

# Alternatively spliced products of the maize *P* gene encode proteins with homology to the DNA-binding domain of *myb*-like transcription factors

(pericarp/phlobaphene/flavonoids/regulation/*CI* gene)

ERICH GROTEWOLD, PRASANNA ATHMA, AND THOMAS PETERSON<sup>†</sup>

Cold Spring Harbor Laboratory, Box 100, Cold Spring Harbor, NY 11724

Communicated by James D. Watson, March 4, 1991

**ABSTRACT** The *Zea mays P* gene has been postulated to regulate the biosynthetic pathway of a flavonoid-derived pigment in certain floral tissues [Styles, E. D. & Ceska, O. (1977) *Can. J. Genet. Cytol.* 19, 289–302]. We have characterized two *P* transcripts that are alternatively spliced at their 3' ends. One message of 1802 nucleotides encodes a 43.7-kDa protein with an N-terminal region showing ≈40% homology to the DNA-binding domain of several members of the *myb* family of protooncogene proteins. A second message of 945 nucleotides encodes a 17.3-kDa protein that contains most of the *myb*-homologous domain but differs from the first protein at the C terminus. The deduced *P*-encoded proteins show an even higher homology (70%) in the *myb*-homologous domain to the maize regulatory gene *CI*. Additionally, the *P* and *CI* genes are structurally similar in the sizes and positions of the first and second exons and first intron. We show that *P* is required for accumulation in the pericarp of transcripts of two genes (*A1* and *C2*) encoding enzymes for flavonoid biosynthesis—genes also regulated by *CI* in the aleurone.

Coordinated development of multicellular organisms requires precisely regulated differential gene expression; yet the molecular mechanisms involved in organ- and cell-specific regulation remain largely unknown. Analysis of regulatory genes and their targets is difficult if the genes' products are essential for viability. Flavonoid pigment biosynthesis in plants is an ideal system for studying gene regulation in higher eukaryotes, because the presence or absence of pigment has no deleterious effects, pigmentation is a convenient visual indicator of gene expression, and the biochemical steps to pigment synthesis are well defined (1).

The flavonoid pigments commonly found in maize are either anthocyanins or phlobaphenes; although these two pigments have their initial synthetic steps in common, there are significant differences between them. Anthocyanins are derived from flavan-3,4-diol and can be produced in most tissues of the maize plant. Phlobaphenes are formed by the nonenzymatic polymerization of flavan-4-ol and are found only in certain floral tissues, including the pericarp and the glumes of the cob (2). [The pericarp is the outer covering of the maize kernel, derived from the ovary wall; the cob glumes are floral bracts that subtend the kernel (3)]. Four genes (*R*, *B*, *CI*, and *PI*) are known to regulate anthocyanin biosynthesis in specific parts of the plant (recently reviewed in ref. 4). The *R* gene family is involved in anthocyanin pigmentation in the aleurone, scutellum, coleoptile, roots, and anthers (1); *R*-like proteins have homology to the helix-loop-helix domain of the *myc* oncogene products (4). The *B* gene family, which is required for anthocyanin pigmentation in several other plant parts, has homology with *R* (5). The *CI* gene is

required for anthocyanin pigmentation of the kernel aleurone and embryo, while the *PI* gene, which is homologous to *CI* (6), is required for pigmentation in most vegetative plant parts (1). Genetic and biochemical studies indicate that the biosynthetic pathway leading to flavan-4-ol is regulated by the *P* gene of maize (2); to date, no other regulator of this pathway has been described.

Due to its conspicuous red pigmentation phenotype, the *P* gene has been the object of extensive genetic analysis since the pioneering work of Emerson (7). The *P-vv* allele, which specifies variegated pericarp and cob, carries the transposable element *Ac* inserted in the *P* gene (8, 9). In a series of classic experiments, Brink and Greenblatt (10, 11) inferred important parameters of *Ac* transposition from the patterns of clonal sectors observed on *P-vv* ears. Using the same *P-vv* allele, Lechelt *et al.* (9) cloned genomic DNA from the *P* locus by using *Ac* as a probe, and they identified several *P*-specific RNAs transcribed from a 7-kilobase (kb) region flanked by two 5.2-kb direct repeats. Here, we present the sequences of two *P* transcripts that arise by alternative splicing.<sup>‡</sup> The *P*-encoded proteins share significant homology with products of the vertebrate protooncogene *c-myb* and other transcriptional activators. We also show that *P* regulates the accumulation in the pericarp of RNA from two unlinked genes encoding enzymes for flavonoid biosynthesis. The finding that *P* and other maize genes that regulate anthocyanin biosynthesis (4) are homologous to mammalian transcription factors supports the idea that the molecular mechanisms of gene regulation in the plant and animal kingdoms are fundamentally similar.

## MATERIALS AND METHODS

**Maize Stocks.** *P* alleles are identified by a two-letter suffix indicating their expression in the pericarp and cob glumes: *P-rr* specifies red pericarp and cob, and *P-ww* specifies white (colorless) pericarp and cob (see figure 1 in ref. 12). Stocks carrying *P-vv* (variegated pericarp and cob) were obtained from Tony Pryor (Commonwealth Scientific and Industrial Research Organisation, Canberra, Australia). *P-ovov-1114* (orange variegated pericarp and cob) was derived from *P-vv* by intragenic transposition of *Ac* (12). *P-rr-4B2* and *P-rr-4026* are *P-rr* revertants derived by excision of *Ac* from *P-vv* and *P-ovov-1114*, respectively (E.G. and T.P., unpublished data). *P-ww-1112* is a *P-ww* allele derived from *P-vv* by a deletion of *Ac* and 12.5 kb of the *P* gene including the *P* coding sequences (13). *P-ww-10:443-3* is a *P-ww* allele derived by *Ac* excision from *P-vv\*-8393-4*, which in turn was derived from

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "advertisement" in accordance with 18 U.S.C. §1734 solely to indicate this fact.

Abbreviation: nt, nucleotide(s).

<sup>†</sup>To whom reprint requests should be addressed.

<sup>‡</sup>The sequences reported in this paper have been deposited in the GenBank data base (accession nos. M62878 and M62879).

*P-ovov-1114* by intragenic transposition of *Ac* (E.G. and T.P., unpublished data).

**Primers.** The oligodeoxynucleotide primers were as follows (numbers in parentheses indicate position in sequence shown in Fig. 3A): EP1PE, 5'-GGATACACGCTGGCAGTCG-3' (69-51); EP5-8, 5'-ACGCGCACCAGCTGCTAACCGTG-3' (130-153); PA-B7, 5'-CACACCGGAGTCGTATGTGG-3' (233-214); EP3-10, 5'-CTGTCGGCCTCCCCCAGACTAGG-3' (1032-1008); EP3-7, 5'-ACGAATTCACGCACCTAAAGCAGAAGCGAAC-3' (1617-1594); EP5-2, 5'-GCCTCGAGAATTCAAGCTTTTTTTTTTT-3'; EP5-4, 5'-GCCTCGAGAATTCAAGCTT-3'.

**Primer Extension.** Three micrograms of poly(A)<sup>+</sup> pericarp RNA was reverse-transcribed using 10 pmol of primer PA-B7. The product was tailed with deoxyadenylate residues, and 5% of the tailed product was amplified by PCR using 10 pmol of the internal primer EP1PE and 10 pmol of EP5-2 and of EP5-4, in the reaction conditions described by Perkin-Elmer/Cetus. A 5-min denaturation step at 94°C was followed by 30 cycles of 1 min at 94°C, 1 min at 50°C, and 1 min at 72°C, plus a final extension of 20 min at 72°C. Products were cloned into plasmid vectors, and the sequences of 17 independent clones were determined.

**cDNA Cloning.** Three independent cDNA libraries were constructed (Stratagene cDNA synthesis kit); from 3 × 10<sup>6</sup> clones, 15 *P* cDNA clones were obtained and sequenced. cDNA clones from the 5' region of *P* were obtained as follows. A first strand of cDNA was synthesized as above, but using (dT)<sub>17</sub> as a primer. The RNase A-treated first strand cDNA was PCR-amplified in two separate 25-μl reactions using 10 pmol of each primer from the primer pairs EP5-8/EP3-7 or EP5-8/EP3-10 in the above conditions, except that for the reaction using EP5-8/EP3-10, 50% of the dGTP in the reaction mixture was replaced by 7-deaza-dGTP. Reaction conditions were as above, except that annealing was at 60°C and polymerizations were for 2 min. PCR products were run on agarose gels and analyzed by Southern hybridizations with *P* gene probes. The sequence of several independent clones was compared with the sequence of the total PCR product reamplified by asymmetric PCR (14).

## RESULTS

**Localization of the 5' End of the *P* Gene.** Previous studies (9) indicated that the most 5' restriction fragment hybridizing to *P*-specific transcripts was a 3-kb *Sal* I-*Pst* I fragment (termed fragment 8 in ref. 9). This fragment was divided into two pieces (8A and 8B) by an internal *Kpn* I site. In Northern blot hybridizations (Fig. 1), 8A hybridized to four *P* transcripts of 7.0, 6.5, 2.0, and 1.4 kb (the 7.0- and 6.5-kb transcripts are poorly resolved in Fig. 1). In contrast, 8B hybridized only to the two larger transcripts. The exons

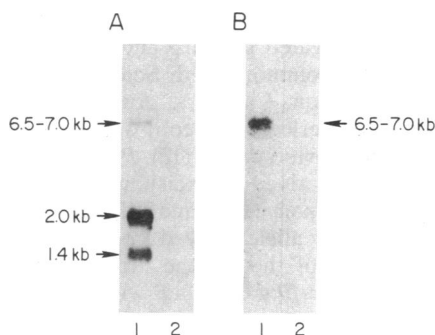


FIG. 1. Northern analysis of *P* gene transcripts. Poly(A)<sup>+</sup> RNA from plants homozygous for *P-rr-4B2* (lane 1) or *P-ww-1112* (lane 2) was hybridized with *P* genomic fragments 8A (A) and 8B (B). Arrows indicate the *P*-specific transcripts detected with each probe.

within 8A were mapped by RNase protection (Fig. 2). Probe A gave protected bands of ≈450 and ≈130 nt (Fig. 2B, lanes 1, 2, and 6). The 450-nt band was replaced by a band of ≈200 nt when probe B was used (lanes 3 and 9), suggesting that the 3' end of an ≈450-nt exon is located 200 nt 3' of the *Bal* I site. The 130-nt protected fragment obtained with probes A, B, and C (lanes 1, 3, 6, and 12) indicates that a second exon lies between the *Stu* I and *Kpn* I sites in fragment 8A; this exon was confirmed by cDNA cloning (see below). The 5' end of the 450-nt protected band was mapped by nuclease S1 protection using probe D in Fig. 2A. The 69- to 70-nt protected fragment obtained (Fig. 2C, lane 1) places the 5' end of the 450-nt exon at position -1 or +1 of the sequence shown in Fig. 3A.

Northern blot hybridizations using as probe a 295-bp fragment contained within the 450-nt exon (5' of the *Bal* I site) gave a pattern identical to that obtained when the whole of 8A was used as probe (data not shown), indicating that the 450-nt exon is present in all the *P* transcripts detected on Northern blots. This result, together with the previous finding that probes from a 13.5-kb contiguous region upstream of probe 8A did not detect RNAs specific for the *P* gene (9), suggests that the 450-nt exon is the most 5' exon of the *P* gene. This hypothesis was tested by primer extension experiments using a primer located ≈230 bp downstream of the putative 5' end of the exon. Seventeen independent clones arising from the primer extension and amplification were sequenced (*Materials and Methods*); nine clones reached nt +1 in the sequence shown in Fig. 3A, three ended at nt +5, and the rest ended at various sites between nt +15 and +1, possibly due to random termination of the reverse transcriptase. These results confirm the initiation of transcription of the *P* transcripts at position +1 in Fig. 3A and B.

**Isolation of *P* cDNA Clones.** To obtain cDNA clones, 3 × 10<sup>6</sup> independent recombinant clones from three oligo(dT)-primed libraries were screened using a previously cloned genomic probe that contains the 3' end of the *P* gene (fragment 14 in Fig. 4). Fifteen independent clones were isolated ranging between 250 and 800 bp in size, and their sequences were colinear with the genomic sequence 5' of fragment 14. Three alternative polyadenylation sites were determined as almost equally distributed among the cDNA clones (Fig. 3). Fully extended cDNAs were PCR-amplified using a primer (EP5-8) complementary to a sequence located 130 bp downstream of the transcription start site, within the 452-nt 5' exon; based on the previous results, this sequence is present in all *P* transcripts. Two primers for the 3' end were synthesized, one (EP3-7) located ≈180 bp from the polyadenylation site and the other (EP3-10) located near the 5' end of the 800-bp cDNA clone obtained from the conventional library. The primers located closer to each other (EP3-10 and EP5-8) produced an amplification product of ≈900 bp, while the primers located further apart (EP3-7 and EP5-8) produced a product of ≈450 bp (data not shown). The sequence of the products indicated that they were derived from two different transcripts (Fig. 3).

A comparison of the cDNA sequences with the genomic sequence of the *P* locus (unpublished data) shows that the cDNAs each contain three exons (Fig. 4). Both transcripts contain the first two exons of 452 and 130 nt, which are separated by an intron of 118 nt. However, the 3' exon is different due to alternative splicing. The longer transcript (1802 nt) arises by splicing of the 130-nt exon to a 1220-nt exon located 4.9 kb downstream. The shorter transcript (945 nt) arises by splicing of the 130-nt exon to a 363-nt exon located 7 kb downstream. The two *P* transcripts share an identical 318-nt sequence at their 3' ends (nt 1485 to the end in Fig. 3A; nt 628 to the end in Fig. 3B). This identity is due to a 1.2-kb direct repeat in the genomic sequence, which partially overlaps with the two alternative 3' exons (Fig. 4).

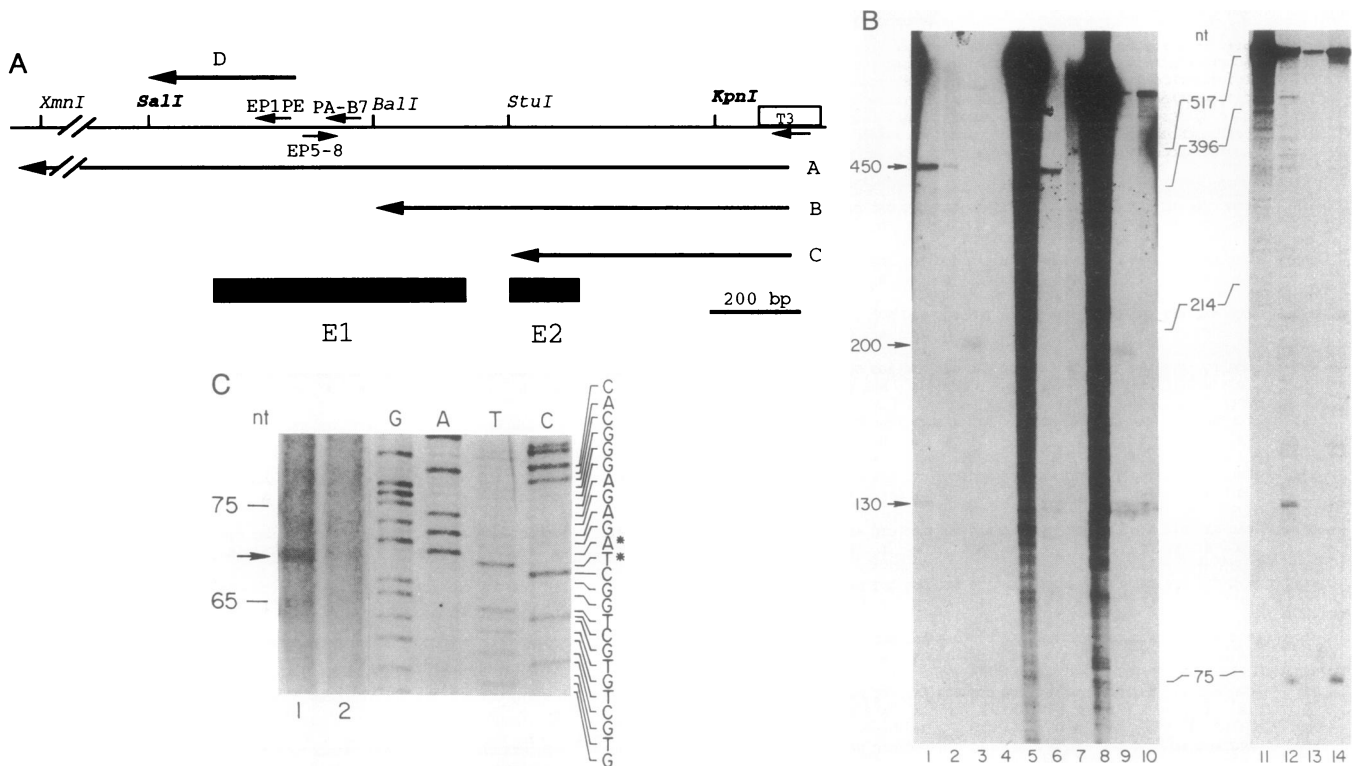


FIG. 2. RNase A and S1 nuclease protection mapping of the *P*-gene 5' region. (A) Map of genomic fragment 8A and probes. RNA probes used for RNase A protection experiments are labeled A, B, and C. The single-stranded DNA probe obtained by extension of primer EP1PE is labeled D. The black boxes show the deduced exons (E1 and E2). bp, Base pairs. (B) RNase A protection. Total and poly(A)<sup>+</sup> RNA samples were hybridized with the indicated RNA probes and digested with RNaseA (14). All reaction mixtures contained RNA purified from pericarps homozygous for *P-rr-4B2*, except for lane 13, which contained *P-ww-1112* RNA. Reaction mixtures were supplemented with tRNA to a total of 25  $\mu$ g of RNA, except lanes 2 and 4. Probe A was hybridized with 25  $\mu$ g of total RNA (lane 1), 10  $\mu$ g of total RNA (lane 2), 2  $\mu$ g of poly(A)<sup>+</sup> RNA (lane 6), 10  $\mu$ g tRNA (lane 4), or 25  $\mu$ g tRNA (lane 7). Probe B was hybridized with 25  $\mu$ g of total RNA (lane 3), 2  $\mu$ g of poly(A)<sup>+</sup> RNA (lane 9), or 25  $\mu$ g of tRNA (lane 10). Probe C was hybridized with 25  $\mu$ g of total RNA (lane 12), 25  $\mu$ g of total *P-ww* RNA (lane 13), or 25  $\mu$ g of tRNA (lane 14). Lanes 5, 8, and 11 correspond to undigested RNA probes A, B, and C, respectively. Arrows indicate the 450-, 200-, and 130-nucleotide (nt) protected fragments. (C) S1 nuclease mapping. Single-stranded DNA probe D was hybridized with  $\approx$ 2  $\mu$ g of *P-rr-4B2* poly(A)<sup>+</sup> RNA supplemented to 25  $\mu$ g with tRNA (lane 1) or 25  $\mu$ g of tRNA (lane 2), and treated with 300 units of S1 nuclease. Arrow shows the 69- to 70-nt protected fragment. The sequence obtained by using EP1PE as primer is shown at right. Asterisks indicate the deduced 5' end of the 450-nt exon.

**Homology of the *P*-Encoded Products to *myb*.** The first ATG is located at 320 bp from the transcription start and its sequence, GCGATGG, is in agreement with Kozak's consensus for translational initiation (18). The longer transcript encodes a 43.7-kDa polypeptide (Fig. 3A), whereas the shorter transcript encodes a 17.3-kDa protein (Fig. 3B). The N termini of the *P* translation products show striking homology ( $\approx$ 40% sequence identity) with the DNA-binding domain of several members of the *c-myb* family (15, 16). The homology is even greater (70%) with the *myb*-homologous domain of the protein encoded by *C1*, a maize gene that regulates anthocyanin biosynthesis in the kernel aleurone and embryo (17, 19) (Fig. 3C).

**Regulatory Effect of *P*.** *P* has been proposed to regulate the biosynthetic pathway leading to flavan-4-ol, the precursor of the red phlobaphene pigments (2). Since the *myb* motif has been implicated in transcriptional activation (20), we asked whether *P* expression is correlated with accumulation of RNA from genes encoding enzymes for flavonoid biosynthesis. The two genes tested were *C2* and *A1*: *C2* encodes chalcone synthase (21), the first enzyme in the pathway, and *A1* encodes an NADPH-dependent reductase (22) that converts flavanone into flavan-4-ol. Poly(A)<sup>+</sup> RNA from pericarps carrying different *P* alleles was tested by Northern hybridizations with *C2*, *A1*, and *P* probes (Fig. 5). Lanes 1 and 5 in each blot contain RNA from *P-rr-4B2* and *P-rr-4026*; these RNAs contain significant levels of message from *P*, *C2*, and *A1*. In contrast, lane 2 in each blot contains RNA from

the null allele *P-ww-1112*, which shows no detectable *P*, *C2*, or *A1* messages, although longer exposures show a very low level of *C2*-hybridizing RNA. Lanes 3 and 4 contain RNA from *P-vv* and *P-ovov-1114*, respectively; both alleles carry insertions of *Ac* within the large intron (Fig. 4). *P-vv* RNA contains low levels of *P*, *C2*, and *A1* message, while *P-ovov-1114* RNA contains moderate levels of *P* transcripts and significant levels of *C2* and *A1* RNA. These results indicate that accumulation of *A1* and *C2* RNA in the pericarp requires a functional *P* gene; thus, the role of *P* as a regulator is deduced not only from its homology to other transcriptional activators but also by a functional test.

## DISCUSSION

**Analysis of Two *P* Transcripts.** The *P* coding sequences are  $\approx$ 9 kb long and the start of transcription seems to be the same for all *P* transcripts. We have characterized two transcripts, 1802 and 945 nt long, as deduced from the amplified internal segments, the position of the 5' end determined by S1 mapping and primer extension, and the sequences of the 3' ends obtained from cDNA clones. The 5' region is composed of two exons of 452 bp and 130 bp, separated by an intron of 118 bp. The second intron is 4.9 kb and 7.0 kb for the longer and shorter transcripts, respectively, being the largest introns so far reported in the plant kingdom. Because of a direct sequence repetition at the *P* locus (Fig. 4; T.P., unpublished data), the 1802- and 945-nt transcripts have 318 nt in common

A

```

-59  gtacgtacgtacgtactccgtccgctgctatattatggccggccggtggcgtccctctct
1    AGCCAGCACAGCACACACTGGAAGTGCAGAGCTGTAGTGAGACCTGGCGACTGCCAG
61  CGTGTATCCGGCGGCAAGGAGCGTACGGCCGGTCTGGCCCGCACGGCCACCAACTCC
121 CTTGGACCGACGGCGACCACTGCTAAACCGTGCAGAGTAGTAGTGCGACTTCGCCGCC
181 GCGCCGGATCGTAGCTCGATCGATCGGCGGACCAATACGACTCCGGTGTGGCCAGCG
241 GCGCCCGGGCGGGAAACGACGCTGCTGGCGGAGCGAGGGCAGACGCTAGCTGTTGCC
301 GGGAGCTAGCCGGCCGGCGG  ATG GGG AGG ACG CCG TGC TGC GAG AAG GTG
      M G R T P C C E A K V

350  GGG CTC AAG CGA GGG AGG TGG ACG GCG GAA GAG GAC CAG TTA CTT
      G L K R G R W T A E E D Q L L

395  GCC AAC TAC ATT GCG GAG CAC GGC GAG GGG TCC TGG AGG TCG CTG
      A N Y I A E H G E G S W R S L

440  CCC AAG AAT GCA GGC CTG CTC CGG TGC GGC AAG AGC TGC CGG CTC
      P K N A G L L R C G K S C R L

485  CGG TGG ATC AAC TAC CTT CGG GCG GAC GTC AAG AGG GGG AAC ATC
      R W I N Y L R A D V K R G N I

530  TCC AAG GAG GAA GAA GAC ATC ATC ATC AAG CTC CAC GCC ACC CTC
      S K E E E D I I I K L H A T L

575  GGC AAC AGG TGG TCC CTG ATC GGC AGC CAC CTC CCC GGC CGA ACA
      G N R W S L I A S H L P G R T

620  GAC AAC GAG ATC AAG AAC TAC TGG AAC TCG CAC CTC AGC CGG CAG
      D N E I K N Y W N S H L S R Q

665  ATC CAC ACG TAC CGC CGG AAA TAC ACC GCC GGG CCT GAC GAC ACC
      I H T Y R R K Y T A G P D D T

710  GCC ATC GCC ATC GAC ATG AGC AAG CTG CAG AGC GCC GAC AGG CGG
      A I A I D M C K L Q S A D R R

755  CGC GGC GGC AGG ACC CCG GGC CGG CCG CCG AAG GCT AGC GCC AGC
      R G G R T P G R P P K A S A S

800  AGG ACC AAG CAG GCG GAC GCC GAT CAG CCC GGC GGC GAG GCG AAA
      R T K Q A D A D Q P G G E A K

845  GGC CCG GCC GCG GCG GCG TCG AGC CCG CGG CAC AGC GAC GTG GTG
      G P A A A A S S P R H S D D V V

890  AAC CCG GGC CCG AAC CAG CCC AAC AGC AGC AGC GGC AGC ACG GGC
      N P G P N Q P N S S S G S T G

935  ACG GCC GAG GAG GAG GGG CCC AGC AGC GAG GAC GCG AGC GGG CGG
      T A E E E G P S S E D A S G P

980  TGG GTG CTG GAG CCG ATA GAG CTC GGG GAC CTA GTC TGG GGG GAG
      W V L E P I E L G D L V W G E
      ↓
1025 GCC GAC AGC GAG ATG GAC GCC CTG ATG CCT ATC GGG CCC GGC GGC
      A D S E M D A L M P I G P G G

1070 ACG ACT CGG CTG CCC TCG AAG GGC TTG GCG CGG TCG GCT GCG AGG
      T T R L P S K G L A R S A A R

1115 CCC AGG TGG ACG ACC TGT TCG ACA TGG ACT GGG ATG GCT TCG CGG
      P R W T C T M T G M A S R

1160 CCC ATC TGT GGG GCG GGC CGG AGC AGG ACG AGC ACA GCG CGC AGC
      P I C G A G R S R T S T A R S

1205 TGC GGC AGG CCG CCG AGC CGC TGG AAG TTG CTG CTG CTG CTG
      C G R P P S R W K L L L L L L

1250 CTG CGA CCG CCG CCC GCA CCC CGG ACG ATC GCG AGC TGG AGG CGT
      L R R R P A P R T I A S W R R

1295 TCG AGA CTT GGC TCC TGT CCG ACT CGT TCT GAC GGC TCC GGT CAC
      S R L G S C P T R S D G G S G H

1340 CGG ACC GAT CAG ACA GAC CAA ATA ATT GGG TCA CGT GTG CTC GCT
      R T D Q T D Q I I G S R V L A

1385 CGC TCG CTG CCG TCG CGT GGG TCT TGG TTC AGA TGG CCA AAT AAT
      R S L P S R G S W F R W P N N

1430 TGG GAA AAA AAT TCT ACG GCC AGG GCC GTA AAG CCA CCA CCG TGC
      W E K N A R A R A V K P P C

1475 GCT CCT GAT GTC GAT GCC TGC CGC GTG GAG CTC TTG CGT ATC TAA
      A P D V D A C R V E L L R I

1520 CGCTCCACGACAAATCCCTTCCAGACGGCTCGAATTACATACGACAGGATCGGCTCCG
1580 CTTACTCCGTTCTGTTCGCTTCTGCTTTAGTGTGGTGCCTAGCAGATGCTGAGGCGGG
1640 TCGCCGCGCCCTCCGACGCTCGCCGCGCCGCTACGGGCGCTGCTGACGACGCCCTC
1700 CTCACGCGCTGAAAAGAGCTTGTATTTTACCTGTTTGTGTGCTTTTGTGCAATGGAA
1760 TAAACAATGATATTACTGAAATAACAATGAATGTTCTGAGAC
    
```

B

```

-59  gtacgtacgtacgtactccgtccgctgctatattatggccggccggtggcgtccctctct
1    AGCCAGCACAGCACACACTGGAAGTGCAGAGCTGTAGTGAGACCTGGCGACTGCCAG
61  CGTGTATCCGGCGGCAAGGAGCGTACGGCCGGTCTGGCCCGCACGGCCACCAACTCC
121 CTTGGACCGACGGCGACCACTGCTAAACCGTGCAGAGTAGTAGTGCGACTTCGCCGCC
181 GCGCCGGATCGTAGCTCGATCGATCGGCGGACCAATACGACTCCGGTGTGGCCAGCG
241 GCGCCCGGGCGGGAAACGACGCTGCTGGCGGAGCGAGGGCAGACGCTAGCTGTTGCC
301 GGGAGCTAGCCGGCCGGCGG  ATG GGG AGG ACG CCG TGC TGC GAG AAG GTG
      M G R T P C C E A K V

350  GGG CTC AAG CGA GGG AGG TGG ACG GCG GAA GAG GAC CAG TTA CTT
      G L K R G R W T A E E D Q L L

395  GCC AAC TAC ATT GCG GAG CAC GGC GAG GGG TCC TGG AGG TCG CTG
      A N Y I A E H G E G S W R S L

440  CCC AAG AAT GCA GGC CTG CTC CGG TGC GGC AAG AGC TGC CGG CTC
      P K N A G L L R C G K S C R L

485  CGG TGG ATC AAC TAC CTT CGG GCG GAC GTC AAG AGG GGG AAC ATC
      R W I N Y L R A D V K R G N I

530  TCC AAG GAG GAA GAA GAC ATC ATC ATC AAG CTC CAC GCC ACC CTC
      S K E E E D I I I K L H A T L

575  GGC AAC AGG CGC CAC CTG ATG ATG GAA GCG GAT TAC TCA CCG CCC
      G N R R R H L M I E A D Y S P P

620  TCG ACT GTT CGA TGC CTG CCG CGT GGA GCT CTT GCG TAT CTA ACG
      S T V R C L P R G A L A Y L T

665  CTC CCA CGA CAA TCA CCC TTC CAG ACG GCT CGA ATT ACA TAC GAC
      L P R Q S P F Q T A R I T Y D

710  AGG ATC GGC TCC GCT CTA CTC RG TCT GTT CGC TTC TGC TTT AFR
      R I G S A L L R C S V R G C C F R

755  TGC GTG CCT AGC AGA TGG TGAGCGCGGTGCGCGGGCCCTCCGACGGCTCGC
      C V P S R W

808  CGCCCGCGCTACGGGGCCTGCTGCAGCAGCCCTCCTCCAGCCTGAAAAGAGCTTTG
867  TATTACCTGTTTGTGCTTTGTGCAATGGAATAACAATGATATTACTGAA
926  AAACATGAATGTTCTGAGAC
    
```

C

123	QQRWAKVLNP	ELIKGPWTRD	EDDMVIKLVLR	NFGPKKWT		Dm.myb
33	QHRWQKVLNP	ELIKGPWTKK	EDDRVIELVQ	KYGEKRWV	.V	Hs.myb
1	MGRRACCPKK	GVKRGAWTSK	EDDALAAVVK	AHGEKRWRE		Zm.C1
1	MGRTPCCPKV	GLKRGRTIAE	EDQLLANVKA	EHGEGSWRS		Zm.P1
1	*****	*****	*****	*****		Zm.P2
162	IARYLNGRIG	KOCRERWLNH	LNPNIKKTAW	TEKEDETIYQ		Dm.myb
72	IAKHLKGRIG	KOCRERWLNH	LNPEVKKTSW	TEEDRKIYQ		Hs.myb
41	IQKAGLRGCG	KSCRLRWLNH	LREINIRGNI	SYDEEDLIIR		Zm.C1
41	PKNAGLRGCG	KSCRLRWLNH	LADYVIRGNI	SKEEEDLIIR		Zm.P1
41	*****	*****	*****	*****		Zm.P2
202	ARLELGNRWA	KIAKRLPGRT	DNAIKNHNS	TMRKRYDV		Dm.myb
112	AHKRLGNRWA	ETAKLLPGRT	DNAIKNHNS	TMRKRYEQ		Hs.myb
81	LHRTLGNRWS	LIAGRLPGRT	DNEIKNYWNS	TLGRRAGA		Zm.C1
81	LHATLGNRWS	LIASHLPGR	DNEIKNYWNS	HLSPQIHT		Zm.P1
81	*****RH	LMIEADYSPP	STVRCILPRGA	LAYITLPR		Zm.P2

Fig. 3. Nucleotide and amino acid sequences of the 1802-nt (A) and 945-nt (B) *P* gene transcripts. The transcription initiation site is labeled +1. The ends of the nucleotide sequences correspond to the site at which the oligo(dA) tails were located in about half of the cDNA clones, and the stars indicate the two alternative polyadenylation sites. The nucleotides flanking each intron are underlined (452–453 and 582–583). Arrow at position 1020 in A shows the site of the 7-bp insertion found in *P-wv-10:443-3*. (C) Comparison of amino acid sequences of the predicted proteins encoded by *Drosophila melanogaster myb* (Dm.myb) (15), human *c-myb* (Hs.myb) (16), maize *C1* (Zm.C1) (17), the 1802-nt *P* RNA (Zm.P1), and the 945-nt *P* RNA (Zm.P2). Shaded boxes indicate identical amino acids. Amino acids conserved in Dm.myb, Hs.myb, Zm.C1, and Zm.P1 are marked with a star.

at their 3' ends. However, the two encoded proteins would differ at their C termini due to a shift in the reading frame at the boundary between exons 2 and 3 (Fig. 3C).

The two large *P* transcripts of 6.5 and 7 kb are most likely unprocessed or incompletely processed transcripts, since they hybridize with probes spanning the length of the tran-

scribed region, including four probes from within the large intron (Fig. 1 and ref. 9). The 1802- and 945-nt transcripts probably represent the 2.0- and 1.4-kb mRNAs identified on Northern blots (Fig. 1). However, it is possible that the smaller cDNA corresponds to a 1-kb transcript previously reported (9), and that the 1.4-kb transcript has not been

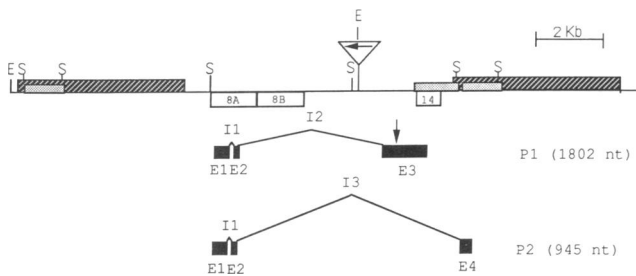


FIG. 4. Structure of *P* gene transcripts. The map of the *P* locus shows restriction sites for *EcoRI* (E) and *Sal I* (S). Hatched and stippled boxes indicate 5.2-kb and 1.2-kb direct repeats, respectively. The triangle containing the arrow marks the insertion and transcriptional orientation of the transposable element *Ac* in the *P-ovov-1114* allele; in *P-vv* *Ac* is inserted 161 bp 3' of the site in *P-ovov-1114*, in the opposite orientation (12). The numbered open boxes below the line indicate genomic fragments 8A, 8B, and 14 used as hybridization probes. Black boxes indicate the exons giving rise to the 1802-nt transcript (E1 + E2 + E3) or the 945-nt transcript (E1 + E2 + E4). Arrow above E3 indicates the position of the 7-bp insertion in *P-wv-10:443-3*. I1, intron 1.

cloned. In any case there may be additional transcripts arising from *P* that have yet to be characterized.

Which *P* transcript encodes a functional product? We have isolated a new *P-wv* allele (*P-wv-10:443-3*) that has a 7-bp insertion in the *P* gene corresponding to position 1020 in the 1802-nt transcript (arrow in Fig. 3A); the 7-bp insertion is a footprint remaining after excision of an *Ac* inserted at that site (see *Materials and Methods*). RNA from *P-wv-10:443-3* produces a normal pattern of *P* transcripts in Northern hybridizations (data not shown). The sequence containing the 7-bp footprint is not present in the 945-nt RNA, suggesting that the 1802-nt transcript is necessary for *P* function. However, we cannot rule out the possibility that the 7-bp insertion affects an as yet uncharacterized transcript that is required for *P* function. The product of the 945-nt transcript might act as a competitive inhibitor of the functional *P* protein, as proposed for an alternatively spliced product of the *c-myc* gene (24).

**Coding Potential of the *P* Transcripts.** The N terminus of the protein encoded by the 1802-nt transcript has about 40% identity with the DNA-binding domain of several members of the *c-myc* family of oncoproteins (Fig. 3C). The translated open reading frame also has a negatively charged region, amino acids 207–242, that overlaps a region, amino acids 221–244, predicted to be  $\alpha$ -helical by the Chou and Fasman algorithm (25). These characteristics are reportedly necessary for an activating domain (26). The hypothesis that *P* encodes a transcriptional activator is further supported by the finding that two structural genes of the flavonoid biosynthetic pathway (*A1* and *C2*) are expressed only in the presence of a functional *P* allele (Fig. 5).

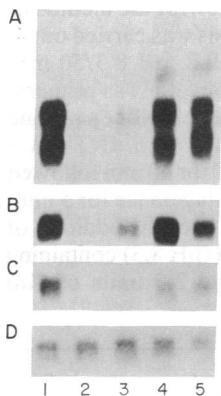


FIG. 5. Northern analysis of pericarp RNA from *P-rr-4B2* (lane 1), *P-wv-1112* (lane 2), *P-vv* (lane 3), *P-ovov-1114* (lane 4), and *P-rr-4026* (lane 5). The blots were hybridized with the following probes: *P* gene fragment 14 (Fig. 4) (A); a *C2* cDNA probe (21) (B); an *A1* cDNA probe (22) (C), and a maize actin probe (23) (D). RNA was prepared from pericarps dissected from kernels 21 days after pollination.

**Similarities of *P* and *C1*.** The N-terminal 118-amino acid *myb*-homologous domain of the predicted *P* proteins shows >70% identity with the product of the maize gene *C1*, which regulates anthocyanin biosynthesis in the kernel aleurone (17, 19). *P* and *C1* also share some striking similarities in gene structure. In both *P* and *C1*, the first intron is inserted between two guanine nucleotides of a glycine codon located 133 bp from the first ATG. The second exon is exactly equal in size in both genes. Although the structural similarity between the two genes is restricted to the first two exons, the homology between the encoded proteins extends into the third exon; this would rule out any easy association of homologous protein domains with gene intron/exon structure.

Since the biosynthetic pathways for anthocyanin and phlobaphene pigments have some steps in common, the four anthocyanin regulatory genes (*R*, *B*, *C1*, and *P1*) and *P* control an overlapping set of target genes, including *A1* and *C2*. Classic genetic (1) and molecular (27) experiments indicate that *C1* and its homolog *P1* require an active *R* or *B* allele for their function. The relatedness of *C1* and *P* described here raises the question of whether *P* also requires an accessory factor for transcription activation; no other required factor has been reported. Alternatively, *P* might encode a product that combines the transcription-activating functions of *C1/R* and *P1/B* in a single protein.

We thank Susan Allan for technical assistance and Rob Martienssen and Kim Arndt for helpful discussions and critical reading of the manuscript. We also thank R. Meagher for the maize actin probe, U. Wienand for the *C2* probe, and Z. Schwarz-Sommer for the *A1* probe. This research was supported in part by Pioneer Hi-Bred International, Inc., and National Institutes of Health Grant GM39832 to T.P.

- Coe, E. H., Jr., Neuffer, M. G. & Hoisington, D. A. (1988) in *Corn and Corn Improvement*, eds. Sprague, G. F. & Dudley, J. W. (Am. Soc. Agron./Crop Sci. Soc. Am./Soil Sci. Soc. Am., Madison, WI), pp. 81–258.
- Styles, E. D. & Ceska, O. (1977) *Can. J. Genet. Cytol.* **19**, 289–302.
- Miller, E. C. (1919) *J. Agric. Res.* **18**, 255–265.
- Ludwig, S. R. & Wessler, S. R. (1990) *Cell* **62**, 849–851.
- Chandler, V. L., Radicella, J. P., Robbins, T. P., Chen, J. & Turks, D. (1989) *Plant Cell* **1**, 1175–1183.
- Cone, K. C. & Burr, B. (1988) in *The Genetics of Flavonoids*, eds. Styles, D. E., Gavazzi, G. A. & Racchi, M. L. (Edizioni Unicopli, Milan), pp. 143–145.
- Emerson, R. A. (1917) *Genetics* **2**, 1–35.
- Barclay, P. C. & Brink, R. A. (1954) *Proc. Natl. Acad. Sci. USA* **40**, 1118–1126.
- Lechelt, C., Peterson, T., Laird, A., Chen, J., Dellaporta, S. L., Dennis, E., Peacock, W. J. & Starlinger, P. (1989) *Mol. Gen. Genet.* **219**, 225–234.
- Greenblatt, I. M. & Brink, R. A. (1963) *Nature (London)* **197**, 412–413.
- Greenblatt, I. M. (1984) *Genetics* **108**, 471–485.
- Peterson, T. (1990) *Genetics* **126**, 469–476.
- Athma, P. & Peterson, T. (1991) *Genetics*, in press.
- Sambrook, J., Fritsch, E. F. & Maniatis, T. (1989) *Molecular Cloning: A Laboratory Manual* (Cold Spring Harbor Lab., Cold Spring Harbor, NY).
- Peters, C. W. B., Sippel, A. E., Vingron, M. & Klempner, K.-H. (1987) *EMBO J.* **6**, 3085–3090.
- Slamon, D. J., Boone, T. C., Murdock, D. C., Keith, D. E., Press, M. F., Larson, R. A. & Souza, L. M. (1986) *Science* **233**, 347–351.
- Paz-Ares, J., Ghosal, D., Wienand, U., Peterson, P. A. & Saedler, H. (1987) *EMBO J.* **6**, 3553–3558.
- Kozak, M. (1984) *Nucleic Acids Res.* **12**, 857–872.
- Cone, K. C., Burr, F. A. & Burr, B. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 9631–9635.
- Dasgupta, P., Saikumar, P., Reddy, C. D. & Reddy, E. P. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8090–8094.
- Wienand, U., Weydemann, U., Niesbach-Klöggen, U., Peterson, P. A. & Saedler, H. (1986) *Mol. Gen. Genet.* **203**, 202–207.
- Schwarz-Sommer, Z., Shepherd, N., Tacke, E., Gierl, A., Rohde, W., Leclercq, L., Mattes, M., Berndtgen, R., Peterson, P. A. & Saedler, H. (1987) *EMBO J.* **6**, 287–294.
- Shah, D. M., Hightower, R. C. & Meagher, R. B. (1983) *J. Mol. Appl. Genet.* **2**, 111–126.
- Weber, B. L., Westin, E. H. & Clarke, M. F. (1990) *Science* **249**, 1291–1293.
- Chou, P. Y. & Fasman, G. D. (1974) *Biochemistry* **13**, 222–245.
- Ptashne, M. (1988) *Nature (London)* **335**, 683–689.
- Goff, S. A., Klein, T. M., Roth, B. A., Fromm, M. E., Cone, K. C., Radicella, J. P. & Chandler, V. L. (1990) *EMBO J.* **9**, 2517–2522.