
An enlarged largest subunit of *Plasmodium falciparum* RNA polymerase II defines conserved and variable RNA polymerase domains

Wu-Bo Li, David J.Bzik, Haoming Gu, Manami Tanaka, Barbara A.Fox and Joseph Inselburg

Department of Microbiology, Dartmouth Medical School, Hanover, NH 03756, USA

Received September 11, 1989; Revised and Accepted November 1, 1989

EMBL accession no. X16561

ABSTRACT

We have isolated the gene encoding the largest subunit of RNA polymerase II from *Plasmodium falciparum*. The RPII gene is expressed in the asexual erythrocytic stages of the parasite as a 9 kb mRNA, and is present as a single copy gene located on chromosome 3. The *P. falciparum* RPII subunit is the largest (2452 amino acids) eukaryotic RPII subunit, and it contains enlarged variable regions that clearly separate and define five conserved regions of the eukaryotic RPII largest subunits. A distinctive carboxyl-terminal domain contains a short highly conserved heptapeptide repeat domain which is bounded on its 5' side by a highly diverged heptapeptide repeat domain, and is bounded on its 3' side by a long carboxyl-terminal extension.

INTRODUCTION

Transcription in eukaryotes is directed by three classes of nuclear RNA polymerases (1). The genes encoding the largest subunits of eukaryotic RNA polymerases I (RPI) (2), II (RPII) (3–8), and III (RPIII) (3) have been isolated and are single copy genes with the exception of *Trypanosoma brucei* RPII, which contains two alleles (7,8). RPII is a complex enzyme consisting of multiple subunits, and is responsible for the transcription of protein coding genes. DNA sequence analysis of the cloned largest RPII subunit from different species (3–8) revealed nucleotide and amino acid homology with the largest subunit of *E. coli* RNA polymerase (β') (9), yeast RPI (2), and yeast RPIII (3). The largest subunit of eukaryotic RPII (to be subsequently referred to as RPII subunit) contains a repeated heptapeptide sequence in its carboxyl-terminal domain (CTD), which is repeated 52 times in mouse (4), 42 times in *Drosophila* (6), 27 times in yeast (3), and is absent only in the *T. brucei* RPII subunit (7,8). The CTD structure is specific for the RPII subunit and is essential for cell viability (5,10). The CTD may function in initiation of transcription by mediating phosphorylation and destabilization of histone and DNA interactions, thus facilitating transcription through nucleosomes (3,4,11,12,13). Alternatively, the CTD may function as a receptor for essential transcription factors, or to anchor the RNA polymerase II to a structure within the nucleus (14).

Little is known about transcription in *P. falciparum*, or its regulation. RNA synthesis in *P. falciparum*, during the 48 hour intraerythrocytic growth cycle *in vitro* (15), starts within 6 hr after infection of red blood cells (RBC), rapidly rises during the early trophozoite stage and peaks during the schizont stage (16). Many species of the parasite's proteins, including most of the known parasite antigens, are synthesized in synchronously infected cultures during the trophozoite and schizont stages [25–42 hr after infection of RBC] (16).

We are initiating a molecular approach to the study of transcription in *P. falciparum* with this report of the isolation, sequence, and structure of the gene that encodes the RPII

subunit of *P. falciparum*. This is the first characterization of a transcription factor for *P. falciparum*. We describe the determination of the RPII gene copy number, its chromosomal location, its expression during erythrocytic growth, and identify conserved regions of the protein that are separated by novel enlarged variable regions.

MATERIALS AND METHODS

Cells, DNA, and vectors

P. falciparum strains FCR3 (17) and Honduras-1 (18) were used in the present studies. Parasites were grown in RPMI 1640 medium (15,18) prepared with HEPES buffer (pH 7.2), 0.2% sodium bicarbonate, 10% heat-inactivated human type A, Rh⁺ fresh frozen plasma, penicillin (100 IU/ml), streptomycin (100 mg/ml), and gentamycin (20 mg/ml). α -Amanitin was purchased from Boehringer Mannheim Biochemicals. Honduras-1 DNA was isolated from asynchronous cultures of parasites by lysis in a GuITC solution (4 M guanidium isothiocyanate, 5 mM sodium citrate (pH 7.0), 0.1 M β -mercaptoethanol, and 0.5% Sarkosyl) and centrifugation to separate the RNA by pelleting. The DNA was collected and was purified by banding in CsCl (19). The cloning vectors λ gt11, λ ZAP, and pBluescript were obtained from Stratagene, Inc.

Oligonucleotides

Four oligonucleotide probes, A, B, C, and D were synthesized and used in hybridization experiments. The sequence and location of the probes were as follows: probe A (a mixture containing 16,384 35-mers) 5'-GGA(orT)C(orG)AATAA(orT)G(orC)AA(orT)GGA(orT)G(orC)AA(orT)GTA(orT)GGA(orT)G(orC)AATAA(orT)G(orC)AA(orT)GG (based on the consensus heptapeptide repeat region); probe B (a 30-mer) 5'-ACCTCTATCATTAGCAA CTAAGTGTCTTC (nucleotide 4860 to 4831); probe C (a 30-mer) 5'-TCC-AAACAAAGCCGATGTATTAGAATCACC (nucleotide 5190 to 5161); and probe D (a 22 mer) 5'-TTCTACAGAATGATCACATGCT (nucleotide 462 to 441). The oligonucleotides were treated in concentrated NH₄OH at 55°C for 5 hr, lyophilized, and purified by NENSORB™ PREP (DuPont). The oligonucleotides were end-labeled according to standard procedures (19).

Isolation of cDNA and genomic DNA clones

The cDNA library in λ gt11 was constructed from trophozoite and schizont stage Honduras-1 mRNA (20). Honduras-1 genomic EcoRI and XbaI libraries were constructed in λ gt11 and λ ZAP as previously described (21). The DraI genomic DNA library was constructed in pBluescript. Honduras-1 genomic DNA was digested to completion with DraI, electrophoresed on a 1.0% agarose gel, and DNA of 400 to 1000 bp was collected and purified (22). The DraI digested size-fractionated DNA was blunt-end ligated into SmaI digested pBluescript. End-labeled probe A was hybridized to replica filters of λ phage plaques in 6 \times SSC at 42°C and washed in 1 \times SSC at 45°C. All of the clones isolated, C1, C3, and C8, gave very strong hybridization signals. A specific 169 bp DpnI restriction fragment from the 5' part of cDNA clone, C8, was cut out of a 1.0% agarose gel and oligo-labeled to a specific activity of 1 \times 10⁹ cpm/ μ g (23). Filters of cDNA phage plaques were hybridized with the 169 bp DpnI probe in 6 \times SSC at 42°C and washed in 1 \times SSC at 60°C. All of the 14 cDNA clones isolated gave strong hybridization signals and were overlapping cDNA clones. Probe B and C were end-labeled and used to screen filters of genomic DNA containing phage plaques by hybridization in 6 \times SSC at 42°C and washing in 1 \times SSC at 50°C. Probe D was end-labeled and used to screen filters containing colonies from the DraI genomic DNA library by hybridization in 6 \times SSC at 37°C and washing

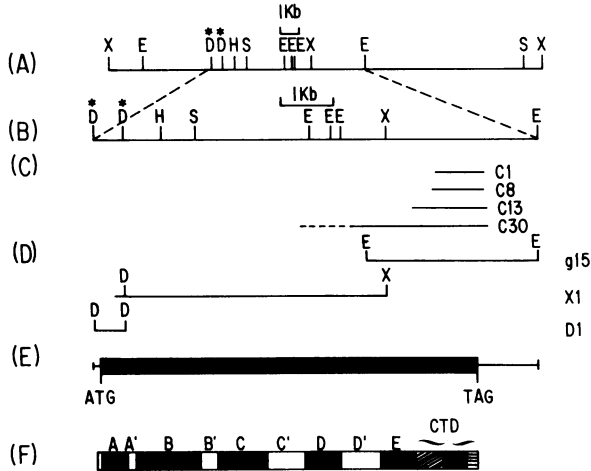


Figure 1. Restriction maps, cDNA, and genomic DNA clones of the *P. falciparum* RPII gene. Restriction sites shown (A) are E = EcoRI, H = HindIII, S = SpeI, X = XbaI, and D = DraI. Only two of the twenty DraI sites are shown (*). Enlarged restriction map encompassing the RPII subunit gene (B). Restriction sites are the same as in (A). Location of cDNA clones C1, C8, C13, and C30 (C). The dashed line on cDNA C30 indicates a region of sequence non-collinearity with the RPII gene due to a double ligation. Location of genomic DNA clones g15, X1, and D1 (D). Location of the RPII subunit open reading frame (E). The ATG start, and TAG termination codons are shown. Domains of the RPII subunit (F). The five conserved regions, A, B, C, D, and E, of the RPII subunit are indicated as solid blocks. The enlarged variable regions, A', B', C', and D', are indicated as open blocks. The CTD is represented as three domains (see text), where the middle region (solid block) represents the heptapeptide repeat domain.

in $1 \times$ SSC at 45°C . All of the clones isolated with probes A, B, C, D, and the 169 bp DpnI fragment were determined to be RPII gene clones by DNA sequence analysis.

DNA sequencing and computer analysis

The cDNA and genomic DNA clones were sequenced as previously described (24). Both DNA strands were completely sequenced using the dideoxy chain termination method (25). The complete DNA sequences were established using the DNA Inspector II programs (Textco Inc., Lebanon, NH). Amino acid sequence comparisons with other RNA polymerase subunits was done by a dot matrix analysis using the MacGene Plus™ program. The prediction of secondary structure of the RPII subunits was performed using the IBI-Pustell programs (International Biotechnologies Incorporated), and the MacGene Plus™ programs.

Northern blot analysis

Total *P. falciparum* Honduras-1 RNA was prepared as previously described (20) from synchronized cultures of trophozoite and schizont stage parasites. Poly(A⁺) RNA was purified from total RNA using an oligo d(T)-cellulose column (19). Total poly(A⁺) RNA [mRNA] (4 μg per lane) was electrophoresed in a 1.2% agarose-formaldehyde gel, blotted to Zetabind nylon membrane (CUNO, Inc), and hybridized to oligo-labeled cDNA C1 as previously described (24,26). Washing was done in $0.1 \times$ SSC at 65°C . cDNA C1 is specific for the RPII gene.

Pulsed field gradient (PFG) electrophoresis of chromosomes

PFG electrophoresis (27) was performed using a contour clamped homogeneous electric

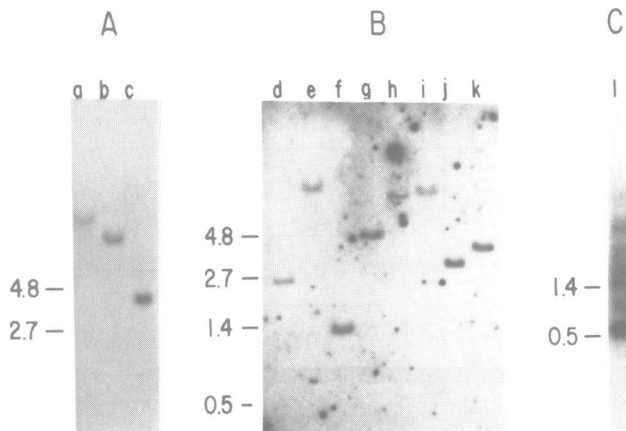


Figure 2. Genomic southern analysis of the *P. falciparum* RPII gene. *P. falciparum* genomic DNA (5 μ g per lane) was digested with restriction enzymes SpeI (a), XbaI (b and e), EcoRI (c), BclI (d), XbaI/BclI (f), XbaI/HindIII (g), XbaI/EcoRI (h), XbaI/ClaI (i), XbaI/BglII (j), XbaI/SpeI (k), and DraI (l). The digested DNA was electrophoresed on a 1% agarose gel, blotted to nylon, and hybridized with probe A (A), a mixture of probe B and C (B), or probe D (C). Each probe was hybridized in 6 \times SSC at 42 $^{\circ}$ C, and was washed in 1 \times SSC at 45 $^{\circ}$ C (probe A), 50 $^{\circ}$ C (probe B and C), or 45 $^{\circ}$ C (probe D). The sizes of the molecular weight markers are indicated in kb.

field (CHEF) apparatus (28). Agarose blocks containing parasite chromosomes were inserted into a low melting point agarose gel (0.6%) as described (29). Briefly, the electrophoresis was done in 0.75 \times TBE (67 mM Tris, 67 mM boric acid, and 1.5 mM EDTA, pH 8.0) at 14 $^{\circ}$ C for 48 hr using a 3 min pulse time, and an 8 V/cm field strength. Separated chromosomes were transferred to a Zeta probe nylon membrane (BioRad), and hybridized with the RPII gene specific probe, cDNA C1, in 6 \times SSC at 42 $^{\circ}$ C, and was washed in 1 \times SSC at 65 $^{\circ}$ C.

RESULTS AND DISCUSSION

Cloning and sequencing of the P. falciparum RPII gene

A 35-base oligonucleotide mixture (probe A) was synthesized, based on the consensus carboxyl-terminal heptapeptide repeat sequence of the eukaryotic RPII subunits, and used to detect the *P. falciparum* RPII gene. A cDNA library was screened with probe A and clones, C1, C8 and C13 were isolated and sequenced. C1, C8, and C13 cDNA inserts were 800 bp, 1 kb and 820 bp, respectively (Fig. 1C.), and contained overlapping DNA sequences. Fourteen additional cDNA clones were selected with a 169 bp DpnI fragment derived from the 5' region of cDNA C8. The largest cDNA clone isolated, C30, contained a 3.5 kb cDNA insert (Fig. 1C.).

The difficulty of isolating large cDNA inserts compelled us to begin selecting clones from genomic DNA libraries. This approach is reasonable in *P. falciparum* because relatively few of the characterized protein coding genes contain introns. In addition, *P. falciparum* introns are always short sequences (107 to 430 bp) that are easily identified by their characteristic structure(s) (24). A unique 3.8 kb genomic EcoRI fragment hybridized with probe A (Fig. 2A.). Six λ gt11 genomic EcoRI clones, including clone g15 (Fig. 1D.), were selected with probe A and contained the same 3.8 kb EcoRI insert. The DNA

sequences of clone g15 and the cDNA clones were colinear (Fig. 1C. and 1D.). Oligonucleotide probes B and C, derived from the 5' region of clone g15, hybridized with a 12 kb genomic XbaI fragment (Fig. 2B.). Probe B and C selected a clone, X1, from a XbaI genomic DNA library. Clone X1 contained the 3', 5.3 kb, sequence of the original 12 kb fragment, and had deleted the 5', 6.7 kb (Fig. 1D.) [Gene flanking sequences of *P. falciparum* commonly delete in recombination deficient *E. coli* (21).]. Oligonucleotide probe D, derived from the 5' region of clone X1, identified a 580 bp genomic DraI fragment (Fig. 2C.). Probe D selected clone D1, which contained 580 bp (Fig. 1D.), from a DraI genomic DNA library.

The restriction maps of HindIII, EcoRI, SpeI, and XbaI for the genomic DNA clones spanning clones D1, X1, and g15 are shown (see Fig. 1A. and 1B.). The overlapping genomic DNA clones (g15, X1, and D1) were sequenced and collectively represented 8631 bp (Fig. 3.) [The 5' 80 bp of clone D1 are not shown]. A large open reading frame (Fig. 1E.) begins with the ATG at bp 1 and ends at the TAG at bp 7357 (Fig. 3.), and the deduced protein contains 2452 amino acids. The distribution of A + T content was typical for a *P. falciparum* coding region (30), as the 5' and 3' flanking sequences contained higher A + T content [90% for the 5' flanking region (data not shown) and 83% for the 3' region] than the region encoding the long open reading frame (72%). Introns were not considered to be present in the long open reading frame based on the observations that all characterized *P. falciparum* introns (i) have a minimum A + T content of 85%, (ii) have no intron-length open reading frame, (iii) are between 107 to 430 bp long, and (iv) have a highly conserved 5' and 3' junction sequence (24).

Amino acid sequence comparisons with other RNA polymerase subunits

Amino acid comparisons of the 2452 amino acid *P. falciparum* protein with various RNA polymerase largest subunits from different species was done by a dot matrix analysis (Fig. 4.). Because the two RPII subunits of *T. brucei*, RPIIA and RPIIB, differ by only 4 amino acid substitutions over 1765 amino acids (7), only the RPIIA subunit was compared. The homology of the 2452 amino acid *P. falciparum* protein to the mouse, *Drosophila*, yeast, and *T. brucei* RPII subunits was significantly higher than the homology to *E. coli* RNA polymerase (β' subunit), and the largest subunit of yeast RPI and RPIII. This observation, along with the presence of the heptapeptide repeat in the CTD of the *P. falciparum* protein (Fig. 3.), shows that the gene we characterized is a form of the *P. falciparum* RPII subunit. The homologous sequences of the *P. falciparum* RPII subunit and the other RPII subunits were partitioned into only 5 colinear regions (A through E), followed by the CTD. The heptapeptide repeat is absent in both RPII genes of *T. brucei* (7,8). Each conserved region (A through E) of the other RPII subunits was shifted to the carboxyl-terminus (right) when compared with the *P. falciparum* RPII subunit. The shifts show that the *P. falciparum* RPII subunit is the largest characterized RPII subunit, and that it contains an enlarged variable region [regions A' to D' (Fig. 1F. and Fig. 4.)] between each of the conserved regions. The longest isolated cDNA clone, C30, included most of variable region D' (Fig. 1C), demonstrating that this variable region is expressed. The presence of the enlarged variable regions helps to clearly define the conserved regions of the eukaryotic RPII subunits. Previous comparisons indicated that the RPII subunits could be divided into 6, 7, or 8 conserved domains (6,7,31), depending on the specific criteria used for the interpretation of those alignments.

The five conserved regions (A to E) consisted of 148, 416, 312, 239, and 195 amino acids, respectively (Fig. 5.). The conserved regions (total of 1310 amino acids) of the

Nucleic Acids Research

ATGACGGTTGATTGTAATTCATTCATTCAGCATTCCAGTAATTAAGAGACTAAAACGATAGTAGCTAGGCTGTTTGGATCCGAAATAAATAAGGATAGCTCTTTCGAAATGTAAT 120
M T V D L X N I P Y S A C E L K A G R V K R L E L G V L D P E I L V X N 40

GTAGATATATATAAAGATGTTTTCCAAGAGAAAGGTTGAAATTAATGATATACGTATGGCTACTATTGATTATAGGACCTTATGTGGTACATGTAATGAATGTAATAATTCTCTGGT 240
V D I Y K K D G F P R E G G L N D I R M G T I D Y R T L C G T C N M X V K Y C P G 80

CATTTGGTCATATAGAATTAGCGAAACCTATGTCATTATGTTTGAATGTGATTAATAATGTTTGAAGATGTGATGTATCATGTGGTGTGATATTATGTAATGTGGAACAT 360
H F F G H I E L A K P M Y H Y G F M N V V L N V L R C V C Y H C G R L L C X V N S 120

TCTAAGTTAAATATTTCAAAGATTAAGTAATAGTTTAAAGTACGAAATTAAGCTGAACCTGTGTTAGGATAAAGACATGTGATCATCTGTAGAGAAGAAGGATTAATATTT 480
S K V Y E K I E K I K V N S L R L R K L A E L C L G I R A C D H S V E E G L N I 160

AACGATAATCTTTAAATAATTTTATAACAATGATTTAAGTAAATTAATATGAATCAACAATGCTTTTAAATAAAGCTAATTAACGAACATATTTGAATGGTACTAAGAAGAT 600
N D N S L N N F Y R N N D L S N L N M N Q S M L L N K S N Y N N T P N I F E M V S K E D 200

GTAGCTTGGATGCTCAACAAAATATAGTACAGAAGACCAAAATGATATATCAATTTTACATAGTACTGAGAAGATATTGATGAGATAAAAGAAAATTAAGTCTGAAGAA 720
V D C G C V Q P K E V S R E G P N M Y I Q F L H S S E E D I D E S K R K L S A E E 240

CCATTAGAAATATTAAGAAGAAAAGAAAGAAAGTATTAAGGATTTAATCTGATAGGCTGTACCAGCTTTTAAATATTAACATGTATACCTATACCTCCACCATGTGGC 840
A L E I L K K I R K E E H S I L G F N S D R C V P A S L I L T C I P I P P P C A 280

AGACCTTATGTTCAATATGGAATCAAGAACTGAAGTATTAACCTTAAATATTAGATATCTAAAAGAAATACACATTTAAAAGGCAACCGGATGACGAGCAAAATCAAT 960
R P Y V Q Y G N Q R S E D D L T L K L L D I V I T T N I Q L K R Q T D R C A K S H 320

GTATTACAGGATTTATGTTCCCTTACAATTTCAATAAAGTCTTTGATAATGATTTCCAGGATGCGGACACCAACAGCATCAAGAACCTATAAAGCTATAAAGCA 1080
V L Q D L L C S L Q F H I T T L F D N D I P G M P I A T T R S K K P I K A I R T 360

AGCTTAAAGGTAAAGAAAGCACTAAGAGCTAATTTGATGGTAAAAGAGTGGACTTTTCAAGCAAGACCGTTATTCAGGAGATCAAAATTTAATGATTATGATAGGTGCTCT 1200
R L K G K E R L R K R D F S A R T V I T T G D P N L N I F E M I D Y I G V 400

AAATCGGTAGCTATGACATTAACATTTTGTGAGACAGTAAACCTTTTAAATATGATAATTTAAAGAAGCTGTAGAAAGGGTCCATTAATGATGGCTGGAGCAAAAATATTATTAGA 1320
K S V A H M T L T F C E T V T P L N Y D N L K K L V E R G P Y E E W P G A K Y I I R 440

GATAAGTGTACAAATATGATTAAGACATGTACGAAGAAGTATGAGATATGAGATATAAGTATAAAGTATAAAGTATGACCGATGAAATTTAATTTTAAACAGACAG 1440
D N G T K Y I D L R H V R R N S E K E L E Y G Y K V E R H M T D E D Y I L F N R Q 480

CCTTCATACATAAAGTATTAAGGCTAAGGCAAAAATTAACCTTATCAACATTTCCGTTTAAATTTACGTCACCTTCCCGGTAATGCTGATTTGATGGAGCAAGAAATG 1560
P S L H K H S I M G H K A K I L P Y S T F R L N L S V T S P Y N N N N N N N N N N 520

AACCTACATTAAGCTCAGTCACATGAACAAGACTGTGATTAACATTTAATGATAGTACAAGCAAAATGTTTACCACCAAGCTAAACACCAATTTGGGATAGTACAGATTTCC 1680
N I K H L A K S H E A K T E S I K H L I V Q R Q I V S P Q C N K P V M G I G D S 560

TTTATAGCTATAAGAAAATTTACAGAAGACATAATTTCCCTACAAAAGAAAGTATGTCCTTTAATTTGGATTCATATGGAATCAATGTATACCAACACCAGCAATAAATAAAA 1800
L L L A I R K F T R R D N L F L T K E E V H S L L I W I P Y V N N H V I P T P A I I K 600

CCAAGGATATTTGGACAGAAATAAAATTTTTCGATTTTACAATTTGATGATATAGAAGATATAAAGTATAAAGTATAAAGTATAAAGTATAAAGTATAAAGTATAAAG 1920
P R A L L W T G K Q I F S M L L Q F D D M N I E D D K N D T A N N K N V G R D V N T 240

AATGCTAAAGCAACTAGTCAAAATGAATGACTGTGCTAATTTATGCTAATGCTAATGCTAATGCTAATGCTAATGCTAATGCTAATGCTAATGCTAATGCTAATGCTAAT 2640
N V N K D S S K A H N T S G N Y Y Y G N S T N D N T D D D Y L E K G N A Y S R S G N 680

AATCACTAATAGCTCTTATCTATGGGATAATAAAGTACGAAATGACAGCAAAATGATAGCTCCCAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAA 2160
N H P N S P L S I G D N N V P N Q N D M H S G T N N N N N N N N N N N N N N 720

AATAAATAAATAAATAAATGCGGAGTTAATGATTTAAACGTTTTAATGATGCTAAAATAAATAAATAAAGAGAGCTTCAACATGATCAAGATGATATCCATATGCTTCAAT 2280
N N N N N I G G G I N S K E T R S F N M V K I N L M R D S T S K C K D D N P Y C S I 760

AATGATGCTAAGGTTATAAATAAATAAAGCAATTAATTAAGTGTATCATATGTAAGAAGCTGTGGTCTCTTGTAGGTGCTGTTAATTCATGTTTTATGGCATGAAATGGGTCAGAT 2400
N D G K Y I I X N N E L L S G I C K R T V G S S G S S L I V L W H E M G G P D 800

AAAACGAAGATTTTATGCACTTCAAAAAGTACAATAATGCGTTGAAATGTTGCTTATGAGTATGTTGCTTATGCAAGTATTTGGCAAGTAAAGTATGGCCAGGTCGCA 2520
K T K D F L T S A L Q K V T N N W L E Y V G F T V S C S D I A S N K V L G K V R 840

GAATATGATGAAATCTAAAGTGAAGCTCAAAATCTGTGAAGAAGCAGAAAGGGGAATGAAATGAGTGTACGCCAGGCAAAATGATTAATGATGATTTGAACTAGAGTAAAT 2640
E I L D E K S K S E V S K L V E K A Q K G E L E C Q P G K S L Y E S F E T R V N M 880

GAATTAATGCTGCTAGCAAGAAATGGCTGAAAAGTGCATCTGAGACTTTAGATGAAAGAAATAAATTTTACTATGGGCTAGTGGGCTAAAAGCTTCTAATTAATAATATCCAA 2760
E L L A N C A R E H A G K V A S E S L D E R N N I F S M V A V S G S K G S I I N I S Q 920

ATTAATCATGCTAGCTCAAGAAATGGAAGCAAAATGCAATTTGCTTATATATGATGATCTTACCTCAATTTAATAATTTGATTTGCTGCTGAGTACAGAGTTGCTA 2880
I I S C V G Q Q N T V E G K R I P F G F N H R S L P H I T F F G G P E S R G C F V 960

TCAAATCTTTAATTAAGTGAATAACACCACAAGAAGTATTTTCCATGCTATGGGAGTACAGAAAGTATTAATGATGATGATGATGATGATGATGATGATGATGATGATGATGAT 3000
S N S Y L D L S G L T P Q E V F F H A M G G R E G I D T A C K T S T E T G Y I Q R 1000

TAAATAAAGCCATGGAGAGCTTTATGGTCAATATGATAGCACTGTAAGAAATGATATGAGATATTAATGAAATTTTGTATGGAGAAGTGTGCTGCGGATATATAGAACAT 3120
L I X A M E D V M V Q Y D R T V R N S Y G D I I Q F L Y G E D G M A G E Y I E D 1040

CAATATGATGATTAAGAAATGATATAAAGAGATTAATAAATATATAAATAAATTTGATGAGAAGCAATTTGGAAGGATTTATTAATGATGATAAATAAATAAATAAATGCTGATGATA 3240
Q I I D L H K K L D N K E I N K L Y K Y N F D E E P F A T T K D Y I I G N K I K A D G S R 1080

AATCTAGTATATAGATTAATAAAGCAAAATTTTAAATCAAGAAATTTGAAGATATATAAAGTAAATAATTTATGTAAGAAGAAATTTCCAGATGGAGATTAAGCAACAT 3360
N T T Y I D Y N Q K N I L N Q E F E E L Y K C K E I F P D G D I R Q H 1120

TTCCAATATATGATGATGATTAATGAAATGCAAAATTTCCATCTATACCATTTGTAAGTAAATAAATAACTACAACAATAAATAAATAAATAAATAAATAAATAAATAA 3480
L P I N H N R L I E Y A K S C A I P F V S N N N S T N N N N N N N N N N N I 1160

AGTAATGCTAGAAAATTTGATGAAAGGTAATTTATCGCTACACATAATCATAAGGAAATAAAGAAGCAAGCAAGCAAGCAAGCAAGCAAGCAAGCAAGCAAGCAAGCAAG 3600
S N S R K L H M D K G N L S T H N H K E N K K R R K R R R K R R R K R R R K E E 1200

AATAAGCACTTATCTCAATTAAGAAGGATGATAAATAATGATCTTAAATAATGATGATGATGATGATGATGATGATGATGATGATGATGATGATGATGATGATGATGAT 3720
N N E L M S E I K K E Y E N N D L N N M H I S K G D Q S P F K G M H E F H M G V 1240

GCAGATATGAT 3840
A D N D H G S D L G N N N N N N N N N D F V D D D Y D N D D D D Y D D D D 1280

TATGAT 3960
Y D D D D L D D D D L G N S D N I N I G G N R K Y Y G N T L K N K Y I N I N I N I N I N I N I N I N I N I N I N I N I N I N I N I N I N I N I N 1320

CCAATGAT 4080
P I D V V H K V L E N L V I E K L V I I K Q I N S N D T L S V E A Q N N A T I L K 1360

GCACATTAAGAATTTTGAATTTCAAACCTTTAAGCTCAAACCTATAAAGTTAGTTTAAAGGATAGATGCTTATTACAAGAAATAGAAAAAATTTTAAATACCTTATGCTCAT 4200
A H L R T Y L N S K G L L T Q T H K V S V K G L D W L L Q G G I E K I F Y K S L C H 1400

CCAGCAGAAATGCTAGCAGCTGCTGAATCAATGGGAGCTGCAAGCTCAGATGACATTTACCTTCCGCTGAGTTTCAAAAATTTGATGATGATGATGATGATGATGATGAT 4320
P G E C V G A L A Q S I G E P A T Q M T L N T F H F A G V G S K N V T L G V P 1440

AGATTAAGAAATTAATAAATAGTAAAAAATGAAAGACTCCATCAACAACAATATATTAGATGATCGTTTCAATGATCAACAACAAAGCTAAAGATATTTAAACAAAATAGAA 4440
R L K E L I N I V K N V X T P S T T I Y L D D M V S N D Q Q K A K D I L T K L E 1480

TATACTACATTTGAACAATTAACCTCAGTGCACAAATTTATGTAGTCTTAACAACAACAATTTTGGAGGAAGATAAATCGTGGTAAATGAATTTATGAATTTCCGACGATGAG 4560
Y T T L K Q L T L S H A Q I I T D P N T T P T I L E E D K M S W N E F Y E F P D E 1520

GATGATCACTAAATTTAGTGAATGGGTTAAAGAAATCAAAATCAACATATACATGTAATGAAAAAAAATTAAGTAAAGAAAATTTGTTATATATATCTGTCTTTTCA 4680
D D T Q Y S L G E W V L R I Q L T N I H V N E K K L T M K E I V Y I I Y S V F S 1560

ACTGATGAATAGATATATATACAGATGTAACCTCAGAAAGTACTTTTAAAGATTCGGGTGAAATTTAAAGTCTGGAATATTTTGAATTTATCTGCTGATAGATAGCT 4800
S D E L D D I I Y T T D N N S E D L V L R I E R K T L N G E Y T F M N Y D V V D G N A 1600

AATGAACAAGTTGATGAACAAGAAGAGTGAACAACACTTACTGCTAATGATAGAGGTAATACGATGAACAACAAAATACTACTCATCTCATGATTTAATAACAATACTACA 4920
N E Q V D E Q E E D E E H L V A N D R G N Y D E T K N S T H P H H D Y N N N T T 1640

AATATATTTAAGTCAAGTAAAAAATAATATATCATCAGATATAAATGACAAAGATGAGAGTACTTACTGATAAATGACTAACAATGACAACTAAAAAATTTAATTCATCACC 5040
N I F K S K V R K N N I S S D I N T K N E D S I S I N S S N N E Q V K N I N S S P 1680

CTTTCAATAATATGCAATAATAATAATAATAATAATGACTAGCAATATTAATGATTTAAGTGAAGAATATAAAAAAGAAGATGCAAAATGAAGTCCATTAAGAGCGGT 5160
V S N N H H N N N N H N N N D S S N I N D I K V K N I K E E D G K V N E F Y E F P D E 1720

GGTGAATCTACATCGGCTTTGTTGCAAAATAAAAATGCAAAAAGAGATAAATTCGTGAACAATAATGATAAATGATGATGATGATGATGATGATGATGATGATGATGATGAT 5280
G D S N T S A L F T G N A L K N S Q K E D N I V N N N D N N D D D D E E E E E E E 1760

TTGTTTGGTACCATAATGTATCTCCAAAAATACGAAGATGCAAAAAATAAGAATCAACAACAACAAGTAAATAATAAGAAAACAAAAAAGCGGAATAATAATAGTAA 5400
L F G D H N V S P K N T X D G K N K N T H N K S N N H E N K N K S G M N N S N 1800

AATAGTAATAGCTGATGATGCTGAT 5520
N S N T T A A T D D G D V D D N D N D D D D N K S D I T I K E D N D V A F M K T S T 1840

AAAAATGAGAAAGATTTAAGCAATTAAGAATTAAGCAATTAAGCAATTAAGCAATTAAGCAATTAAGCAATTAAGCAATTAAGCAATTAAGCAATTAAGCAATTAAG 5640
K N A E E E D L E L E N K N H I E N I S E E D T E D P F L E K L M E Q C L S T L 1880

AAATTAAGACCTTTGAAATTAAGTAAATATGACAGACCAATCCAAAAATCAATCAATCAATCAATCAATCAATCAATCAATCAATCAATCAATCAATCAATCAAT 5760
K L R G I E N I T K V Y H R E E S K I T Y Y D S D N H G K F V R S N H W V L D P D G 1920

TGTAATTTAGAAAATATTTTGGCAGCAGCAATGATTTTAAAAAAGCAATGTAATGATTTGTAAGAATTTTGAAGTATGCTATGAGAACAGTAAAGAAAGCTTTATAAAA 5880
C N L E N I F C A P V N D F P E K T Y S N D I V E I F E V L G I E A V R A L L K 1960

GAATTAAGCACTAATATCAATTTGATAGTTCATATGTTAATTCGACATTTATCAATATTTATGATGTTATGACACAAAAGGCTTTAATGCTATAACAAGACATGCTATAAT 6000
E L R T Y I S F D S S Y V N Y R H L S I L C D V M T Q K G Y L M S I T R H G I N 2000

ACAGCTTAAGACCAATTAAGTAAATGATGATAAGTAAATGATGATAAGTAAATGATGATAAGTAAATGATGATAAGTAAATGATGATAAGTAAATGATGATAAGTAAAT 6120
R V D K G P L I K C S F E E T V E I L L E A A F A Q V D N L R G I T E N I M L 2040

GGTCAATTTGCTAAATAGGAATGGTTCATTGATATAATAGATAATCAAAAATGAAATGATGCAAAATCAAAAATTTAGAAAATTTAAGAACTTTAAGAACTGGCGGTTTACACA 6240
Q Q L C K I G T G S F D I I I D N Q K L N D A N Q M L E T I Q D L T S A G F T 2080

CCAGATGTTTACATGTTAATCACTGATGTTTAACTACCTGGGCAATTAATACGATAAATTCGCTTTACCAATTTACCAACATATAATGCTAATTTATCTCAACAGCA 6360
P D S L H V I T P D G L Q S P V A I N T I N S P L P F S P T Y N A N L L S P T A 2120

CCTATAGATAATGTAATAATTTATATACCACAAATATAATTTACAAAATTTAGGAGATAATGTAATGCTCCCAACATCAAAAAGATATAAATAATTTAGATACATTAATTTAGTGG 6480
P I D N V N N L L S P Q Y N L Q N Y G D N V H S P T S K D I N N L D T L K L G G 2160

AAATTTCCACCAACAATCACCTAAATCACCAACATCTGTTATGCACTCAGCATTTCTCTCTTTGATCATCAAAAACCAACCAACCAAGTCAACCAATTTAATTTTCTCCGAAA 6600
K F S P T Q S P K S P T S V M H S P F S P F D H Q N Q Q P V D A T N L L F S P K 2200

AATAATAATTAATGATAATTAATGATTTCACTCAACCAAAATTAATAATGTTATTCACCACTAATATATTTCCCAATGCTATGTTAGATATTTTTCACCTAATKACCT 6720
H N N I N H N Y H V F S P K P H I N N H Y I Q S P H I Y S P N H V L D I F S P K P 2240

CAAAATATCAATAATTTTACCTTCATATCCCAACATCACTCACTGATTAATGCAAAATGCTTATTTTCCCAACCTCCCAAAAATCAAAAATGATCAAAATGATGATAAT 6840
Q I N H N I Y S P S Y S P T S Y N N H N A Y Y S P T S P K N Q D Q H N V T 2280

TGGCAGTAAATCTTATGTCACCTGTTTATCACTAATCACTACCAAAAATTTACCTCACTACCAAAAATTTCCGCACTACCAAAAATTTCCGCACTACCAAAAATTTCCG 6960
S Q Y N M S P V Y S V T S P K Y S P T S P K Y S P T S P K Y S P T S P K Y S P 2320

ACATCACCAAAATTTGCGCCGATCACCAAAATTTGCGCCGATCACCAAAATTTGCGCCGATCACCAAAATTTGCGCCGATCACCAAAATTTGCGCCGATCACCAAAATTTG 7080
T S P K T Y S P T S P K Y S P T S P K Y S P T S P K Y S P T S P V A Q N I A S P N 2360

TATTCACCTTATCAATAAATCACTCACTAATTTTCCCAACATCTCCAGCATTTGATAAGTTCACCTGTCGACGCAAAAGCGGTGATGATGATGATGATGATGATGATGATGAT 7200
Y S P Y S I T S P K F S P T S P A Y S I S S P V Y D K S G V V N A H Q P M S P A 2400

TATATTTACAATCAGCTGCGAGATAAAAAATGTAAGATGTAATGTTTCCGCCATACAGCAGGCATGCGATGAAAGCAAAAATGACGACCCATTTCTCCAAATGCCT 7320
Y I L Q S P V T I E Q K V Q D V N V S P I Q Q A H V D E A K N D P P S P M 2440

TACAACATGACGACGCAAAATGAAAGAAAATATGAGGATGCACTAGTAAAGCAATAAAATAAAATTTAGTCCCAATACACACATATATATATATATATATATATATAT 7440
Y N I D E N K E D H N *** 2452

AT 7560
TTTTTTTTTTTTTTTTTCTGCAAAATGCTGAATACAT 7680
TTTTTTGGTAAAAAATA 7800
AAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAAAATAAAAA 7920
ACATTTAACAATATAAATGCTTGTGATGCAATTTGCTGTTTAAAAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATAAATA 8040
AAAGCAATTTCAAAATA 8160
CTTTATGATAAACAATCTTTCTCTTTATGATATGACATATGCTGAT 8280
CTGACTTGAATAACAATATGACCTGCTATTTTATTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTTATTTT 8400
TTAATCTCTATTTTAAATTTCTGCTACAGTAAAAAAGTGGAGTATGCGACTGCTAGGTAATTTTTTTCACCAAAAAATATAAATACTCAAAAATCAATTTTCTCTAATA 8520
TAACCTCCTAATTTTTCTGACAGGAATTC 8551

Figure 3. Nucleotide and predicted amino acid sequence of the *P. falciparum* RPII gene. The complete nucleotide sequence of the coding and 3' flanking portion of the RPII gene is shown. The ATG start codon begins at nucleotide 1 and the RPII gene open reading frame ends at nucleotide 7357. Notable features of the RPII gene nucleotide and amino acid sequence are detailed in the text according to the nucleotide and amino acid numbering shown. The RNA Pol II sequence has been deposited with the EMBL data library (accession number X16561).

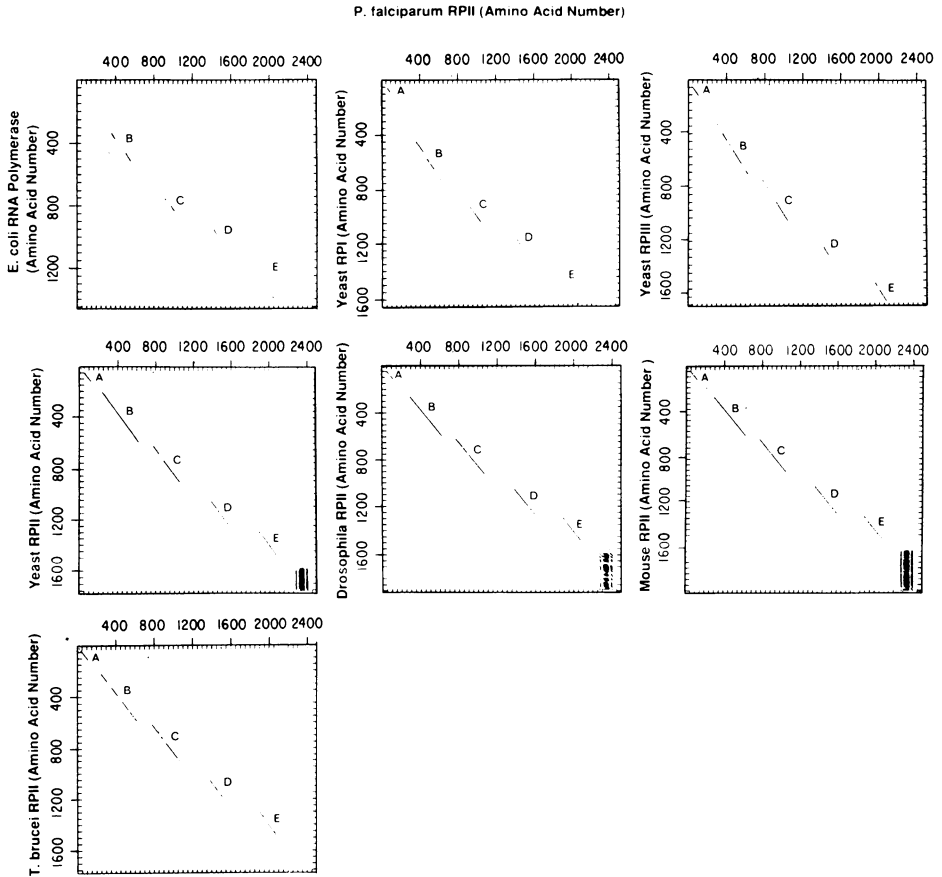


Figure 4. Amino acid comparison of the *P. falciparum* RPII subunit with other RNA polymerase subunits. The complete amino acid sequence of the *P. falciparum* RPII subunit (2452 amino acids), represented by the horizontal axis (top), was aligned with the amino acid sequence of other RNA polymerase largest subunits, represented by the vertical axis (left). The subunits aligned to the *P. falciparum* RPII subunit were *E. coli* RNA polymerase (β' subunit) (1407 amino acids), yeast RPI large subunit (1664 amino acids), yeast RPIII large subunit (1460 amino acids), yeast RPII subunit (1726 amino acids), *Drosophila* RPII subunit (1896 amino acids), mouse RPII subunit (1932 amino acids), and *T. brucei* RPII subunit (RPIIA allele) (1765 amino acids). A dot was plotted in the matrix if at least 9, or more, amino acid residues matched when windows of 15 amino acids were compared between the two sequences. The five colinear and conserved regions of the RPII subunits are marked as regions, A, B, C, D, and E. The unmarked vertical boxes at the CTD in some of the comparisons represents a region of homology between the RPII subunits that contain the conserved heptapeptide repeat domain. The different vertical sizes of the CTD homology box shows that the heptapeptide repeat domains are of different sizes.

P. falciparum RPII subunit shared 54% amino acid homology with mouse, 52% with *Drosophila*, 49% with yeast, and 42% with *T. brucei* RPII subunits (Fig. 5.). The protozoan *P. falciparum* RPII subunit thus showed more homology to higher eukaryotic forms than to either the lower eukaryote yeast, or the protozoan *T. brucei*. The same homology gradient was previously observed for the parasite dihydrofolate reductase-thymidylate synthase (DHFR-TS) gene (21).

Notable sequences and structures in the conserved regions A to E

The DNA-dependent RNA polymerases require tightly bound Zn atoms, or a divalent metal atom such as cobalt, for activity (32). An amino acid sequence thought to form a zinc finger structure that can fold about a Zn atom was present in the largest yeast RPI, RPII, and RPIII subunits as a consensus sequence Cys-X₂-Cys-X₆-His-X₂-His (2). This sequence was found in conserved region A of the *P. falciparum* RPII subunit (position 68 to 84), and was conserved in the mouse (4) and *Drosophila* RPII subunits (6) (Fig. 5.). The last histidine residue was, however, replaced with a tyrosine residue in the *T. brucei* RPII subunit (7).

The sequence, Tyr-Asn-Ala-Asp-Phe-Asp-Gly-Asp-Glu-Met-Asn (position 511 to 521), was found in conserved region B. This sequence motif is the longest conserved region in all eukaryotic RPI, RPII, and RPIII largest subunits (31).

A sequence motif predicted to form a helix-turn-helix structure is found in region B (position 378 to 411) of the *P. falciparum* RPII subunit (Fig. 5.). This conserved structure has been described in *E. coli* DNA polymerase I and T7 DNA polymerase (3), and is probably involved in DNA binding.

A common amino acid loop, Gly-X₄-Gly-Lys, is found in purine nucleotide binding, or processing, proteins such as elongation factor Tu (EF-Tu) (33). Gly-X₄-Gly forms a relatively large loop in EF-Tu (14). A similar sequence motif, Gly-Asn-Leu-Met-Gly-Lys (position 371 to 376), was found in region B of all the RPII subunits (Fig. 5.) and it may also form a loop for purine nucleotide binding.

The DD-domain consists of two aspartate residues followed by at least five uncharged residues (34). DD-domains are present in RNA-dependent RNA, and DNA polymerases. A DD-domain is found in the β' subunit of *E. coli* RNA polymerase and in conserved region B (starting at position 293) of the RPII subunits (Fig. 5.) (14). A second *P. falciparum* DD-domain was found in the unique carboxyl-terminal extension (position 2433) (Fig. 3.).

Notable sequences and structures in the variable regions A' to D'

The regions, A' to D', that separate the conserved regions have little overall amino acid homology when compared to the corresponding regions of other RPII subunits, and we refer to them as variable regions. The variable regions do, however, contain short sequences that are conserved among species. Each of the *P. falciparum* variable regions were significantly enlarged in comparison to the variable regions of other RPII subunits (Fig. 4.), and consisted of 43, 136, 278, and 284 amino acid residues, respectively (Fig. 3. and Fig. 5.). The enlarged variable regions could contain unique regulatory domains involved in the control of stage-specific, or species-specific, gene transcription during the developmentally complex life cycle of the parasite.

The variable region C' contained 5 leucine repeats beginning at position 1093, with an interruption of one aspartate (Fig. 3.), and this domain may form a leucine zipper (35,36,37). This structure is thought to facilitate protein-protein interaction (35). Similar sequences have been observed in the corresponding RPII subunit variable domain in *Drosophila*, yeast, *T. brucei*, and the mouse RPII subunit.

A remarkable basic domain was found in region C' (position 1182 to 1193) that contained 11 basic amino acids out of 12 residues. The DNA binding domain of the transcription/replication factor CTF/NF1 contains a high density of basic amino acids thought to form an α -helical structure (38,39). The basic domain present in region C' could form an α -helical structure that probably interacts with DNA in the transcriptional complex.

Region C:

P. f. II: 753 DDMPYC--SLMDCKVILKWNELLSGILCRKRVGSSGLIHLWHBQKDKTKDPLSALQKQVWNNWVYVGRITVSCSDITLASKVGLKQVRELLNKSNSKLVKA 857
 M. II: 614 DSGPKILHSPODKVWVENGELTNGILCKXKSLTSSAGSIVHTSYLHSDIITRUPYSQVWVWNNWLEDEHITIGCTGSDTASQVQNTQTKAKQNDIEMT EKA 720
 D. m. II: 506 BDFPKVLSPODKVWVHENGELTNGILCKXKSLTSSAGSILHCFJELJLHDIAGREYQVWVWNNWLEDEHITIGCTGSDTASQVQNTQTKAKQNDIEMT EKA 712
 S. c. II: 602 BDFPKVLSPODKVWVHENGELTNGILCKXKSLTSSAGSILHCFJELJLHDIAGREYQVWVWNNWLEDEHITIGCTGSDTASQVQNTQTKAKQNDIEMT EKA 697
 T. b. IIIA: 595 QDRPFP--PHNDSVYVLRKRLQLDQPITHSVNGAAGSGLHVFVFNHNSGDEVAHPIKGRVWVITFFLHNFSTVGVVQVWASDPLRQMNNDVLYVETRRNVEKTEGAA 699
 P. f. II: 858 QKGELECPGKSLVBSFETRWVNEELNCAREHAKVVASBLSLDRNNITSHVWGGKSGKSHINISQVLSVGGQVVECKRIPFGFNHRSPLPHFIRKEDYGPESRGFVSNYS 964
 M. II: 721 HNNLEBPPTQNTLRQFENQVNRHLDNCAEDRNGSSAOKLSLSVYNNFTSHVWGGKSGKSHINISQVLSVGGQVVECKRIPFGFNHRSPLPHFIRKEDYGPESRGFVSNYS 827
 D. m. II: 713 HNNLEBPPTQNTLRQFENQVNRHLDNCAEDRNGSSAOKLSLSVYNNFTSHVWGGKSGKSHINISQVLSVGGQVVECKRIPFGFNHRSPLPHFIRKEDYGPESRGFVSNYS 814
 S. c. II: 698 QANVITAKRITLRESEFEDNVDVDFLMEADKAGLAEVNRKLDLNNWVWVAGSKGSHINISQVLSVGGQVVECKRIPFGFNHRSPLPHFIRKEDYGPESRGFVSNYS 804
 T. b. IIIA: 700 WNRVLRNARKAGVLIQSFPAWNSALNKCEBEMAKKALSVRRTQSPKVFTEAGSKGFDNLIQVAVTVGQVWVASCIPFEGFRNRTLPHFMDDYGETSGMANRQV 806
 P. f. II: 965 LSGILTRVFFHANGREGLEDITAVKTAETGYQRRLLKAMEDVWVATVYVRSVSDIQLQVGEDGAGSVEFQNTIDLMKLNKELMLKVFYVDFD 1064
 M. II: 828 LAGLITVFFHANGREGLEDITAVKTAETGYQRRLLKAMEDVWVATVYVRSVSDIQLQVGEDGAGSVEFQNTIDLMKLNKELMLKVFYVDFD 927
 D. m. II: 820 LAGLITVFFHANGREGLEDITAVKTAETGYQRRLLKAMEDVWVATVYVRSVSDIQLQVGEDGAGSVEFQNTIDLMKLNKELMLKVFYVDFD 914
 S. c. II: 805 LAGLITVFFHANGREGLEDITAVKTAETGYQRRLLKAMEDVWVATVYVRSVSDIQLQVGEDGAGSVEFQNTIDLMKLNKELMLKVFYVDFD 909
 T. b. IIIA: 807 VEGVLRPHFEFFHANGREGLEDITAVKTAETGYQRRLLKAMEDVWVATVYVRSVSDIQLQVGEDGAGSVEFQNTIDLMKLNKELMLKVFYVDFD 907

Region D:

P. f. II: 1343 SNETLISVEAQQMATHLKLHLRITVLSLITQTHKVSVKGLDMLLQAEKELFYKSLCHHPGEOGALAAQSICEPATOMTLNTHFAVGSKNKVTLVGVRPKELINI 1448
 M. II: 1025 QDPLSRQAQEMATLLFNHLRSHLCSRRMAEERTSKGATDMLLQAEKELTESFNQCAAPGEOGALAAQSICEPATOMTLNTHFAVGSKNKVTLVGVRPKELINI 1130
 D. m. II: 1017 QNDRTSKQAEEMATLLFQCLTRSHLCSRRMAEERTLSKATDMLLQAEKELTESFNQCAAPGEOGALAAQSICEPATOMTLNTHFAVGSKNKVTLVGVRPKELINI 1122
 S. c. II: 1002 GKVELTQMAQRDAVILLFQCLTRSHLCSRRMAEERTLSKATDMLLQAEKELTESFNQCAAPGEOGALAAQSICEPATOMTLNTHFAVGSKNKVTLVGVRPKELINI 1107
 T. b. IIIA: 1022 TRAVTSREKTESALTLFNVHLRQLLAKRVLKEDVLDNDRAFEYLKELRTYVHOSLTPGENTCALAAQSICEPATOMTLNTHFAVGSKNKVTLVGVRPKELINI 1127
 P. f. II: 1449 VKNMTPTSTLVLDDHVSNDQOKAQTLLTLYLMTLLOVSHAOVLDPMVTTTLLLEBKSWWNEFVDEPDDTOYSL----GBMVLRLQTLNHVNEKCLTM 1348
 M. II: 1131 SKRPHKPTSLVETLLGOSRDQAEKADLCLREHTLIRAVTANTAIYDPMFOSTVVAEDKSWWVYVDEPDDVARI----SPVLLRVELDRKRMVTKCLTM 1228
 D. m. II: 1123 SKRPHKPTSLVETLLGMAARDAEKANVLCRLREHTLIRAVTANTAIYDPPORTVLSDEKSWWVYVDEPDDTRI----SPVLLRVELDRKRMVTKCLTM 1220
 S. c. II: 1108 AKSMKPTSLVETLLPECHANDQOAKLITSAETHLTKSVTLASBELYDPPDRSTVLPEDDEITLQLHFSLLDDEAQSDFDQO----SPVLLRVELDRKRMVTKCLTM 1209
 T. b. IIIA: 1128 SRNQLFASVLSLFPYDEXRVAQKAOHL--LHETQLLESIDRRIQVTPDPRHFWEDRDILELEMVWVDESDAELRQEVVAGSPVAVRVELDMDVTDMAIDM 1232
 P. f. II: 1549 KDFVYVHSVFSSEDDIIVTDNDSMDVLRIR 1581
 M. II: 1229 EQLAKINMAGGDDLNGEINDNDEKLVLRIR 1260
 D. m. II: 1221 EQLAKINMAGGDDLNGEINDNDEKLVLRIR 1252
 S. c. II: 1210 GQGERIKUTEKVLVYVHNSBDEKLVLRIR 1241
 T. b. IIIA: 1233 KDKVKNLVRDES--IIEIGMANNRQRTIR 1264

Region E:

P. f. II: 1866 DNFLEKLMEOCLSTLKLKRGTEMTKVVYH----REESKITYSDNGLVRRSSRHHVDFDQCMLENIFACQVDEK-----TWSNDIWEIPEVLGTEAVRHALK 1960
 M. II: 1282 DNFLEKLMEOCLSTLKLKRGTEMTKVVYHLLPDDNKKKIITIDEGKRALQEMILLDTDQVSLMHHMLSEKQVDFVH-----TWSNDIWEIPEVLGTEAVRHALK 1380
 D. m. II: 1274 DNFLEKLMEOCLSTLKLKRGTEMTKVVYHLLPDDNKKKIITIDEGKRALQEMILLDTDQVSLMHHMLSEKQVDFVH-----TWSNDIWEIPEVLGTEAVRHALK 1372
 S. c. II: 1257 DNFLEKLMEOCLSTLKLKRGTEMTKVVYHLLPDDNKKKIITIDEGKRALQEMILLDTDQVSLMHHMLSEKQVDFVH-----TWSNDIWEIPEVLGTEAVRHALK 1350
 T. b. IIIA: 1274 TPVLRKRETPALLARVLRGIFGVRRALL----MPTTEFVDDOVLGMSGNLTHADTDLRRAEATGVGVEDOKNIIIAVAVTSSNKKVPEVCSLLGTEAANSKHL 1375
 P. f. II: 1961 ELRVTISFVSVYVNRHLSLGDVMTQGLMSTIRHGNNVTK--EPLIKCSFEEVWLEAAKAVQVLRJENHMLCOLCKICTGSGSEPIIIFDK 2060
 M. II: 1381 ELRVTISFVSVYVNRHLSLGDVMTQGLMSTIRHGNNVTK--EPLIKCSFEEVWLEAAKAVQVLRJENHMLCOLCKICTGSGSEPIIIFDK 1480
 D. m. II: 1373 ENMVALEQVSVYVNRHLLALGDVMTQGLMSTIRHGNNVTK--EPLIKCSFEEVWLEAAKAVQVLRJENHMLCOLCKICTGSGSEPIIIFDK 1472
 S. c. II: 1351 ENMVALEQVSVYVNRHLLALGDVMTQGLMSTIRHGNNVTK--EPLIKCSFEEVWLEAAKAVQVLRJENHMLCOLCKICTGSGSEPIIIFDK 1450
 T. b. IIIA: 1376 ELREAVLALGEMVNRHRTLLMATTIQGRLMVA--SFGSGVNSDSSPLMRRSFEETWLVLMVAVSAGSDPVRGVSNNVLCGRVLRVGTQLFDVLDVNRMAAL 1476

Figure 5. Amino acid alignments of the five conserved regions in the RPII subunits. The amino acid alignments of the conserved regions, A, B, C, D, and E, is shown for *P. falciparum* (P.f.II), mouse (M.II), *Drosophila* (D.m.II), *Saccharomyces cerevisiae* (S.c.II), and *T. brucei* (T.b.II) (allele RPIIA). Amino acids were boxed when three out of five amino acids matched in the vertical column. The amino acid location is given in the left and right column.

The composition of regions C' and D' was enriched in acidic amino acids (21%). Four extremely acidic domains were found in C' (position 1258 to 1290) and D' position 1602 to 1612, 1746 to 1759, and 1806 to 1820). The longest acidic domain, in region C', contained a novel Tyr-Asp-Asp-Asp-Asp repeat sequence.

The variable regions C' and D' contained high overall charge densities with 54 acidic and 42 basic residues in C' (out of 278 residues), and 64 acidic and 38 basic residues in D' (out of 284 residues). Because charge clusters are apparently associated with functional domains of many cellular transcription factors and regulatory proteins (40), we propose that the enlarged variable domains C' and D' could provide specific and unique regulatory functions to the *P. falciparum* RPII subunit.

The notable sequence, Asn-Lys-X-Asp, was found in both variable domains C' (position 1074 to 1077) and D' (position 1821 to 1824). This sequence motif is found in *E. coli* EF-Tu, as well as in other elongation factors, and forms a loop for guanine base interaction (33). This sequence motif is unique to the *P. falciparum* RPII subunit.

The variable regions in the *P. falciparum* RPII subunit were extremely asparagine rich (24%) compared to the conserved regions (5%). Asparagine repeats were located in regions B' (position 707 to 725), C' (position 1144 to 1159), and D' (position 1687 to 1694). Two glutamine rich transcriptional activation domains in the transcription factor Sp1 (41), unlike other charged activation domains (42,43), are largely devoid of charged residues in their primary sequence. Amide moieties of the glutamine side chains may participate in hydrogen bonding to the RNA polymerase, or another component of the transcriptional complex (41). The asparagine repeats in the *P. falciparum* RPII subunit could be similarly involved in an interaction with another component of the transcriptional complex via hydrogen bonding.

1	Y	S	P	S	-	-	-											
2	Y	S	P	T	S	P	T	Y	N	A	N	N	A	Y				
3	Y	S	P	T	S	P	K	N	Q	N	D	Q	M	N	V	N	S	Q
4	Y	N	V	M	S	P	V											
5	Y	S	V	T	S	P	K											
6	Y	S	P	T	S	P	K											
7	Y	S	P	T	S	P	K											
8	Y	S	P	T	S	P	K											
9	Y	S	P	T	S	P	K											
10	Y	S	P	T	S	P	K											
11	Y	S	P	T	S	P	K											
12	Y	S	P	T	S	P	K											
13	Y	S	P	T	S	P	V	A	Q	N	I							
14	A	S	P	N	Y	S	P	-										
15	Y	S	I	T	S	P	K											
16	F	S	P	T	S	P	A											
17	Y	S	I	S	S	P	V											

(Followed by 68 amino acid residues)

Figure 6. The *P. falciparum* RPII subunit heptapeptide repeat domain. The *P. falciparum* heptapeptide repeat domain is represented by 17 highly conserved repeats with the consensus sequence Tyr-Ser-Pro-Thr-Ser-Pro-Lys. The repeat region corresponds to amino acids 2247 to 2384 (Fig. 3.).

Remarkably, certain amino acid residues were not found in some of the variable regions. Cysteine was not found in variable region A', B', and D', but it was found in each conserved region, and in variable region C'. Tryptophan was absent from all of the variable regions, but it was found in conserved region B, C, D, and E.

Structure of the carboxyl-terminal domain

A remarkable feature of the eukaryotic RPII subunits is the presence of a heptapeptide repeat (consensus Tyr-Ser-Pro-Thr-Ser-Pro-Ser) in the CTD. The *P. falciparum* RPII subunit contains a small cluster of 17 highly conserved heptapeptide repeats (position 2247 to 2384) (Fig. 6.). The amino acid at position 7 of the repeat unit in *P. falciparum* contained a lysine rather than a serine, which is present in all other RPII subunits with the repeats. This same replacement was found on the carboxyl-terminal side of the repeat element in the mouse and hamster repeats (5), where the serine residue tended to be replaced by a charged amino acid (particularly lysine).

The CTD of the *P. falciparum* RPII subunit contains 392 amino acids that extend from the 3' side of conserved region E (position 2061) to the carboxyl-terminal amino acid

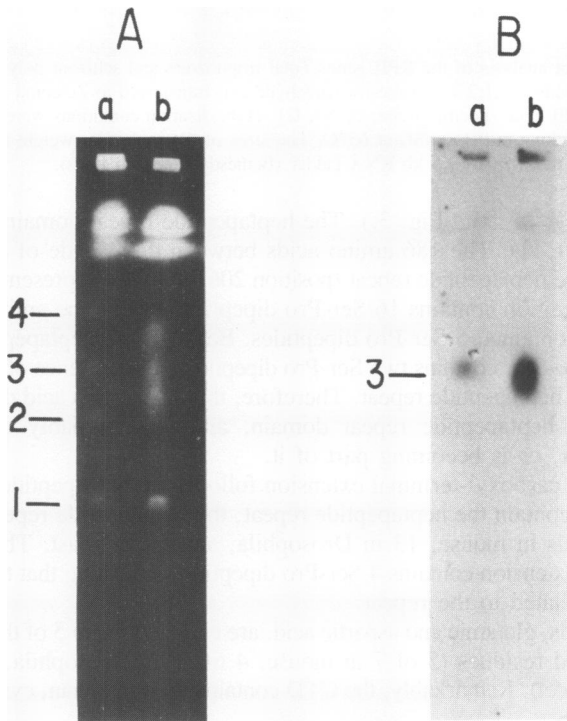


Figure 7. Chromosome assignment of the *P. falciparum* RPII gene. Chromosomes 1, 2, 3, and 4 of *P. falciparum* strains FCR3 (a) and Honduras-1 (b) were separated by pulsed field gradient gel electrophoresis and stained with ethidium bromide (A). The electrophoresis parameters used for this gel were intended to resolve chromosomes 1 to 4, only (28). The chromosomal DNA was transferred to nylon and hybridized with the RPII gene specific probe, cDNA C1, in $6\times$ SSC at 42°C , and was washed in $1\times$ SSC at 65°C (B). The chromosome number assignments are shown.

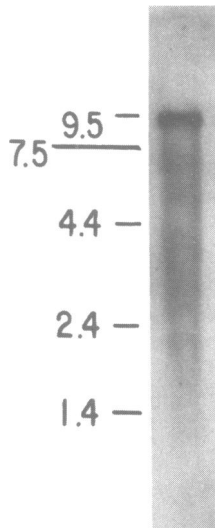


Figure 8. Northern blot analysis of the RPII gene. Total trophozoite and schizont poly(A⁺) RNA (4 μ g per lane) was electrophoresed on a 1.2% agarose-formaldehyde gel, transferred to Zetabind nylon membrane, and hybridized with the RPII gene specific probe, cDNA C1. Hybridization conditions were previously described (25), and washing was done in 0.1 \times SSC at 65°C. The sizes of the molecular weight markers are shown in kb. The markers were the 0.24 to 9.5 kb RNA ladder (Bethesda Research Labs).

(position 2452) (Fig. 1F. and Fig. 3.). The heptapeptide repeat domain is located in the middle part of the CTD. The 186 amino acids between the 3' side of conserved region E and the start of the heptapeptide repeat (position 2061 to 2246) represents another domain of the CTD. This region contains 16 Ser-Pro dipeptides, while the amino-terminal 2060 amino acids only contained 5 Ser-Pro dipeptides. Because one heptapeptide repeat, Tyr-Ser-Pro-Thr-Ser-Pro-Ser, contains two Ser-Pro dipeptides, the presence of Ser-Pro indicates a relatedness to the heptapeptide repeat. Therefore, this 186 amino acid domain represents a highly diverged heptapeptide repeat domain, and was probably once part of the heptapeptide repeat, or is becoming part of it.

A 68 amino acid carboxyl-terminal extension follows the heptapeptide repeat. In other RPII subunits that contain the heptapeptide repeat, the heptapeptide repeat is followed by only 10 amino acids in mouse, 13 in *Drosophila*, and 10 in yeast. The 68 amino acid carboxyl-terminal extension contains 4 Ser-Pro dipeptides indicating that the extension may also be distantly related to the repeat.

Acidic amino acids, glutamic and aspartic acid, are enriched for in 5 of the last 9 carboxyl-terminal amino acid residues (5 of 7 in mouse, 4 of 12 in *Drosophila*, 2 of 6 in yeast, and 6 of 7 in *T. brucei*). Remarkably, the CTD contains no tryptophan, cysteine, or arginine amino acid residues.

The abundant serine residues in the CTD were previously suggested to provide phosphorylation sites used for regulating RPII activity (10). A protein kinase, CTD kinase, that phosphorylates *in vitro* the CTD of the RPII subunit from mouse cells was recently isolated (44). The CTD kinase phosphorylates one or more serine residues in the CTD

heptapeptide repeat. The purification of CTD kinase now provides direct evidence that the CTD of the RPII subunit may serve as a phosphorylation domain which modulates RPII activity.

Chromosomal location and copy number of the RPII gene

Recently, we clearly separated and determined the sizes of 13 DNA bands, or chromosomes, by pulsed field gradient gel (PFG) electrophoresis of *P. falciparum* strain FCR3 (29). cDNA C1, which is specific for the RPII gene, was hybridized to Southern blot of separated chromosomes 1 to 4, and the RPII gene was located on chromosome 3, based on the number assignment of FCR3 chromosomes, and on a Honduras-1 chromosome that corresponded in size to chromosome 3 of FCR3 (Fig. 7.).

Based on our more precise measurements of the size of the malaria chromosomes by PFG electrophoresis, and an estimated total haploid genome content of 2.64×10^7 bp (29), the copy number of the RPII gene in both Honduras-1 and FCR3 was found to be one copy per parasite (data not shown).

Expression of the RPII gene

Northern blot analysis of a mixture of trophozoite and schizont stage mRNA revealed a single species of RPII mRNA of 9 kb by hybridization to the RPII gene specific probe, cDNA C1 (Fig. 8.). Because of the single copy nature of the RPII subunit gene we conclude that the same RPII gene must be active throughout the life cycle of the parasite both in man and mosquitoes.

We determined the sensitivity of parasite growth to α -amanitin in vitro, and were unable to detect a growth inhibitory effect of α -amanitin at concentrations as high as 20 μ g/ml. Because RPII subunit α -amanitin sensitivity should be measured by an in vitro transcription assay, which is not yet developed for this organism, we do not know if the RPII subunit is in fact resistant, or if the cell is impermeable to the drug.

ACKNOWLEDGEMENTS

This work was funded by NIH grants #5R01 AI-20437 and #5R01 AI-22038 to J.I., and #5R29 AI-26651 to D.J.B. The authors are grateful to Thomas Ciardelli, the Structural Biology Laboratory at Dartmouth, and the Dartmouth Molecular Genetics Center for synthesizing the oligonucleotides used in this work.

REFERENCES

1. Sentenac, A. (1985) *CRTC Crit. Rev. Biochem.* **18**, 31–90.
2. Memet, S., Gouy, M., Marck, C., Sentenac, A., and Buhler, J.-M. (1989) *J. Biol. Chem.* **263**, 2830–2839.
3. Allison, L. A., Moyle, M., Shales, M., and Ingles, C. J. (1985) *Cell* **42**, 599–610.
4. Ahearn, J. M., Jr., Bartolomei, M. S., West, M. L., Cisek, L. J., and Corden, J. L. (1987) *J. Biol. Chem.* **262**, 10695–10705.
5. Allison, L. A., Wong, J. K.-C., Fitzpatrick, V. D., Moyle, M., and Ingles, C. J. (1988) *Mol. Cell. Biol.* **8**, 321–329.
6. Jokerst, R. S., Weeks, J. R., Zehring, W. A., and Greenleaf, A. L. (1989) *Mol. Gen. Genet.* **215**, 266–275.
7. Smith, J. L., Levin, J. R., Ingles, C. J., and Agabian, N. (1989) *Cell* **56**, 815–827.
8. Evers, R., Hammer, A., Kock, J., Jess, W., Borst, P., Memet, S., and Cornelissen, A. W. C. A. (1989) *Cell* **56**, 585–597.
9. Ovchinnikov, Y. A., Monastyrskaya, G. S., Gubanov, V. V., Guryev, S. O., Salomatina, I. S., Shvaea, T. M., Lipkin, V. M., and Sverdlov, E. D. (1982) *Nucl. Acids Res.* **10**, 4035–4044.
10. Zehring, W. A., Lee, J. M., Weeks, J. M., Jokerst, R. S., and Greenleaf, A. L. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3698–3702.

11. Corden, J. L., Cadena, D. L., Ahearn, J. M., and Dahmus, M. E. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 7934–7938.
12. Ingles, C. J., Moyle, M., Allison, L. A., Wong, J. K.-C., Archambault, J., and Freisen, J. D. (1987) *Mol. Cell Biol.* **52**, 383–393.
13. Nonet, M., Sweetser, D., and Young, R. A. (1987) *Cell* **50**, 909–915.
14. Cornelissen, A. W.C.A., Evers, R., and Kock, J. (1988) *Oxford Surveys Euk. Gen.* **5**, 91–131.
15. Trager, W., and Jensen, J. B. (1976) *Science* **193**, 673–675.
16. Gritzmacher, C. A., and Reese, R. T. (1984) *J. Bacteriol.* **160**, 1165–1167.
17. Banyal, H. S., and Inselburg, J. (1985) *Am. J. Trop. Med. Hyg.* **34**, 1055–1064.
18. Inselburg, J. (1983) *J. Parasitol.* **69**, 584–591.
19. Maniatis, T., Fritsch, E. F., and Sambrook, J. (1982) *Molecular Cloning: A Laboratory Manual*. (Cold Spring Harbor, New York: Cold Spring Harbor Laboratory).
20. Horii, T., Bzik, D. J., and Inselburg, J. (1988) *Mol. Biochem. Parasitol.* **30**, 9–18.
21. Bzik, D. J., Li, W.-B., Horri, T. and Inselburg, J. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 8360–8364.
22. Dretzen, G., Bellard, P., Sassone-Corsi, P., and Chambon, P. (1981) *Anal. Biochem.* **112**, 295–298.
23. Feinberg, A. P., and Vogelstein, B. (1983) *Anal. Biochem.* **132**, 6–13.
24. Li, W.-B., Bzik, D. J., Horii, T., and Inselburg, J. (1989) *Mol. Biochem. Parasitol.* **33**, 13–26.
25. Sanger, F., Nicklen, S., and Coulson, A. R. (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463–5467.
26. Bzik, D. J., Li, W.-B., Horii, T., and Inselburg, J. (1988) *Mol. Biochem. Parasitol.* **30**, 279–288.
27. Schwartz, D. C., and Cantor, C. R. (1984) *Cell* **37**, 67–75.
28. Chu, G., Vollrath, D., and Davis, R. W. (1986) *Science* **234**, 1582–1585.
29. Gu, H., Inselburg, J., Bzik, D. J., and Li, W.-B. (1989) *Exp. Parasitol.* (in press).
30. Weber, J. L. (1987) *Gene* **52**, 103–109.
31. Memet, S., Saurin, W., and Sentenac, A. (1988) *J. Biol. Chem.* **263**, 10048–10051.
32. Wu, F. Y.-H., and Wu, C.-W. (1981) In *Advances in Inorganic Biochemistry*, Vol. 3, G. L. Eichhorn and L. G. Marzilli, (eds.), Elsevier/North-Holland, New York, pp. 143–166.
33. La Cour, T. F. M., Nyborg, J., Thirup, S., and Clark, B. F. C. (1985) *EMBO J.* **4**, 2385–2388.
34. Zavriev, S. K., and Borisova, O. V. (1987) [translated from] *Mol. Biol.* **21**, 229–241.
35. Landschulz, W. H., Johnson, P. F., and McKnight, S. L. (1988) *Science* **240**, 1759–1764.
36. Kouzarides, T., and Ziff, E. (1988) *Nature* **336**, 646–651.
37. Struhl, K. (1989) *TIBS* **14**, 137–140.
38. Santoro, C., Mermod, N., Andrews, P. C., and Tjian, R. (1988) *Nature* **334**, 218–224.
39. Mitchell, P. J., and Tjian, R. (1989) *Science* **245**, 371–378.
40. Brendel, V., and Karlin, S. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5698–5702.
41. Courey, A. J., and Tjian, R. (1988) *Cell* **55**, 887–898.
42. Ma, J., and Ptashne, M. (1987) *Cell* **48**, 847–853.
43. Hope, I. A., and Struhl, K. (1986) *Cell* **46**, 885–894.
44. Cisek, L. J., and Corden, J. L. (1989) *Nature* **339**, 679–684.

This article, submitted on disc, has been automatically converted into this typeset format by the publisher.