure 6c**). In
248    contrast, Pol II initiates transcription less precisely at TATA-less mRNA promoters, as shown
249    by a median precision score of 0.36, as exemplified by the *HAT1* promoter. Pol II also initiates
250    mRNA transcription more precisely when promoter DNA contains a TATA box motif, with a
251    median precision score of 0.71, as exemplified by the *TUBB4B* promoter (**Figure 6c**). These
252    large differences in TSS precision are also observed in genome browser views of representative
253    promoters (**Figure 6d**). The observed high TSS precision of snRNA promoters is consistent
254    with our model that SNAPc defines TSS position. In summary, SNAPc binding to the PSE
255    likely serves as a ruler for positioning of TBP at TATA-less snRNA promoters, leading to
256    initiation at a defined distance downstream of the PSE.
257
258    **DISCUSSION**
259    Here we report structures of SNAPc-containing Pol II PICs on two different snRNA gene
260    promoters and in two different states, the CC and OC states. Together with biochemical results
261    and published literature, our structures suggest the mechanism of SNAPc-mediated snRNA
262    transcription initiation by Pol II (**Figure 7**). SNAPc uses its conserved core to recognize the
263    PSE motif in snRNA promoters, whereas its two wings position TFIIA and TFIIB. Since TFIIA
264    and TFIIB form a rigid complex with TBP, SNAPc can indirectly position TBP at a defined
265    location on snRNA promoters despite the absence of a consensus TATA box motif. This is
266    consistent with the evidence that TFIIB-TBP complexes can be effectively recruited to snRNA
267    promoters exclusively as part of a ternary TFIIA-TFIIB-TBP complex[18]. Positioning of the

268    TFIIA-TFIIB-TBP complex on promoter DNA in turn recruits the Pol II-TFIIF complex to the
269    IMR of the promoter. The low DNA duplex stability at the IMR enables spontaneous DNA
270    opening and occurs with the use of binding energy independent of TFIIE and TFIIH. The
271    emerging DNA template strand then binds in the Pol II active center cleft and RNA chain
272    synthesis is initiated at an Inr dinucleotide CA[33], thereby setting the TSS at a defined distance
273    from the PSE.

274         Comparison of our results with published data also provides insights into the evolution
275    of the three different eukaryotic transcription systems. A distinguishing feature of transcription
276    initiation by Pol II, with respect to Pol I and Pol III, is that the latter two machineries can open
277    promoter DNA spontaneously[36-40], whereas the Pol II machinery generally requires the help of
278    an ATP-dependent translocase subunit in TFIIH and its accessory factor TFIIE[22,41]. However,
279    we show here that on snRNA promoters, mammalian Pol II, together with the factors that form
280    the core PIC, can open DNA spontaneously without the help of TFIIE and TFIIH. Such
281    spontaneous DNA opening has also been observed for yeast Pol II at a subset of promoters[31]
282    and also in the related archaeal transcription system[42]. Whereas spontaneous DNA opening
283    occurs in the upstream-to-downstream direction, TFIIH-assisted DNA opening occurs in the
284    downstream-to-upstream direction[22,41]. Our work thus provides evidence that, depending on the
285    promoter, Pol II can use both types of DNA opening mechanisms, and suggests that TFIIH-
286    assisted DNA opening originated later in the evolution of cellular DNA-dependent RNA
287    polymerase machineries.

288         Several open questions remain to be addressed for a better understanding of snRNA
289    gene transcription. In particular, SNAPc has been identified to be regulated by its direct
290    interaction with activators that localize ~200 bp upstream of the PSE at the so-called distal
291    sequence element (DSE)[7]. The intervening genomic region between PSE and DSE may be
292    decorated by a nucleosome[8]. In the future, our work may be expanded to studying how DSE
293    binding activators interact with the SNAPc-containing Pol II PIC described here, and how a
294    nucleosome may enable or modulate this interaction. Additionally, our work also serves as
295    stepping stone towards addressing the function of SNAPc in U6 snRNA transcription by Pol
296    III. Such work should provide insights into how SNAPc can interact with both, the Pol II and
297    the Pol III initiation machineries, providing further insights into the evolution of eukaryotic
298    transcription systems.

299

**ONLINE METHODS**

**Cloning and protein expression**

cDNA constructs of SNAPc-FL containing SNAPC4 with an N-terminal StrepTwin-tag and a C-terminal His-tag, SNAPC1, SNAPC2, SNAPC3 and SNAPC5 were subcloned into the pLIB vector. The genes were assembled into a pBIG2ab vector using the biGBac system[43]. The cloned construct was transformed into DH10 EMBacY cells to generate bacmids. Next, the purified bacmid was mixed with CelfectinTM II reagent (Thermo Fisher Scientific) and transfected into 2 ml (density: 0.5 million cells/ml) of adherent Sf9 cells in a 6 well plate. After incubating the plate at 27 °C for 72 h, the resulting supernatant (P1 virus) was collected. To amplify the viral stock, 2 ml of P1 virus was added to 25 ml of Sf9 cells (0.5 million cells/ml) and incubated at 27 °C with shaking at 130 rpm. The supernatant (P2 virus) was collected after 4-5 days of infection when the cell viability dropped to <85% and stored at 4 °C. Large scale protein expression was carried out using 3 x 400 ml of High5 cells (0.5 million cells/ml) by adding 2 ml of P2 virus in each flask and incubated at 27 °C for 4 days at 130 rpm. Cells were then harvested by centrifugation at 250 x g for 10 mins at 4 °C, and pellets were stored at -80 °C. SNAPc-core (SNAPC4 1-516 and lack of the SNAPC2 subunit) was expressed as previously described[18].

**Protein purification**

The insect cells pellet of SNAPc-FL were resuspended in buffer A containing 50 mM HEPES pH 7.8, 750 mM NaCl, 10% glycerol, 15 mM imidazole, 10 mM β-mercaptoethanol, 2 mM MgCl$_2$, 1 mM phenylmethylsulfonyl fluoride (PMSF), 1 μg/mL Aprotinin, 1 μg/mL Pepstatin, and 1 μg/mL Leupeptin, supplemented with four EDTA-free protease inhibitor tablets (Pierce), a scoop of DNAse I, and 10 μl benzonase. Lysis was performed using a dounce homogeniser followed by sonication and the lysate was clarified by centrifugation at 48,000 x *g* at 4 °C for 40 mins. The supernatant was filtered using a 0.45-μm filter, and applied onto a HisTrap HP 5 ml column (GE Healthcare), pre-equilibrated with buffer A. The column was washed with 10 CV of buffer A1 (50 mM HEPES pH 7.8, 500 mM NaCl, 10% glycerol, 50 mM imidazole, 10 mM β-mercaptoethanol, 0.5 mM PMSF and 10 mM O-Phospho-L-serine), and then with 5 CV of buffer A2 (50 mM HEPES pH 7.8, 1250 mM NaCl, 10% glycerol, 50 mM imidazole, 10 mM β-mercaptoethanol, and 0.5 mM PMSF). The column was again washed with 5 CV buffer A1, and the bound protein complex was eluted in buffer B (50 mM HEPES pH 7.8, 500 mM NaCl, 10% glycerol, 300 mM imidazole, 10 mM β-mercaptoethanol, and 0.5 mM PMSF). Next, the sample was diluted to 250 mM NaCl with buffer Heparin A (50 mM HEPES pH 7.8, 10% glycerol, 1 mM TCEP, and 0.1 mM PMSF). The sample was centrifuged at 13,000 rpm for 15 mins at 4 °C and loaded onto a HiTrap Heparin HP 5 ml column (GE healthcare), pre-equilibrated with 12.5% of buffer Heparin B (50 mM HEPES pH 8, 2 M NaCl, 10% glycerol, 1 mM TCEP, and 0.1 mM PMSF). After washing with 5 CV of 12.5% buffer Heparin B, elution was performed through a linear gradient from 15% to 60% over 10 CV. The eluted fractions were analysed by SDS-PAGE, and fractions containing the SNAPc-FL complex were pooled, and concentrated using a 100 kDa molecular weight cut-off (MWCO) VivaSpin concentrator (Sartorius). The sample was centrifuged at 13,000 rpm for 15 mins at 4 °C and applied onto a Superose 6 PG XK 16/70 column (GE Healthcare), pre-equilibrated with 50 mM HEPES pH

343 7.8, 250 mM KCl, 10% glycerol, and 1 mM TCEP. Peak fractions were pooled, concentrated,
344 flash-frozen and stored at -80 °C.

345      SNAPc-core has been purified as previously described[18] with some modifications. Briefly,
346 after cell lysis and centrifugation, the supernatant was subjected to nickel column purification
347 (GE Healthcare) and eluted with 300 mM imidazole. The elution was then further purified with
348 an heparin column and eluted with a gradient from 250mM to 1.25M NaCl. The fractions of
349 interest were pooled, concentrated and subjected to size exclusion chromatography with a S200
350 16/600 equilibrated with 100 mM NaCl, 50 mM HEPES pH 7.9, 10% glycerol and 1mM TCEP.
351 S. *scrofa* Pol II and human initiation factors TBP, TFIIA, TFIIB, TFIIE, TFIIF and TFIIH were
352 purified as previously described[22].

353

354 **Electrophoretic Mobility Shift Assays**
355 EMSA was performed using a 76 bp fragment of U1 promoter DNA (template:5'-GAA ACG
356 TTG TGC CTC TGC CCC GAC ACA GCC TCA TAC GCC TCA CTC TTT ACA CAC ACG
357 GTC ACT TG CCC CGC GCA CT-3' and its complementary strand) and a 75 bp fragment of
358 U5 promoter DNA (template:5'-ACC AGT TAC TTC TGT AAC  TCA ATT TTC GGG TAA
359 CTG CAA TTC CTA GTA CAC TGA TGG TGT CTA CTA ATC CC AAG G-3' and its
360 complementary strand; Integrated DNA Technologies). 20 pM of SNAPc FL or core were
361 incubated with 5 pM of annealed oligonucleotides in presence or absence of 25 pM of TFIIB
362 and TBP in 20 µL of incubation buffer (250 mM NaCl, 50 mM HEPES pH 7.9, 20% glycerol,
363 1 mM TCEP) at room temperature for 15 min. The complexes were resolved on 5%
364 polyacrylamide (37.5:1 acrylamide/bisacrylamide, 10% glycerol, Tris Borate EDTA 1x) gels
365 in 0.5X Tris Borate EDTA running buffer at 40 mA. After staining with Ethidium bromide, the
366 gels were scanned with a Typhoon FLA9500 (GE Healthcare).

367

368 **Promoter-dependent *in vitro* transcription assay**
369 *In vitro* transcription assays were performed as described previously[22,41] with minor alterations.
370 The DNA scaffold (dsDNA) was prepared as reported using a pUC119 vector into which a 92
371 nucleotide fragment of the native U1 snRNA promoter[20] had been inserted. The scaffold (non-
372 template: 5'-GGG CGT GAC CGT GTG TGT AAA GAG TGA GGC GTA TGA GGC TGT
373 GTC GGG GCA GAG GCA CAA CGT TTC GCC CGA AGA TCT CAT ACT TAC CTG
374 GCA GGG CTA AGC TTG GCG TAA TCA TGG TCA TAG CTG TTT CCT GTG TGA AAT
375 TGT TAT CCG CTC ACA ATT CCG CCC-3', template: 5'-GGG CGG AAT TGT GAG CGG
376 ATA ACA ATT TCA CAC AGG AAA CAG CTA TGA CCA TGA TTA CGC CAA GCT
377 TAG CCC TGC CAG GTA AGT ATG AGA TCT TCG GGC GAA ACG TTG TGC CTC
378 TGC CCC GAC ACA GCC TCA TAC GCC TCA CTC TTT ACA CAC ACG GTC ACG
379 CCC-3') was stored in low salt buffer (60 mM KCl, 10 mM K-HEPES pH 7.5, 8 mM MgCl$_2$,
380 3% (v/v) glycerol).

381      Initiation complexes for *in vitro* transcription were reconstituted on scaffold DNA
382 essentially as described[22,41]. All incubation steps were performed at 25 °C unless indicated
383 otherwise. Per sample, 1.6 pmol scaffold, 1.8 pmol Pol II, TFIIE and TFIIH, 5 pmol TBP and
384 TFIIB, 9 pmol TFIIF and TFIIA and 5 pmol SNAPc-FL or SNAPc-core were used. SNAPc
385 was mixed and added to the sample simultaneously with TFIIB. Reactions were prepared in a
386 sample volume of 23.8 µl with final assay conditions of 60 mM KCl, 3 mM K-HEPES pH 7.9,
387 20 mM Tris-HCl pH 7.9, 8 mM MgCl$_2$, 2% (w/v) PVA, 3% (v/v) glycerol, 0.5 mM 1,4-

388 dithiothreitol, 0.5 mg ml$^{-1}$ BSA and 20 units RNase inhibitor. To achieve complete PIC
389 formation, samples were incubated for 45 min at 30 °C. Transcription was started by adding
390 1.2 µl of 10 mM NTP solution and permitted to proceed for 60 min at 30 °C. Reactions were
391 quenched with 100 µl Stop buffer (300 mM NaCl, 10 mM Tris-HCl pH 7.5, 0.5 mM EDTA)
392 and 14 µl 10% SDS, followed by treatment with 4 µg proteinase K (New England Biolabs) for
393 30 min at 37 °C. RNA products were isolated from the samples as described[41], applied to urea
394 gels (7 M urea, 1x TBE, 6% acrylamide:bis-acrylamide 19:1) and separated by denaturing gel
395 electrophoresis (urea-PAGE) in 1x TBE buffer for 45 minutes at 180 volts. Gels were stained
396 for 30 min with SYBR™ Gold (Thermo Fisher Scientific) and RNA was visualized with a
397 Typhoon 9500 FLA imager (GE Healthcare Life Sciences).
398

**Preparation of the SNAPc-containing Pol II PIC**

400 We performed the assembly of SNAPc containing Pol II PIC on snRNA promoters at 25°C
401 essentially as described previously. We used a 96bp fragment of both the native U1 promoter
402 DNA(template:5'-ATC ATG GTA TCT CCC CTG CCA GGT AAG TAT GAA ACG TTG
403 TGC CTC TGC CCC GAC ACA GCC TCA TAC GCC TCA CTC TTT ACA CAC ACGGTC
404 ACT TGC-3';non-template: 5'-GCA AGT GAC CGT GTG TGT AAA GAG TGA GGC GTA
405 TGA GGC TGT GTC GGG GCA GAG GCA CAA CGT TTC ATA CTT ACC TGG CAG
406 GGG AGA TAC CAT GAT-3') and an engineered U5 promoter with 10bp deleted from the
407 downstream edge of the PSE sequence (template: 5'- CCC TGC CAG GTT TTA TGC GAT
408 CTG AAG AGA AAC CAG AGT ATA CCA GTT ACT TCT GTA ACT CAA TTT TCG
409 GGT CCTAGT ACA CTG ATG GTG TCT ACT-3'; non-template: 5'- AGT AGA CAC CAT
410 CAG TGT ACT AGG ACC CGA AAA TTG AGT TAC AGA AGT AAC TGG TAT ACT
411 CTG GTT TCT CTT CAG ATC GCA TAA AAC CTG GCA GGG- 3'). In summary, SNAPc
412 (FL or Core) was pre-incubated for 5 min with the snRNA promoter (U1 or U5) scaffold. It was
413 then mixed with TFIIA-TFIIB and TBP followed by the pre-formed Pol II-TFIIF complex.
414 TFIIE was then added to this mixture and the assembly was incubated at 25°C for 60 min at
415 300 rpm. This reconstituted SNAPc containing Pol II PIC was subjected to 10-30% sucrose-
416 gradient ultra-centrifugation with simultaneous cross-linking using GraFix (Kastner et al.,
417 2008) at 175,000$g$ for 16h at 4°C.  The assay was then fractionated as 200µl aliquots where the
418 crosslinking reaction was quenched using a cocktail of 10mM aspartate and 30mM lysine for
419 10mins. Fractions with SNAPc containing Pol II PIC were dialysed against the cryo-EM sample
420 buffer (25 mM HEPES pH 7.6, 100 mM KCl, 5 mM MgCl$_2$, 1% glycerol and 3 mM TCEP).
421

**Cryo-EM data collection and processing**

423 Samples for cryo-EM were prepared using Quantifoil R3.5/1 holey carbon grids pre-coated
424 with a homemade 3 nm continuous carbon. Four microlitres of SNAPc containing Pol II PIC
425 sample bound to snRNA promoter (U1/U5) was added to the carbon side and incubated for 2.5
426 min. The grids were blotted for 2.5 s and vitrified by plunging into liquid ethane with a Vitrobot
427 Mark IV (FEI Company) set at 4 °C and 100% humidity. Cryo-EM data were collected on a
428 300-kV FEI Titan Krios with a K3 summit direct detector (Gatan) and a GIF quantum energy
429 filter (Gatan) operated with a slit width of 20 eV. Automated data collection was performed
430 with SerialEM at a nominal magnification of 81,000x, corresponding to a pixel size of 1.05
431 Å/pixel[44]. For the sample containing U1 promoter, a total of 16,854 image stacks, with each

432 stack containing 50 frames, were collected at a defocus range of −0.5 to −3.0 μm. All movie
433 frames were contrast transfer function (CTF)-estimated, motion-corrected and dose-weighted
434 using Warp[45]. Particles were picked by Warp using a trained neural network, resulting in
435 5,181,947 particles as a starting set. Subsequent steps of image processing were performed with
436 cryoSPARC[46] and RELION v.3.1.0[47].

437 Particles were extracted with a binning factor of 2 and a box size of 200 pixels (a pixel
438 size of 2.1 Å/pixel) to perform initial clean-up and sorting. The processing scheme was centered
439 around identifying the best SNAPc-containing particle sets. Iterative rounds of 2D-
440 classification followed by heterogenous and homogenous refinements in cryoSPARC, led to
441 two sets of particles corresponding to CC (set-1: 252,067 particles) and OC (set-2: 240,243
442 particles) promoter states respectively. Each set was re-extracted without binning and processed
443 using RELION v.3.1.0, as follows. For set-1, the particles were further sorted by focused 3D
444 classification with a large spherical mask (Mask-1) encompassing the upstream region of PIC
445 containing SNAPc, TBP, TFIIA and TFIIB. This resulted in identifying the best 47,293
446 SNAPc-containing particles. These particles were again subjected to 3D refinement using
447 Mask-1, giving rise to a reconstruction of SNAPc containing Pol II PIC bound to U1 promoter
448 in CC state at 3.4 Å resolution (map-1). In parallel, focused 3D classification of set-2 with a
449 spherical mask (Mask-2) around the upstream region helped to identify the best 137,246 SNAPc
450 containing particles. These particles were then subjected to 3D refinement followed by CTF
451 refinement and Bayesian polishing. Following this, the particles were subject to refinement
452 with and without mask-1 to obtain of SNAPc containing Pol II PIC bound to U1 promoter in
453 OC state at 3.0 Å (map-2) and a local map spanning the SNAPc containing upstream region at
454 3.7 Å resolution(map-3).

455 For the sample containing U5 promoter dataset, 4,842image stacks, with each stack
456 containing 60 frames, were collected at a defocus range of −0.3 to −2.5 μm. All movie frames
457 were contrast transfer function (CTF)-estimated, motion-corrected and dose-weighted using
458 Warp[45]. Particles were picked by Warp using a trained neural network, resulting in 1,299,523
459 particles. Subsequent steps of image processing were performed with cryoSPARC[46] and
460 RELION v.3.1.0[47]. Particles were extracted with a binning factor of 4 and a box size of 100
461 pixels (a pixel size of 4.2 Å/pixel) to perform initial clean-up and sorting. After sorting in
462 cryoSPARC using 2D-classification followed by heterogenous and homogenous refinements,
463 a particle set (set-3: 443,960 particles) in CC promoter state was re-extracted with 2x binning
464 (a pixel size of 2.1 Å/pixel) and processed using RELION v.3.1.0, as follows. For set-3, the
465 particles were further sorted by 3D classification followed by focused 3D classification using
466 Mask-1. The resulting 159,144 particles were re-extracted without a binning factor and were
467 subjected to CTF refinement and Bayesian polishing. These particles were then subjected to
468 another round of masked classification yielding 85,787 SNAPc-containing particles. These
469 particles were then 3D refined without and with Mask-1, giving rise to a reconstruction of
470 SNAPc containing Pol II PIC bound to U5 promoter in CC state at 3.0 Å resolution (map-4)
471 and a local map of the SNAPc containing upstream complex extending to 3.2 Å(map-5).

472 The reported resolutions were calculated on the basis of the gold standard Fourier shell
473 correlation (FSC) 0.143 criterion. After processing of the final reconstructions, B-factor
474 sharpening was performed for all final maps on the basis of automatic B-factor determination
475 in RELION (−5 Å$^2$ for map-1: SNAPc-PIC bound to U1 promoter in CC state, -10 Å$^2$ for map-
476 2: SNAPc-PIC bound to U1 promoter in OC state and −10 Å$^2$ for map-3: local map of SNAPc

477  containing upstream complex, -10 Å² for map-4: SNAPc-PIC bound to U5 promoter in CC state
478  and −10 Å² for map-5: local map of SNAPc containing upstream complex). Estimates of local
479  resolution were calculated using the in-built local-resolution tool of RELION and the estimated
480  B-factors. To assist model building, a local-resolution-filtered map (but unsharpened) of map-
481  5 was sharpened locally using PHENIX.auto_sharpen[48].

482

### Model building and refinement

484  The PIC was modelled using the core PIC part of the previously published high resolution
485  structures in closed and open promoter states[22]. For SNAPc, the subunits SNAPC1 and
486  SNAPC4 were built using partial homology models generated using TrRosetta[49]. The partial
487  models were rigid body fitted into the density using UCSF Chimera[50] and were manually
488  extended and corrected using Coot[51] to fit the density. The subunit SNAPC3 was modelled
489  entirely de novo using the experimental density in Coot. Ambiguous density corresponding to
490  linker regions were not modelled. The model corresponding to the wing-2 region constituting
491  parts of SNAPC1, SNAPC3 and SNAPC4 was modelled using AlphaFold[27]. The model for
492  promoter DNA in CC and OC states was obtained using the high-resolution structures of human
493  PIC as template where in the sequence register was mutated to fit the U1 and U5 respectively.
494  The models were then subjected to iterative rounds of PHENIX real-space refinement followed
495  by manual adjustment in coot to achieve final models with good stereochemistry as assessed
496  by MolProbity[52]. Figures representing the 3D structures and maps were prepared using PyMOL,
497  UCSF Chimera and UCSF ChimeraX.

498

### Crosslinking mass-spectrometry

500  To prepare a sample for performing crosslinking mass-spectrometry, a stable complex of
501  SNAPc-containing Pol-II PIC bound to U5 promoter was isolated. An assay containing Pol II,
502  TBP, TFIIA, TFIIB, TFIIF and SNAPc-FL was incubated in ratios explained above and was
503  subjected to size-exclusion chromatography using Superose 6 increase 3.2/300 GL column (GE
504  Healthcare) pre-equilibrated with buffer-x (25mM Hepes pH 7.5, 100mM NaCl, 5mM MgCl2,
505  5% glycerol and 2mM TCEP). The peak fractions were then pooled and incubated with 1mM
506  of Bissulfosuccinimidyl suberate (BS3) for 45 min at 4° C. The crosslinking reaction was
507  quenched using a cocktail of 10 mM aspartate and 30 mM lysine.

508  Crosslinked proteins were resuspended in 4 M urea/ 50 mM ammonium bicarbonate for
509  10 min at 25°C and reduced for 30 min at RT with 10 mM dithiotreitol (DTT). Proteins were
510  alkylated for 30 min at RT in the dark by adding iodacetamide (IAA) to a final concentration
511  of 55 mM. Sample was diluted to 1M Urea and digested for 30 min at 37 °C with 4 µl Pierce
512  Universal Nuclease (250 U/µl) in the presence of 2 mM MgCl2. Trypsin (Promega) digest was
513  performed o/n at 37 °C in a 1:50 enzyme/protein ratio, the reaction was terminated with 0.2 %
514  (v/v) FA. Tryptic peptides were desalted on MicroSpin Columns (Harvard Apparatus)
515  following manufacturer´s instruction and vacuum-dried. Cross-linked peptides were
516  resuspended in 50 µl 30 % acetonitrile/0.1 % TFA and enriched by peptide size exclusion
517  chromatography/pSEC (Superdex Peptide PC3.2/300 column, GE Healthcare, flow rate 50
518  µl/min).

519  Crosslinked peptides derived from pSEC were subjected to liquid chromatography mass
520  spectrometry (LC-MS) on a Thermo Obitrap Exploris mass spectrometer. Peptides were loaded
521  in duplicates onto a Dionex Ultimate 3000 RSLCnano equipped with a custom column

522 (ReproSil-Pur 120 C18-AQ, 1.9 µm pore size, 75 µm inner diameter, 30 cm length, Dr. Maisch
523 GmbH). Peptides were separated applying the following gradient: mobile phase A consisted of
524 0.1 % formic acid (FA, v/v), mobile phase B of 80 % ACN/0.08 % FA (v/v). The gradient
525 started at 5 % B, increasing to 10, 15 or 20 % B within 3 min, followed by a continuous increase
526 to 48 % B within 45 min, then keeping B constant at 90 % for 8 min. After each gradient, the
527 column was again equilibrated to 5 % B for 2 min. The flow rate was set to 300 nL/min.

528       MS1 spectra were acquired with a resolution of 120,000 in the orbitrap (OT) covering
529 a mass range of 380–1600 m/z. Dynamic exclusion was set to 30 s. Only precursors with a
530 charge state of 3-8 were included. MS2 spectra were recorded with a resolution of 30,000 in
531 OT and the isolation window to 1.6  m/z. Fragmentation was enforced by higher-energy
532 collisional dissociation (HCD) at 30 %. Raw files were searched against a database containing
533 the sequences of the proteins of the complex and analyzed via pLink 2.3.9 at a false discovery
534 rate (FDR) of 1%[53]. Carbamidomethylation of cysteines was set as fixed modification,
535 oxidation of methionines as variable modification. The database contained all proteins within
536 the complex. For further analysis only interaction sites with 3 cross-linked peptide spectrum
537 matches were taken into account. Cross-links were displayed with xiNET and XlinkAnalyzer
538 in UCSF Chimera.[50,54,55]

539

**TSS precision analyses in cells**

541 We utilized published 5'cap-seq data[35] (GEO: GSE159633) for analyses of TSS precision in
542 cells. The raw data were processed as described previously[35] to obtain the 5'-ends of reads and
543 generate normalized coverage. In brief, we first removed the unique molecular identifier (UMI)
544 from 5cap-seq  reads with UMI-tools[56] and then trimmed adapter sequences with Cutadapt [57]
545 and mapped to the human genome (GRCh38) merged with the D. melanogaster genome (Dm6)
546 with the STAR mapper[58]. We next deduplicated the mapped data with UMI-tools to remove
547 any PCR duplicates and then determined the first transcribed base and used this position in
548 downstream analyses. Normalization factors were obtained from the spike-in reads (processed
549 as above) that mapped to the spike-in genome and used to normalize the human genome
550 coverage profiles. The replicates were combined by summing the normalized coverage per nt.
551 Thus, obtaining genome-wide capped 5'-end signal (5'cap-seq signal) at single-base resolution.
552 We subset the NCBI reference genome annotation[59] (GRCh38.p7) to only contain genes
553 annotated to the primary assembly and included only genes with known transcripts (prefix:
554 "NR" or "NM") and also excluded overlapping genes. To exclude genes with alternative start
555 sites from downstream analyses we included only genes that have a constitutive first or a single
556 exon in our downstream analyses.

557       To determine the main TSS we determined the position with the highest 5'cap-seq signal
558 within constitutive first exons of the reference annotation. To accommodate for reference
559 annotation imprecision, we also included 10 bp upstream of the annotated TSS and set the
560 downstream cutoff to 500 bp downstream of the annotated TSS. We thus obtained the main
561 TSS for each constitutive TSS. We next quantified the 5'cap-seq signal of the main TSS ($\pm$2
562 bp) and the TSS region (main TSS $\pm$50 bp). We excluded genes with less than 10 counts in the
563 TSS region and genes with biotypes that are not either protein-coding or snRNA. From the
564 remaining annotated snRNA subset we also removed know Pol III-transcripts: RN7SK,
565 RNU6ATAC, SNAR-G2, RNU6-2, SNAR-C4, SNAR-G1, SNAR-C3 and identified with
566 protein-coding gene promoters contain a TATA-box motif (JASPAR database, 2020 release:

567 https://jaspar2020.genereg.net/matrix/POL012.1/) within 50 bp upstream of the annotated TSS.
568 Finally, we determined the TSS precision score by dividing the TSS peak counts by the TSS
569 region counts. The maximum TSS precision score is 1, which means that all 5'cap-seq signal
570 is within the TSS peak. The preprocessed 5'cap-seq data was analyzed in RStudio[60] utilizing R
571 version 3.6.1[61] and packages from the Bioconductor repository [62,63] and Tidyverse[64]. Plots were
572 generated with ggplot2 and ggbio [65].
573

### Acknowledgements

584

### Author Contributions

586 S.R. carried out all experiments and data analysis, unless stated otherwise. S.S. performed the
587 *in vitro* transcription assay and quantification. T.K. and J.G. cloned, expressed and purified the
588 SNAPc variants and performed EMSA assays. K.Z. performed the reanalysis of 5'-capseq data,
589 TSS precision plots and the web-logo plots. J.S. performed crosslinking mass-spectrometry and
590 data analysis and was supervised by H.U. S.R. and C.D. collected the cryo-EM datasets. P.C.
591 and A.V. designed and supervised research. S.R. and P.C. interpreted the data and wrote the
592 manuscript, with input from all authors.
593

### Competing interests

595 The authors declare no competing interests.
596

**Supplementary Information** is available online at https//doi.org/XYZ.

598

599 **Correspondence** Correspondence and request of materials and resources should be addressed
600 to P.C. (patrick.cramer@mpinat.mpg.de).
601

### Data availability

603 The cryo-EM density reconstructions were deposited to the EMDB under accession codes
604 EMD-AAAA, -BBBB, -CCCC, -DDDD, -EEEE and atomic coordinates were deposited to the
605 PDB under the accession codes PDB-AAAA, -BBBB, -CCCC, -DDDD, -EEEE. All data is
606 available in the main text or the supplementary materials.

**References**

1    Greber, B. J. & Nogales, E. The Structures of Eukaryotic Transcription Pre-initiation Complexes and Their Functional Implications. *Subcell Biochem* **93**, 143-192, doi:10.1007/978-3-030-28151-9_5 (2019).

2    Nogales, E., Louder, R. K. & He, Y. Structural Insights into the Eukaryotic Transcription Initiation Machinery. *Annu Rev Biophys* **46**, 59-83, doi:10.1146/annurev-biophys-070816-033751 (2017).

3    Osman, S. & Cramer, P. Structural Biology of RNA Polymerase II Transcription: 20 Years On. *Annu Rev Cell Dev Biol* **36**, 1-34, doi:10.1146/annurev-cellbio-042020-021954 (2020).

4    Sainsbury, S., Bernecky, C. & Cramer, P. Structural basis of transcription initiation by RNA polymerase II. *Nat Rev Mol Cell Biol* **16**, 129-143, doi:10.1038/nrm3952 (2015).

5    Will, C. L. & Luhrmann, R. Spliceosome structure and function. *Cold Spring Harb Perspect Biol* **3**, doi:10.1101/cshperspect.a003707 (2011).

6    Egloff, S., O'Reilly, D. & Murphy, S. Expression of human snRNA genes from beginning to end. *Biochem Soc Trans* **36**, 590-594, doi:10.1042/BST0360590 (2008).

7    Dergai, O. & Hernandez, N. How to Recruit the Correct RNA Polymerase? Lessons from snRNA Genes. *Trends Genet* **35**, 457-469, doi:10.1016/j.tig.2019.04.001 (2019).

8    Guiro, J. & Murphy, S. Regulation of expression of human RNA polymerase II-transcribed snRNA genes. *Open Biol* **7**, doi:10.1098/rsob.170073 (2017).

9    Jawdekar, G. W. & Henry, R. W. Transcriptional regulation of human small nuclear RNA genes. *Biochim Biophys Acta* **1779**, 295-305, doi:10.1016/j.bbagrm.2008.04.001 (2008).

10   Mittal, V., Ma, B. & Hernandez, N. SNAP(c): a core promoter factor with a built-in DNA-binding damper that is deactivated by the Oct-1 POU domain. *Genes Dev* **13**, 1807-1821, doi:10.1101/gad.13.14.1807 (1999).

11   Jawdekar, G. W. *et al.* The unorthodox SNAP50 zinc finger domain contributes to cooperative promoter recognition by human SNAPC. *J Biol Chem* **281**, 31050-31060, doi:10.1074/jbc.M603810200 (2006).

12   Wong, M. W. *et al.* The large subunit of basal transcription factor SNAPc is a Myb domain protein that interacts with Oct-1. *Mol Cell Biol* **18**, 368-377, doi:10.1128/mcb.18.1.368 (1998).

13   Das, A. & Bellofatto, V. RNA polymerase II-dependent transcription in trypanosomes is associated with a SNAP complex-like transcription factor. *Proc Natl Acad Sci U S A* **100**, 80-85, doi:10.1073/pnas.262609399 (2003).

14   Su, Y. *et al.* Characterization of a Drosophila proximal-sequence-element-binding protein involved in transcription of small nuclear RNA genes. *Eur J Biochem* **248**, 231-237, doi:10.1111/j.1432-1033.1997.t01-1-00231.x (1997).

15   Henry, R. W., Mittal, V., Ma, B., Kobayashi, R. & Hernandez, N. SNAP19 mediates the assembly of a functional core promoter complex (SNAPc) shared by RNA polymerases II and III. *Genes Dev* **12**, 2664-2672, doi:10.1101/gad.12.17.2664 (1998).

16   Ma, B. & Hernandez, N. A map of protein-protein contacts within the small nuclear RNA-activating protein complex SNAPc. *J Biol Chem* **276**, 5027-5035, doi:10.1074/jbc.M009301200 (2001).

17   Kuhlman, T. C., Cho, H., Reinberg, D. & Hernandez, N. The general transcription factors IIA, IIB, IIF, and IIE are required for RNA polymerase II transcription from the

654        human U1 small nuclear RNA promoter. *Mol Cell Biol* **19**, 2130-2141,
655        doi:10.1128/MCB.19.3.2130 (1999).
656    18    Dergai, O. *et al.* Mechanism of selective recruitment of RNA polymerases II and III to
657        snRNA gene promoters. *Genes Dev* **32**, 711-722, doi:10.1101/gad.314245.118 (2018).
658    19    Compe, E. & Egly, J. M. Nucleotide Excision Repair and Transcriptional Regulation:
659        TFIIH and Beyond. *Annu Rev Biochem* **85**, 265-290, doi:10.1146/annurev-biochem-
660        060815-014857 (2016).
661    20    James Faresse, N. *et al.* Genomic study of RNA polymerase II and III SNAPc-bound
662        promoters reveals a gene transcribed by both enzymes and a broad use of common
663        activators. *PLoS Genet* **8**, e1003028, doi:10.1371/journal.pgen.1003028 (2012).
664    21    Kastner, B. *et al.* GraFix: sample preparation for single-particle electron
665        cryomicroscopy. *Nat Methods* **5**, 53-55, doi:10.1038/nmeth1139 (2008).
666    22    Aibara, S., Schilbach, S. & Cramer, P. Structures of mammalian RNA polymerase II
667        pre-initiation complexes. *Nature* **594**, 124-128, doi:10.1038/s41586-021-03554-8
668        (2021).
669    23    Kim, J. L., Nikolov, D. B. & Burley, S. K. Co-crystal structure of TBP recognizing the
670        minor groove of a TATA element. *Nature* **365**, 520-527, doi:10.1038/365520a0
671        (1993).
672    24    Tan, S., Hunziker, Y., Sargent, D. F. & Richmond, T. J. Crystal structure of a yeast
673        TFIIA/TBP/DNA complex. *Nature* **381**, 127-151, doi:10.1038/381127a0 (1996).
674    25    Archuleta, T. L. *et al.* Structure and evolution of ENTH and VHS/ENTH-like domains in
675        tepsin. *Traffic* **18**, 590-603, doi:10.1111/tra.12499 (2017).
676    26    Hung, K. H. & Stumph, W. E. Regulation of snRNA gene expression by the Drosophila
677        melanogaster small nuclear RNA activating protein complex (DmSNAPc). *Crit Rev*
678        *Biochem Mol Biol* **46**, 11-26, doi:10.3109/10409238.2010.518136 (2011).
679    27    Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature*
680        **596**, 583-589, doi:10.1038/s41586-021-03819-2 (2021).
681    28    Hinkley, C. S., Hirsch, H. A., Gu, L., LaMere, B. & Henry, R. W. The small nuclear RNA-
682        activating protein 190 Myb DNA binding domain stimulates TATA box-binding
683        protein-TATA box recognition. *J Biol Chem* **278**, 18649-18657,
684        doi:10.1074/jbc.M204247200 (2003).
685    29    Lu, X. J., Shakked, Z. & Olson, W. K. A-form conformational motifs in ligand-bound
686        DNA structures. *J Mol Biol* **300**, 819-840, doi:10.1006/jmbi.2000.3690 (2000).
687    30    Patel, A. B. *et al.* Structure of human TFIID and mechanism of TBP loading onto
688        promoter DNA. *Science* **362**, doi:10.1126/science.aau8872 (2018).
689    31    Dienemann, C., Schwalb, B., Schilbach, S. & Cramer, P. Promoter Distortion and
690        Opening in the RNA Polymerase II Cleft. *Mol Cell* **73**, 97-106 e104,
691        doi:10.1016/j.molcel.2018.10.014 (2019).
692    32    Holstege, F. C., Tantin, D., Carey, M., van der Vliet, P. C. & Timmers, H. T. The
693        requirement for the basal transcription factor IIE is determined by the helical stability
694        of promoter DNA. *EMBO J* **14**, 810-819 (1995).
695    33    Haberle, V. & Stark, A. Eukaryotic core promoters and the functional basis of
696        transcription initiation. *Nat Rev Mol Cell Biol* **19**, 621-637, doi:10.1038/s41580-018-
697        0028-8 (2018).
698    34    Wang, W., Carey, M. & Gralla, J. D. Polymerase II promoter activation: closed
699        complex formation and ATP-driven start site opening. *Science* **255**, 450-453,
700        doi:10.1126/science.1310361 (1992).

701  35  Zumer, K. *et al.* Two distinct mechanisms of RNA polymerase II elongation
702      stimulation in vivo. *Mol Cell* **81**, 3096-3109 e3098, doi:10.1016/j.molcel.2021.05.028
703      (2021).
704  36  Abascal-Palacios, G., Ramsay, E. P., Beuron, F., Morris, E. & Vannini, A. Structural
705      basis of RNA polymerase III transcription initiation. *Nature* **553**, 301-306,
706      doi:10.1038/nature25441 (2018).
707  37  Han, Y., Yan, C., Fishbain, S., Ivanov, I. & He, Y. Structural visualization of RNA
708      polymerase III transcription machineries. *Cell Discov* **4**, 40, doi:10.1038/s41421-018-
709      0044-z (2018).
710  38  Pilsl, M. & Engel, C. Structural basis of RNA polymerase I pre-initiation complex
711      formation and promoter melting. *Nat Commun* **11**, 1206, doi:10.1038/s41467-020-
712      15052-y (2020).
713  39  Sadian, Y. *et al.* Molecular insight into RNA polymerase I promoter recognition and
714      promoter melting. *Nat Commun* **10**, 5543, doi:10.1038/s41467-019-13510-w (2019).
715  40  Vorlander, M. K., Khatter, H., Wetzel, R., Hagen, W. J. H. & Muller, C. W. Molecular
716      mechanism of promoter opening by RNA polymerase III. *Nature* **553**, 295-300,
717      doi:10.1038/nature25440 (2018).
718  41  Schilbach, S., Aibara, S., Dienemann, C., Grabbe, F. & Cramer, P. Structure of RNA
719      polymerase II pre-initiation complex at 2.9 A defines initial DNA opening. *Cell* **184**,
720      4064-4072 e4028, doi:10.1016/j.cell.2021.05.012 (2021).
721  42  Fouqueau, T. *et al.* The cutting edge of archaeal transcription. *Emerg Top Life Sci* **2**,
722      517-533, doi:10.1042/ETLS20180014 (2018).
723  43  Weissmann, F. *et al.* biGBac enables rapid gene assembly for the expression of large
724      multisubunit protein complexes. *Proc Natl Acad Sci U S A* **113**, E2564-2569,
725      doi:10.1073/pnas.1604935113 (2016).
726  44  Mastronarde, D. N. Automated electron microscope tomography using robust
727      prediction of specimen movements. *J Struct Biol* **152**, 36-51,
728      doi:10.1016/j.jsb.2005.07.007 (2005).
729  45  Tegunov, D. & Cramer, P. Real-time cryo-electron microscopy data preprocessing
730      with Warp. *Nat Methods* **16**, 1146-1152, doi:10.1038/s41592-019-0580-y (2019).
731  46  Punjani, A., Rubinstein, J. L., Fleet, D. J. & Brubaker, M. A. cryoSPARC: algorithms for
732      rapid unsupervised cryo-EM structure determination. *Nat Methods* **14**, 290-296,
733      doi:10.1038/nmeth.4169 (2017).
734  47  Zivanov, J. *et al.* New tools for automated high-resolution cryo-EM structure
735      determination in RELION-3. *Elife* **7**, doi:10.7554/eLife.42166 (2018).
736  48  Liebschner, D. *et al.* Macromolecular structure determination using X-rays, neutrons
737      and electrons: recent developments in Phenix. *Acta Crystallogr D Struct Biol* **75**, 861-
738      877, doi:10.1107/S2059798319011471 (2019).
739  49  Yang, J. *et al.* Improved protein structure prediction using predicted interresidue
740      orientations. *Proc Natl Acad Sci U S A* **117**, 1496-1503, doi:10.1073/pnas.1914677117
741      (2020).
742  50  Pettersen, E. F. *et al.* UCSF Chimera--a visualization system for exploratory research
743      and analysis. *J Comput Chem* **25**, 1605-1612, doi:10.1002/jcc.20084 (2004).
744  51  Casanal, A., Lohkamp, B. & Emsley, P. Current developments in Coot for
745      macromolecular model building of Electron Cryo-microscopy and Crystallographic
746      Data. *Protein Sci* **29**, 1069-1078, doi:10.1002/pro.3791 (2020).
747  52  Prisant, M. G., Williams, C. J., Chen, V. B., Richardson, J. S. & Richardson, D. C. New
748      tools in MolProbity validation: CaBLAM for CryoEM backbone, UnDowser to rethink

749      "waters," and NGL Viewer to recapture online 3D graphics. *Protein Sci* **29**, 315-329,
750      doi:10.1002/pro.3786 (2020).

751    53    Chen, Z. L. *et al.* A high-speed search engine pLink 2 with systematic evaluation for
752      proteome-scale identification of cross-linked peptides. *Nat Commun* **10**, 3404,
753      doi:10.1038/s41467-019-11337-z (2019).

754    54    Combe, C. W., Fischer, L. & Rappsilber, J. xiNET: cross-link network maps with residue
755      resolution. *Mol Cell Proteomics* **14**, 1137-1147, doi:10.1074/mcp.O114.042259
756      (2015).

757    55    Kosinski, J. *et al.* Xlink Analyzer: software for analysis and visualization of cross-linking
758      data in the context of three-dimensional structures. *J Struct Biol* **189**, 177-183,
759      doi:10.1016/j.jsb.2015.01.014 (2015).

760    56    Smith, T., Heger, A. & Sudbery, I. UMI-tools: modeling sequencing errors in Unique
761      Molecular Identifiers to improve quantification accuracy. *Genome Res* **27**, 491-499,
762      doi:10.1101/gr.209601.116 (2017).

763    57    Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing
764      reads. *2011* **17**, 3, doi:10.14806/ej.17.1.200 (2011).

765    58    Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,
766      doi:10.1093/bioinformatics/bts635 (2013).

767    59    O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,
768      taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**, D733-745,
769      doi:10.1093/nar/gkv1189 (2016).

770    60    Team, R. RStudio: Integrated Development for R. *(Boston, MA, RStudio, PBC)* (2020).

771    61    Team, R. C. R. *(Vienna, Austria, R Foundation for Statistical Computing)* (2019).

772    62    Gentleman, R. C. *et al.* Bioconductor: open software development for computational
773      biology and bioinformatics. *Genome Biol* **5**, R80, doi:10.1186/gb-2004-5-10-r80
774      (2004).

775    63    Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor.
776      *Nat Methods* **12**, 115-121, doi:10.1038/nmeth.3252 (2015).

777    64    Wickham, H. ggplot2: Elegant Graphics for Data Analysis. *(Springer-Verlag New York)*
778      (2016).

779    65    Yin, T., Cook, D. & Lawrence, M. ggbio: an R package for extending the grammar of
780      graphics for genomic data. *Genome Biol* **13**, R77, doi:10.1186/gb-2012-13-8-r77
781      (2012).

782    66    Notredame, C., Higgins, D. G. & Heringa, J. T-Coffee: A novel method for fast and
783      accurate multiple sequence alignment. *J Mol Biol* **302**, 205-217,
784      doi:10.1006/jmbi.2000.4042 (2000).

785    67    Robert, X. & Gouet, P. Deciphering key features in protein structures with the new
786      ENDscript server. *Nucleic Acids Res* **42**, W320-324, doi:10.1093/nar/gku316 (2014).

787
788
789
790
791
792
793
794
795
796
797

798
799
800 **FIGURES**
801

802 **Figure 1 | Preparation of SNAPc-containing Pol II PIC on non-coding RNA promoters**

803 a) SDS-PAGE analysis of SNAPc variants (FL, core) purified to homogeneity.

804 b) EMSA shows the binding of SNAPc (± TBP, TFIIB) to U1 and U5 promoter DNA. The

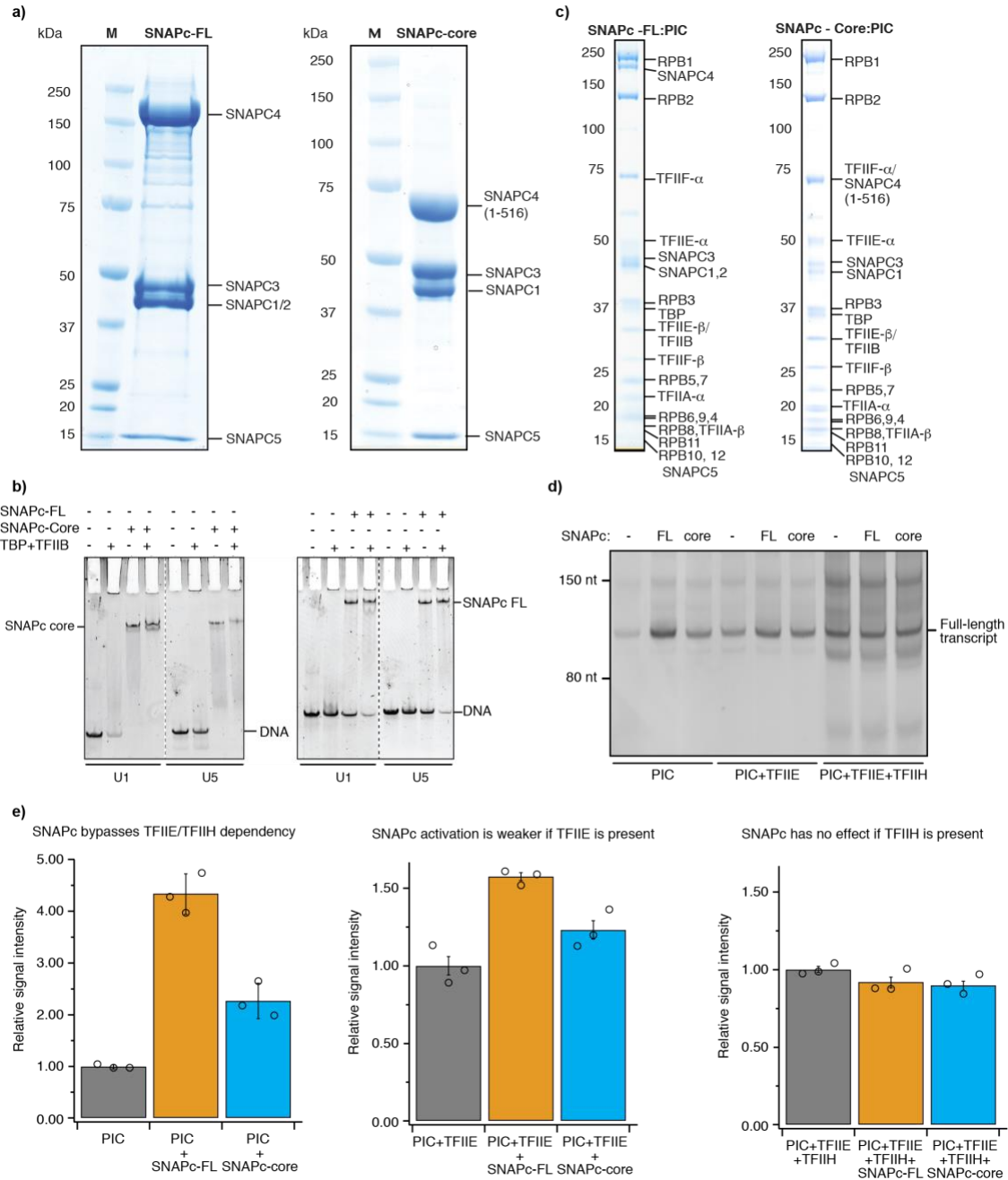805 presence of SNAPc stabilises the binding of TBP-TFIIB to snRNA promoters.

806 c) SDS-PAGE analysis of SNAPc containing Pol II PIC variants isolated through a sucrose

807 gradient ultracentrifugation.

808 d) In vitro transcription assay showing the relative influence of SNAPc variants on Pol II

809 snRNA transcription with different combinations of GTFs'.

810 e) Histogram plots representing the quantification (Methods) of full-length transcripts from the

811 *in vitro* transcription assay in panel d.

812

813

**Figure 2 | Overall structure of SNAPc-containing Pol II PIC**

a) Schematic 2D representation of the U1 and U5 promoter sequences highlighting the binding motifs of the initiation machinery as observed in the cryo-EM structure: PSE (SNAPc), TBP binding site (TBP) and TSS (Pol II). The transcription start site (TSS) is denoted +1 and negative and positive numbers indicate upstream and downstream positions.

b) Cartoon representation of the SNAPc-containing Pol II PIC as viewed from the front and top. The colour codes for Pol II and the GTFs' are consistently used throughout.

**Figure 3 | Structure of SNAPc**

a) 2D-domain schematics of individual SNAPc subunits. The regions visible in the 3D structure are marked by dotted-lines.

b) SNAPc structure in cartoon representation. Domain nomenclature and colours are used as described in panel a. Dashed boxes indicate the interfaces between the subunits.

c, d) Close up view of interfaces 1 and 2 that are formed between SNAPC1 (pink) and SNAPC3 (orange). The residues V115, F120, A123, Y124 of SNAPC1 and V86, L90, L100, C104 of SNAPC3 form mainly hydrophobic interactions, whereas ionic interactions are formed between R34, R47, K96, R128 of SNAPC1 and D89, D107, E153, D332 of SNAPC3. F54 of SNAPC1 and R133 of SNAPC3 form a cation-pi interaction and N49 of SNAPC1 and Y157 of SNAPC3 form polar contacts. Similarly in interface 2: SNAPC1 L101, W104, F137 and A139 form hydrophobic contacts with F47, L50, W51, L55 and L305 of SNAPC3. Salt-bridges involving R98, D105 of SNAPC1 and R54 and D325 of SNAPC3 fortify interface 2.

e, f) Interfaces 3 and 4 between SNAPC3 (orange) and SNAPC4 (chestnut brown). In interface 3, SNAPC3 residues F155, W348, F355, V356, Y360, T361, P372, F377, T409 form the bulk of hydrophobic contacts with F331, L334, L344 and H369 of SNAPC4 (Figure 3e). Likewise in interface 4 the residues Y253, I274, W277, P308 and L310 make hydrophobic contacts with the amino acids F140, Y149, F150, F176 of SNAPC4. Additional salt bridges are formed by R133, R283 of SNAPC3 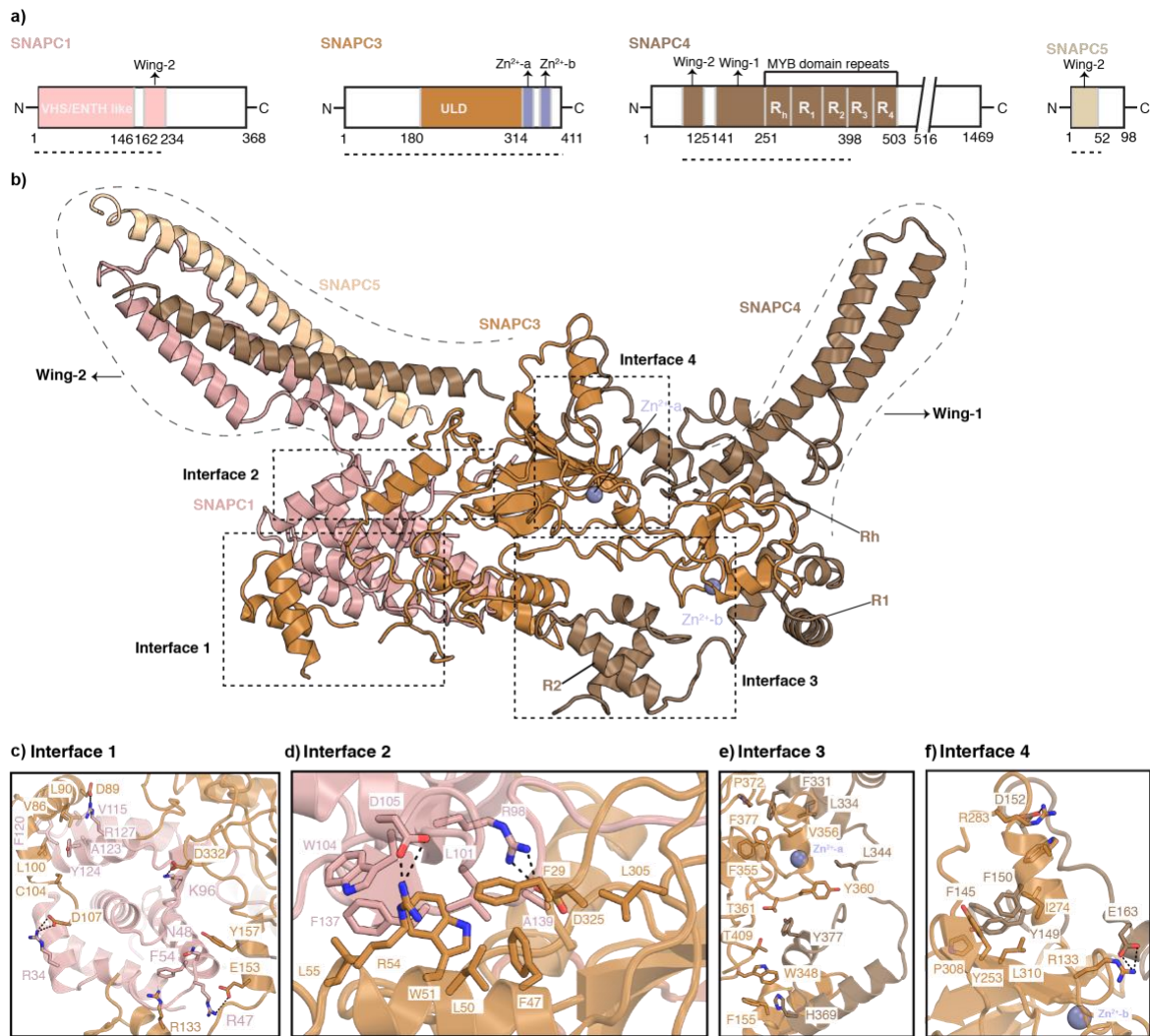with D152 and E153 of SNAPC4. The Zn-fingers (ZF-1, ZF-2) of SNAPC3 are in close proximity to the interfaces 3 and 4, and would be important for the structural integrity of this complex. The residues involved in these protein-protein interaction surfaces are highly conserved across metazoans (Extended Data Figure 7).

21

846



847
848
849
850
851
852
853
854
855
856

**Figure 4 | SNAPc-DNA interactions**
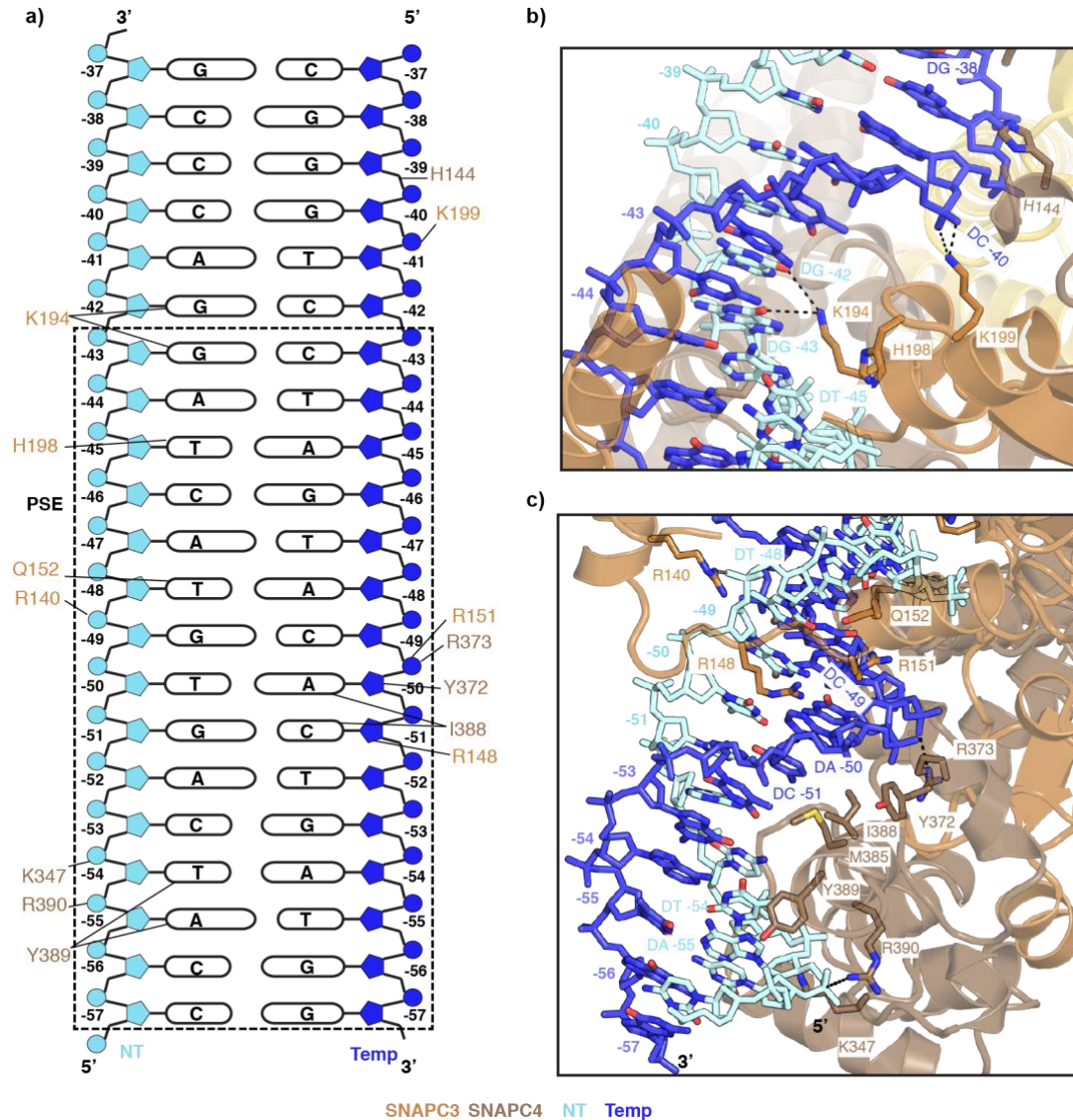
a) Schematic view of the protein-DNA interactions between SNAPc and the PSE motif. Residues interacting with specific regions of the DNA as described in the text are indicated by lines. In panels b and c, nucleotide residues are numbered in atomic colour to indicate the strand and the DNA register

b) DNA-protein interaction network on the preceding major and minor grooves (register: -46 to -35 ) of PSE as bound by SNAPc subunits SNAPC3 and SNAPC4. Colour codes are used uniformly in all panels.

865  c) Close up view of the first major and minor groove (register: -57 to -47 ) interactions between
866  SNAPc and the PSE motif on U5 promoter. The SNAPc subunits are represented as cartoon,
867  whereas the interacting amino acid sidechain residues, DNA chains are depicted as sticks with
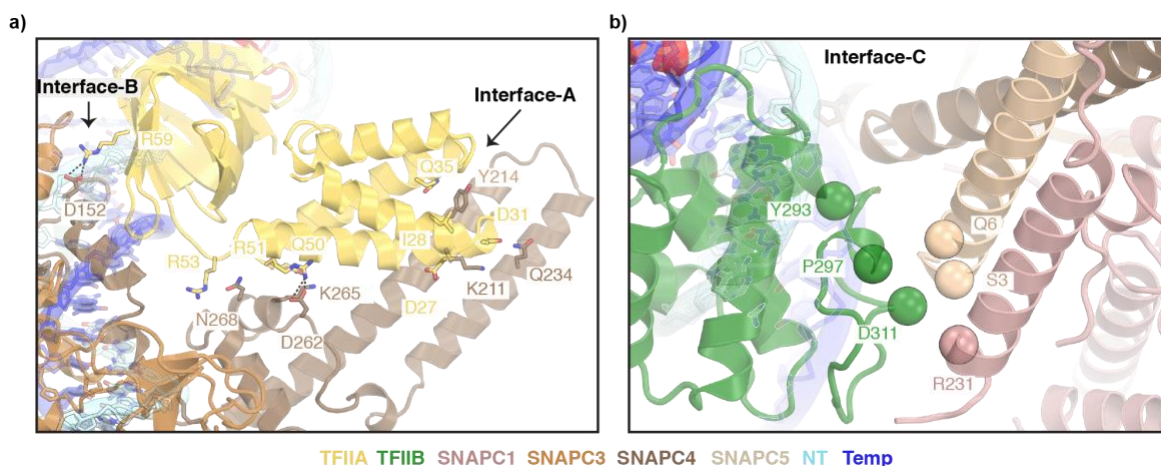868  atomic colours. Dashed lines indicate ionic interactions.
869



870
871
872
873
874
875
876  **Figure 5 | SNAPc-general transcription factors interaction**
877  a) Close up view of wing-1:TFIIA interaction. The amino acid residues involved in the
878  formation of interfaces A and B between TFIIA (yellow orange) and SNAPC4 (chestnut brown)
879  are represented as sticks. Dashed lines indicate salt-bridges.
880  b) Zoomed in view of the interface C formed between wing-2 and TFIIB C-terminal cyclin fold.
881  The Cα atoms of putative residues forming the interaction surface are represented as spheres.
882

23

TFIIA TFIIB SNAPC1 SNAPC3 SNAPC4 SNAPC5 NT Temp
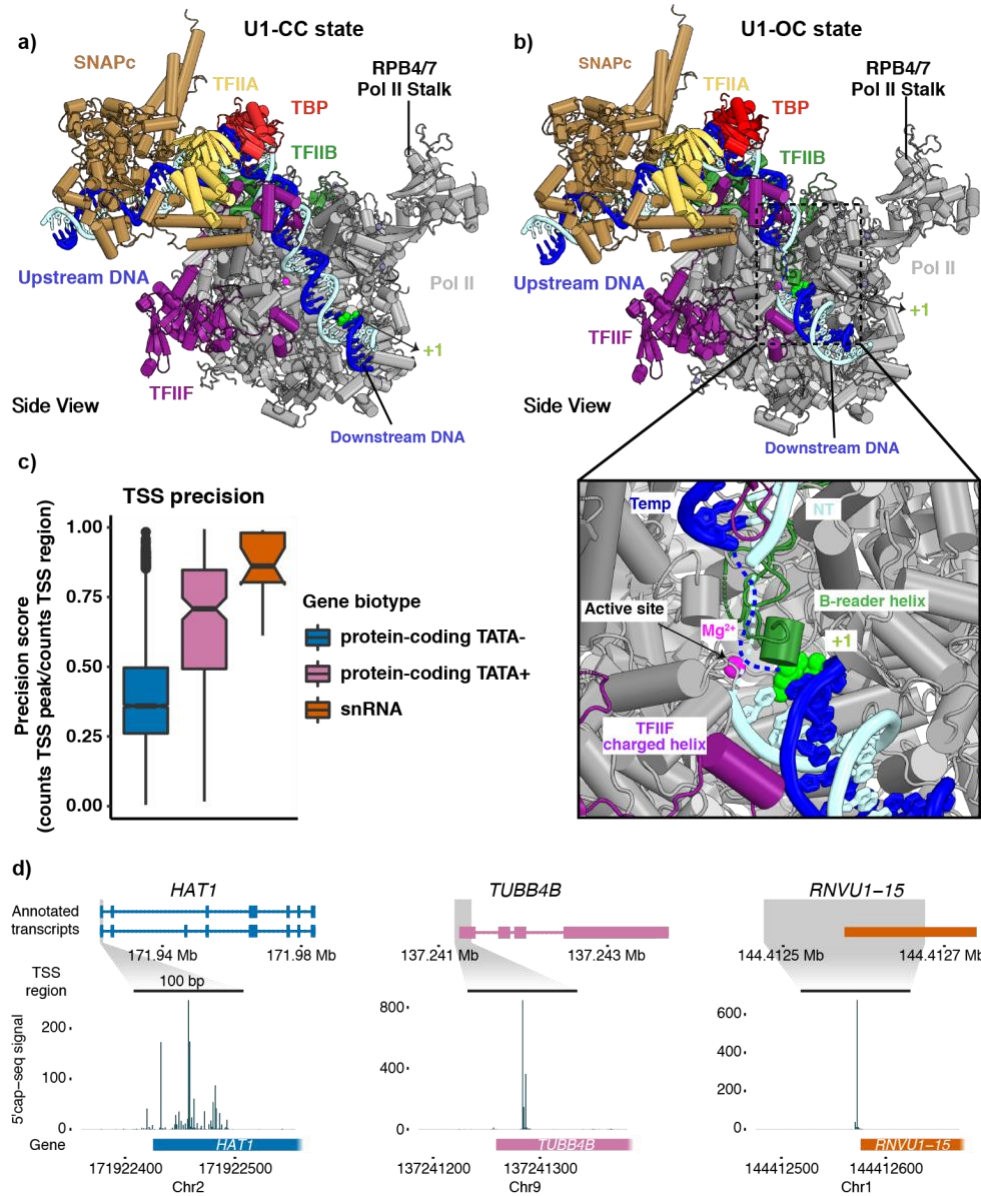
**Figure 6 | Promoter opening**

a) Structure of SNAPc-containing Pol II PIC bound to U1 promoter in closed promoter complex (CC) state. The subunits are coloured as in Figure 1. The nucleotide residue at the TSS (+1) on the template strand (blue) is represented as spheres (green). The Pol II active site metal ion A is depicted as a magenta sphere.

b) Structure of SNAPc-containing Pol II PIC bound to U1 promoter in open promoter complex (OC) state. The inset represents a zoom into the active center containing open promoter DNA. The catalytic $Mg^{2+}$ ion at the active site is represented as a magenta sphere. The B-reader helix of TFIIB and the charged helix of TFIIF are highlighted alongside the +1 nucleotide residue represented as a sphere (green).

c) Box plots showing TSS precision of protein-coding and snRNA genes (N=18) transcribed by Pol II in cells. Protein-coding genes are sub-grouped based on promoter sequence into TATA-less (TATA−, N=4521) and TATA-containing (TATA+, N=200) subsets. The thickened line represents the median value, the hinges correspond to the first and third quartiles, and the notches extend to 1.58 times the inter-quartile range divided by the square root of N. The whiskers represent the largest or smallest value within the 1.5 times inter-quartile range from the hinge, outliers are shown in black. The precision scores were determined from published 5'cap-seq data[35] (Methods).

d) Annotated transcripts of representative examples from subsets in 6C and genome browser views showing 5'cap-seq signal in the magnified region (± 100 bp) centered at the main TSS peak. The annotated gene region is show below the views and only sense strand signal is shown.

910
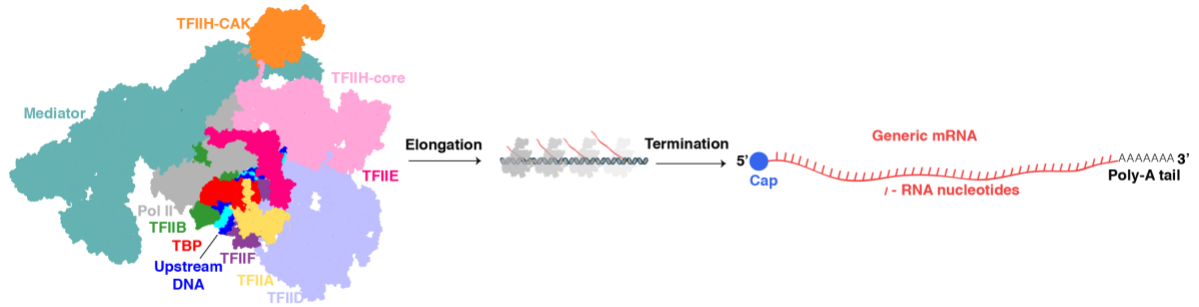911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926

**Figure 7 | Comparisons of Pol II PICs for mRNA and snRNA synthesis.**
a) The Pol II PIC on protein coding genes bound to its elaborate array of initiation factors such as TFIIA, TFIIB, TFIID, TFIIE, TFIIF, TFIIH and Mediator complex.
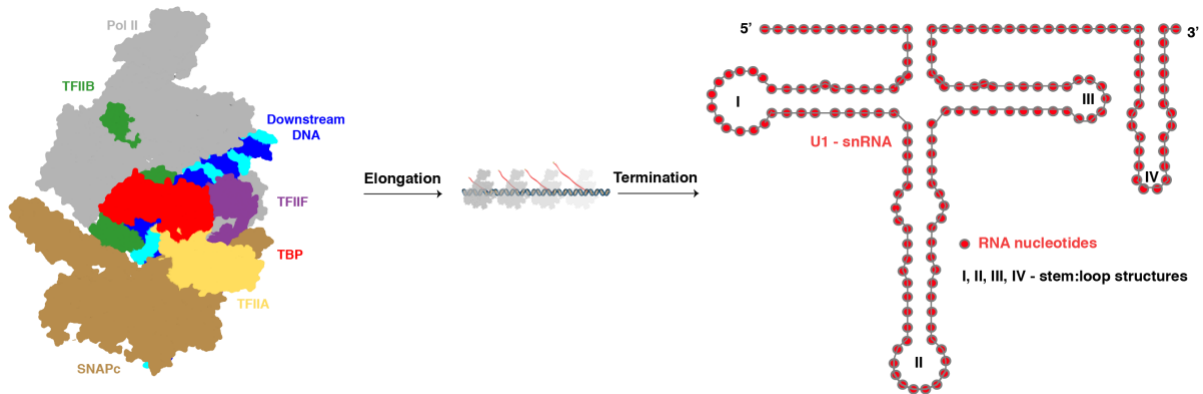b) The Pol II PIC for snRNA transcription requires SNAPc but not TFIIE, TFIIH and Mediator to initiate transcription.
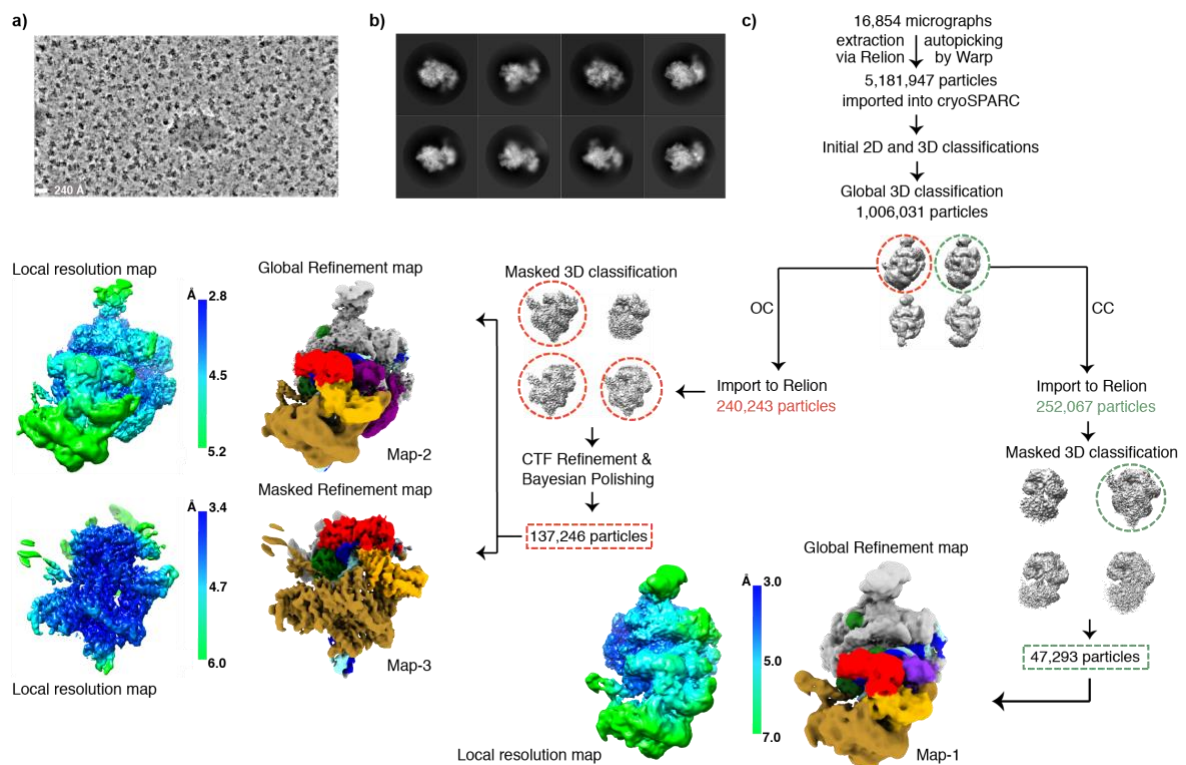
935 **EXTENDED DATA FIGURES**
936
937 **Extended Data Figure 1 | Processing of cryo-EM data for SNAPc-containing Pol II PIC**
938 **bound to U1 promoter. Related to figure 2.**
939 a) Representative cryo-EM micrograph of the SNAPc-containing Pol II PIC bound to U1
940 promoter cryo-EM data collection. Scale bar – 240 Å
941 b) Representative 2D class averages of initially sorted datasets after merging. Adjacent to a
942 well-defined PIC, clear signal for SNAPc is detected.
943 c) Complete processing scheme. After initial clean-up procedures, particles representing
944 SNAPc containing PIC were recovered as two sets. These particle sets were processed
945 separately with respect to the promoter DNA state (CC/OC) and SNAPc occupancy. Final maps
946 are coloured using the subunit color code in Figure 1. The local resolution map indicate the
947 resolution range of final maps (scale bar).
948



949
950
951
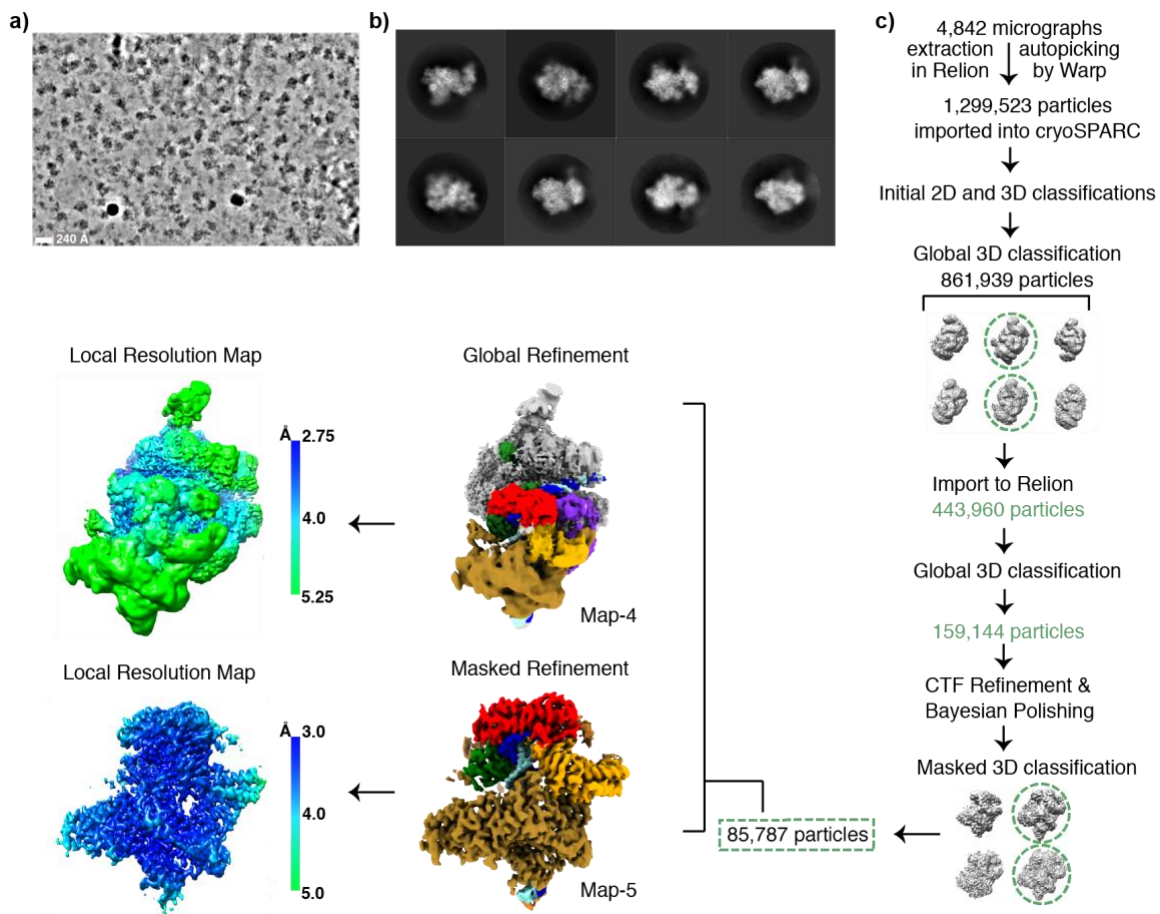952
953
954
955
956
957
958
959
960

961 **Extended Data Figure 2 | Processing of cryo-EM data for SNAPc-containing Pol II PIC**
962 **bound to U5 promoter. Related to Figure 2**
963 a) Representative cryo-EM micrograph of the SNAPc-containing Pol II PIC bound to U5
964 promoter cryo-EM data collection. Scale bar – 240 Å.
965 b) Representative 2D class averages of initially sorted datasets after merging. As in the case of
966 U1 promoter dataset, a clear signal for SNAPc is detected adjacent to a well-defined PIC.
967 c) Complete processing scheme. The optimized strategy from U1 promoter bound SNAPc-PIC
968 dataset was used to obtain high resolution maps of SNAPc-PIC bound to U5 promoter. Final
969 maps are coloured using the subunit color code in Figure 1. The local resolution map of global
970 and locally refined maps indicate the resolution range of final maps (scale bar).
971



972
973
974
975
976
977
978
979
980
981
982
983
984

985 **Extended Data Figure 3 | FSC and angular distribution plot of cryo-EM reconstructions.**
986 **Related to Figure 2**
987 a-e) On the left - FSC plot showing the overall resolution of the reconstructions determined by
988 the gold standard FSC cut-off 0.143, indicated in the graph. In the middle – angular distribution
989 plot of the respective reconstruction showing assignment of particles with respect to various
990 angles. Colour bar indicates number of samples per angular bin (white areas indicate
991 unpopulated angles). On the right - Model-to-map FSCs, showing the fit of modelled structures
992 to their corresponding maps.



993

994 **Extended Data Figure 4 | Structural comparison of TBP bound to TATA containing and**
995 **TATA-less DNA template; Overall of Structure of individual SNAPc subunits. Related to**
996 **Figures 2 and 3**
997 a) Structural super-position of TBP(red) bound TATA-less U1 promoter (cyan/blue) on to TBP
998 (grey) bound to TATA box sequence (PDB: 1YTF)(Tan et al., 1996). The comparison shows
999 that TBP binds to the TATA-less sequence in a canonical fashion and bends the DNA by 90º
1000 b-e) Cartoon representation of the individual structures of SNAPc subunits SNAPC1, 3, 4 and
1001 5 displaying its secondary structure elements as labelled. The N and C termini of all subunits
1002 are indicated.
1003



1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021

1022 **Extended Data Figure 5 | Map quality and map to model fit. Related to Figure 3**
1023 a-h) Sections of cryo-EM density of SNAPc subunits overlaid with their respective atomic
1024 models. Densities are shown as a grey mesh, and sticks are shown for the model as coloured in
1025 Figure 3.
1026 i) cryo-EM density of the TFIIB subunit overlaid to the atomic within the SNAPc containing
1027 Pol II PIC bound to U1 promoter in OC state.
1028 j) cryo-EM density of a region of TFIIF subunit overlaid to the atomic model within the SNAPc
1029 containing Pol II PIC bound to U1 promoter in OC state.
1030 d) Local map of SNAPc containing Pol II PIC bound to U1 promoter in OC state is low pass
1031 filtered to 5Å. The corresponding map is fitted with SNAPc subunits representing map to model
1032 fit, in particular the 'wing-2' region modelled using AlphaFold2[27].
1033
1034
1035



1036
1037
1038
1039
1040
1041
1042
1043
1044
1045

1046 **Extended Data Figure 6: Crosslinking mass-spectrometric analysis of SNAPc containing**
1047 **Pol II PIC. Related to Figures 2 and 3**
1048 a) 2D representation of the overview of BS3 crosslinks. The crosslinks correspond to inter-
1049 protein mono-links that have at least three crosslinked peptide-spectrum matches (CSM). The
1050 subunit colours are consistent with Figure 2.
1051 b) Crosslinks as mapped to SNAPc containing Pol II PIC structure using Xlink analyzer[55]
1052 plugin in UCSF chimera. The inset show the crosslinks observed between SNAPc subunits and
1053 the GTFs' TFIIA and TFIIB respectively.
1054 c) Histogram representing the distribution of Cα pair distances of unique crosslinks mapped to
1055 the structure. Dotted line indicates the 30Å cut-off for BS3 crosslinked Cα pair. A total of
1056 87.8% of the crosslinks were satisfied within this 30 Å cutoff.
1057
1058



1059
1060
1061
1062
1063
1064
1065
1066

32

1067 **Extended Data Figure 7 | Structure based sequence alignment of SNAPc subunits involved**
1068 **in interactions. Related to Figures 3 and 4**
1069 Sequence alignments were performed with the regions of individual subunits for which the
1070 structure has been determined in this study. T-Coffee algorithm[66] was adopted to obtain a
1071 structure based sequence alignment which was then visualized using ESPript[67]. Residues with
1072 identity above 80% are coloured red. Regions involved in interactions are indicated by dashed
1073 boxes and labels.
1074



1075
1076
1077
1078
1079
1080
1081
1082

Rengachari et al.: Structure of SNAPc-containing Pol II PIC

## SNAPC3

```
                                    β1        α1                                    α2
SNAPC3-H.sapiens                    →    ꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁ                            ꞁꞁꞁꞁꞁꞁ
                        30        40        50        60        70        80
SNAPC3-H.sapiens   ...NFPEYELPELNTRAFHVGAFGELWRGRLRGAGDLSLREPPASALPGSQ.AADSDRED
SNAPC3-M.musculus  ...SFPEYELPELHTRVFHVGSFGELWRGRLG.AQDLSLSEPQAAEQPTDGGASNDGFED
SNAPC3-D.rerio     ....VPVYEFVDVNSKEFHIGTFRKLWVDVLN.PEMYSYS.G...........TAPEIED
SNAPC3-X.laevis    MDENIPVYEVTDANTRLIHIGSFGDLWRERLQ.NCDLTLAEK...........DDCLMEN
                                        Interface-2

                                α3                              α4
SNAPC3-H.sapiens   ꞁꞁꞁꞁꞁꞁꞁ   ꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁ                    ꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁ
                        90        100       110       120       130       140
SNAPC3-H.sapiens   AAVARDLDCSLEAAAELRAVCGLDKLKCLE.DGEDPEVIPENTDLVTLGVRKRFLEHREE
SNAPC3-M.musculus  AAVASDLGCSLEAAAELRVVCGLDKLRCLG.EDEDPEVIPENTDLVTLCVRKGLLDYREE
SNAPC3-D.rerio     VELIEEMGIEPAILEELKNICSVDSLRSK...HEDQDIIPSESHLSTLKLRKRRQDYK.E
SNAPC3-X.laevis    SDVAQDLGCSEETAAELRLICGVDALKCSENEEADPENIPEDPSLLTLGIRKKILDRRRE
                        Interface-1                              Interface-4

                                α5                      β2        η1
SNAPC3-H.sapiens   ꞁ   ꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁ                    →        ꞁꞁꞁꞁ       ―
                        150       160       170       180       190       200
SNAPC3-H.sapiens   TITIDRACRQETFVYEMESHAIGKKPENSADMIEEGELILSVNILYPVIFHKHKEHKPYQ
SNAPC3-M.musculus  NITIDRACRQEIFAYEMESHALGKKPENPADMIEEGECILSVNILYPVIFNKHKEHKPYQ
SNAPC3-D.rerio     TLTRDMVDRHEVYANEMEMLSVGKRPDNVRDLIPEGEVILTFNIMYPILFQRFRLVRAFQ
SNAPC3-X.laevis    TLIIERACRQETFLHELEFHAVGKRPEDAADMVDEGELVLTLNIVYPVIFRKHKEYKPYQ
                        Interface-1

                        β3                  η2              α6      β4        β5
SNAPC3-H.sapiens   →                        ꞁꞁꞁꞁ            ꞁꞁꞁꞁ   →         →
                        210       220       230       240       250       260
SNAPC3-H.sapiens   TMLVLGSQKLTQLRDSIRCVSDLQIGGEFSNTPDQAPEHISKDLYKSAFFYFEGTFYNDK
SNAPC3-M.musculus  TMLVLGSQKLTELRDSICCVSDLQIGGEFSNAPDQAPEHISKDLYKSAFFYFEGTFYNDR
SNAPC3-D.rerio     TLHVLGSQKLTDLRDVICCVSDLQVFGEFSNTPDMVPQFISKDHYKSAFFFNGTFYNDT
SNAPC3-X.laevis    TVLVLGSQKLTELRDAINCVSDLQIGGEFSNNPDLAPENICKDLYKSAFFHFEGVFYNDM

                                α7              β6      η3   β7      β8        β9
SNAPC3-H.sapiens   ꞁꞁꞁꞁꞁꞁꞁꞁꞁꞁ              →        ꞁꞁꞁ  →        →         ―
                        270       280       290       300       310       320
SNAPC3-H.sapiens   RYPECRDLSRTIIEWSESHDRGYGKFQTARMEDFTFNDLCIKLGFPYLYCHQGDCEHVIV
SNAPC3-M.musculus  RYPECRDLSRTIIEWSESHDRGYGKFQTARMEDFTFNDLHIKLGFPYLYCHQGDCEHVVV
SNAPC3-D.rerio     RFPECQDISKVIKEWTRSRD..FPDFKTARMEDTSFNDLQMKVGFPYLYTHQGDCEHVVV
SNAPC3-X.laevis    RDPQCRDISRTTIEWAESRDRGYEKFQSAKMEDYTFNDLRLKIGYPYLYCHQGDCEHVIT
                                                                    Interface-4

SNAPC3-H.sapiens   →                       β10            TT        β11       β12
                                            →                        →       →    ꞁꞁ
                        330       340       350       360       370       380
SNAPC3-H.sapiens   ITDIRLVHHDDCLDRTLYPLLIKKHWLWTRKCFVCKMYTARWVTNNDSFAPEDPCFFCDV
SNAPC3-M.musculus  ITDIRLVHHDDCLDRTLYPLLTKKHWLWTRKCFVCKMYTARWVTNNDTFAPEDPCFFCDV
SNAPC3-D.rerio     LTDVRLVHQDDCLDIKLYPLITHKHRVMTRKCSVCHLYISRWITTNDALAPMDPCLFCDQ
SNAPC3-X.laevis    VTDIRLIHHEDCLDRTLYPLLKKRHWFWTRKCSVCLMYTARWVTVNDSLAPDDPCFFCDV
                                                                Interface-3

SNAPC3-H.sapiens   α8              β13
                   ꞁꞁꞁꞁꞁꞁ          →
                        390       400       410
SNAPC3-H.sapiens   CFRMLHYDSEGNKLGEFLAYPYVDPGTFN
SNAPC3-M.musculus  CFRMLHYDSEGNKLGEFLAYPYVDPGTFN
SNAPC3-D.rerio     CFRMFHYDDKGNKVGDFLAYAYVDPGTFN
SNAPC3-X.laevis    CFKMLHYDTDGNKLGEFLAHPYVDPGIFN
```
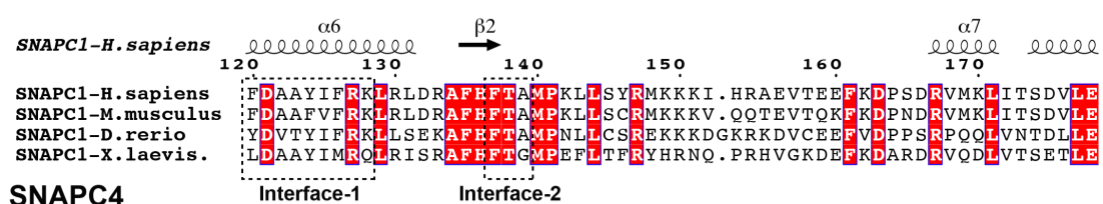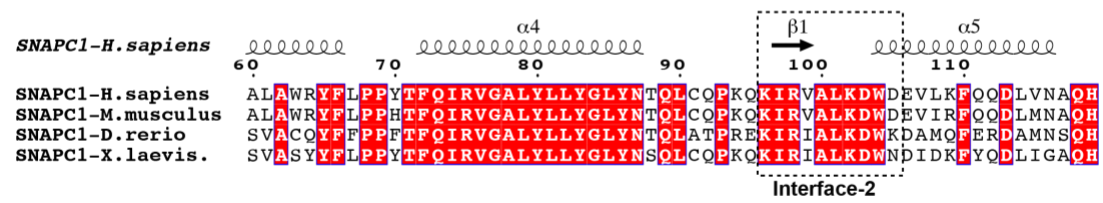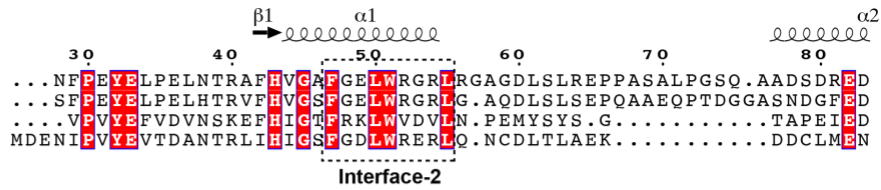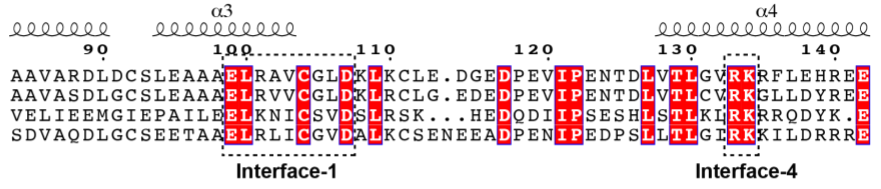
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094

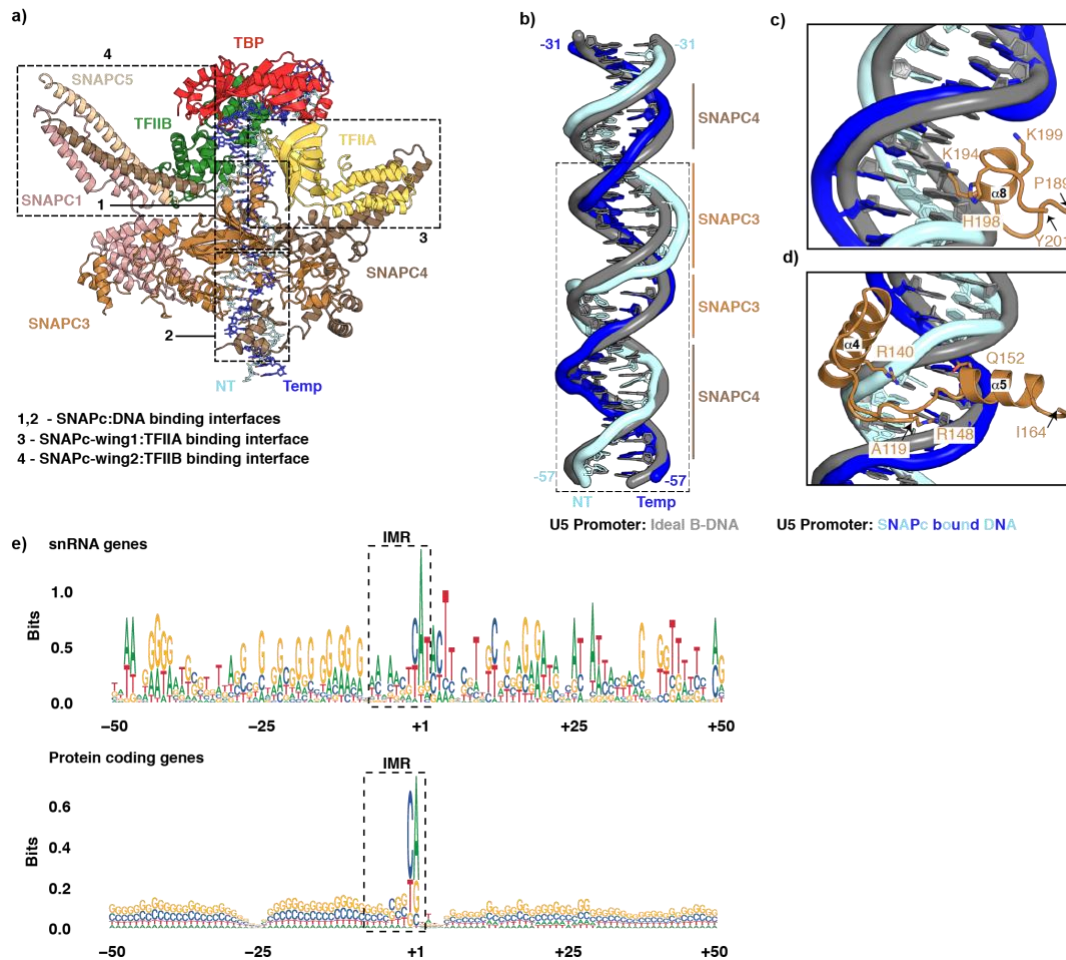1095 **Extended Data Figure 8: Related to Figures 4, 5 and 6**
1096 a) Birds-eye view of the SNAPc interaction with the GTFs' and the PSE motif on U5 snRNA
1097 promoter. The dashed boxes indicate the observed interaction surfaces within the complex (1-
1098 4).
1099 b) Structural super-position of ideal B-DNA of U5 promoter to the SNAPc bound experimental
1100 DNA structure. Major and minor grooves of U5 promoter bound by SNAPC3 and SNAPC4are
1101 labelled and highlighted with lines. Dashed box indicates the PSE region.
1102 c) Close up view of SNAPC3 helix α8 binding to major groove of U5 promoter. The observed
1103 steric clash of K194 with B-DNA highlights the distortion upon SNAPc binding.
1104 d) Close up view of SNAPC3 helices α4, α5 region binding to minor groove of U5 promoter.
1105 The views in panels c and d correspond to Figure 4b, c
1106 e) Sequence logos of DNA sequence surrounding TSS peaks in expressed constitutive
1107 first/single exons for all snRNA genes (n=18) and protein coding genes (n=4721), sorted by
1108 TSS precision scores. The boxes indicate the IMR region (-8 to +2) of promoter flanking the
1109 TSS (+1). While the protein coding genes do not show any enrichment of specific nucleotides,
1110 snRNA genes present a AT-rich profile in the IMR region, indicating its tendency for
1111 spontaneous promoter opening.
1112



1113
1114
1115
1116
1117
1118
1119

1120 **EXTENDED DATA TABLE**

1121

1122 **Extended Data Table 1 | Cryo-EM data collection, refinement and validation statistics.**

| | RNU1-OC (EMDB-xxxx) (PDB xxxx) | RNU1-OC Local map (EMDB-xxxx) | RNU1-CC (EMDB-xxxx) (PDB xxxx) | RNU5-CC (EMDB-xxxx) (PDB xxxx) | RNU5-CC Local map (EMDB-xxxx) |
|---|---|---|---|---|---|
| **Data collection and processing** | | | | | |
| Magnification | | 81,000x | | 81,000x | |
| Voltage (kV) | | 300 | | 300 | |
| Electron exposure (e–/Å$^2$) | | 54.45 | | 51.93 | |
| Defocus range (μm) | | -0.5 to -3.0 | | -0.5 to -2.5 | |
| Pixel size (Å) | | 1.05 | | 1.05 | |
| Micrographs collected | | 16,854 | | 4,842 | |
| Initial particle images (no.) | | 5,181,947 | | 1,299,523 | |
| Final particle images (no.) | 137,246 | 137,246 | 47,293 | 85,787 | 85,787 |
| Map resolution (Å) | 3.0 | 3.5 | 3.4 | 3.0 | 3.2 |
| FSC threshold | 0.143 | 0.143 | 0.143 | 0.143 | 0.143 |
| Map resolution range (Å) | 2.8 – 5.2 | 3.4 – 6.0 | 3.0 – 7.0 | 2.75 – 5.25 | 3.0 – 5.0 |
| **Refinement** | | | | | |
| Initial model used (PDB code) | 7NVU | 7NVU | 7NVS | 7NVS | 7NVS |
| Map sharpening $B$ factor (Å$^2$) | -10 | -10 | -5 | -10 | -10 |
| Model composition | | | | | |
| DNA | 126 | 93 | 132 | 132 | 83 |
| Protein residues | 5842 | 1500 | 5789 | 5789 | 1500 |
| Ligands | 11 | 2 | 12 | 12 | 2 |
| $B$ factors (Å$^2$) | | | | | |
| DNA | 227.04 | 173.83 | 318.34 | 248.19 | 95.70 |
| Protein residues | 111.18 | 185.97 | 215.69 | 124.47 | 109.50 |
| Ligands | 164.08 | 158.45 | 237.49 | 194.42 | 82.54 |
| R.m.s. deviations | | | | | |
| Bond lengths (Å) | 0.005 | 0.003 | 0.004 | 0.007 | 0.003 |
| Bond angles (°) | 0.756 | 0.606 | 0.509 | 0.658 | 0.524 |
| Validation | | | | | |
| MolProbity score | 1.77 | 1.72 | 1.69 | 1.63 | 1.67 |
| Clashscore | 9.25 | 10.84 | 9.85 | 7.55 | 7.63 |
| Poor rotamers (%) | 0.16 | 0.00 | 0.00 | 0.00 | 0.00 |
| CaBLAM outliers | 1.93 | 1.73 | 1.63 | 1.91 | 2.07 |
| Cβ outliers | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Ramachandran plot | | | | | |
| Favored (%) | 95.91 | 97.01 | 97.03 | 96.68 | 96.27 |
| Allowed (%) | 4.09 | 2.99 | 2.97 | 3.32 | 3.73 |
| Disallowed (%) | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |

1123

1124

1125