



US00RE46033E

(19) **United States**  
(12) **Reissued Patent**  
**Boukharov et al.**

(10) **Patent Number:** **US RE46,033 E**  
(45) **Date of Reissued Patent:** **Jun. 21, 2016**

(54) **GENOMIC PLANT SEQUENCES AND USES THEREOF**

(71) Applicant: **Monsanto Technology LLC**, St. Louis, MO (US)

(72) Inventors: **Andrey A. Boukharov**, Chesterfield, MO (US); **Yongwei Cao**, Lexington, MA (US); **Stanton B. Dotson**, Chesterfield, MO (US); **Jeffrey M. Koshi**, Cambridge, MA (US); **David K. Kovalic**, University City, MO (US); **Jingdong Liu**, Ballwin, MO (US); **James D. McIninch**, Burlington, MA (US); **Wei Wu**, St. Louis, MO (US)

(73) Assignee: **Monsanto Technology LLC**, St. Louis, MO (US)

(21) Appl. No.: **14/121,743**

(22) Filed: **Oct. 10, 2014**

**Related U.S. Patent Documents**

Reissue of:

(64) Patent No.: **7,365,185**  
Issued: **Apr. 29, 2008**  
Appl. No.: **09/815,264**  
Filed: **Mar. 23, 2011**

U.S. Applications:

(63) Continuation-in-part of application No. 09/702,134, filed on Oct. 31, 2000, now abandoned, which is a continuation-in-part of application No. 09/620,392, filed on Jul. 19, 2000, now abandoned.

(60) Provisional application No. 60/144,351, filed on Jul. 20, 1999.

(51) **Int. Cl.**  
**C07H 21/04** (2006.01)  
**A01H 5/00** (2006.01)  
**C07H 21/02** (2006.01)  
**C07K 14/415** (2006.01)  
**C12N 15/82** (2006.01)

(52) **U.S. Cl.**  
CPC ..... **C07K 14/415** (2013.01); **C12N 15/8216** (2013.01); **C12N 15/8227** (2013.01)

(58) **Field of Classification Search**  
CPC ..... **C07K 14/415**; **C12N 15/8216**; **C12N 15/8227**  
See application file for complete search history.

(56) **References Cited**  
**PUBLICATIONS**

Adams et al., Complementary DNA sequencing: expressed sequence tags and human genome project, *Science*, 252(5013):1651-1656, 1991.  
Anaviev et al., "Oat-maize chromosome addition lines: a new system for mapping the maize genome," *Proc. Natl. Acad. Sci. USA*, 94:3524-3529, 1997.  
Birkenbihl et al., "Cosmid-derived map of *E. coli* strain BHE2600 in comparison to the map of strain W3110." *Nucleic Acids Res.*, 17(13):5057-5069, 1989.

Bukanov et al., "Ordered cosmid library and high-resolution physical-genetic map of helicobacter pylori strain NCTC11638," *Mol. Microbiol.*, 11(3):509-523, 1994.

Coulson et al., "Toward a physical map of the genome of the nematode caenorhabditis elegans," *Proc. Natl. Acad. Sci. USA*, 83:7821-8725, 1986.

Ebert et al., "Identification of an essential upstream element in the nopaline synthase promoter by stable and transient assays," *Proc. Natl. Acad. Sci. USA*, 84(16):5745-5749, 1987.

Efstratiadis et al., "Enzymatic in vitro synthesis of globin genes," *Cell*, 7:279-288, 1976.

Eiglmeier et al., "Use of an ordered cosmid library to deduce the genomic organization of *Mycobacterium leprae*," *Mol. Microbiol.*, 7(2):197-206, 1993.

Goff, "Rice as a model for cereal genomics," *Curr. Opin. Plant Biol.*, 2:86-89, 1999.

Hong, "A rapid and accurate strategy for rice contig map construction by combination of fingerprinting and hybridization," *Plant Mol. Biol.*, 35:129-133, 1997.

Kidwell et al., "Transposable elements as sources of variation in animals and plants," *Proc. Natl. Acad. Sci. USA*, 94:7704-7711, 1997.

Kim et al., "Construction and characterization of a human bacterial artificial chromosome library," *Genomics*, 34:213-218, 1996.

Knott et al., "Randomly picked cosmid clones overlap the pyrB and oriC gap in the physical map of the *E. coli* chromosome," *Nucleic Acids Res.*, 16:2601-2612, 1988.

Ko et al., "An 'equalized cDNA' library by the reassociation of short double-stranded cDNA," *Nucleic Acids Res.*, 18(19):5705-5711, 1990.

Kurata et al., "A 300 kilobase interval genetic map of rice including 883 expressed sequences," *Natur Gen.*, 8(4):362-372, 1994.

McCombie et al., "Caenorhabditis elegans expressed sequence tags identify gene families and disease gene homologues," *Nature Gen.*, 1:124-131, 1992.

Meinkoth et al., "Hybridization of nucleic acids immobilized on solid supports," *Anal. Biochem.*, 138:267-284, 1984.

Mohan et al., "Genome mapping, molecular markers and marker-assisted selection crop plants," *Mol. Breed.*, 3:87-103, 1997.

(Continued)

*Primary Examiner* — Bruce Campell

(74) *Attorney, Agent, or Firm* — Dentons US LLP; Carine M. Doyle Esq.

(57) **ABSTRACT**

The present invention discloses rice genomic promoter sequences. The promoters are particularly suited for use in rice and other cereal crops. Methods of modifying, producing, and using the promoters are also disclosed. The invention further discloses compositions, transformed host cells, transgenic plants, and seeds containing the rice genomic promoter sequences, and methods for preparing and using the same.

**15 Claims, No Drawings**

(56)

**References Cited**

## PUBLICATIONS

Okubo et al., "Large scale cDNA sequencing for analysis of quantitative and qualitative aspects of gene expression," *Nature Gen.*, 2:173-179, 1992.

Tanksley et al., "Chromosome landing: a paradigm for map-based gene cloning in plants with large genomes," *Trends in Genet.*, 11(2):63-68, 1995.

Venter et al., "A new strategy for genome sequencing," *Nature*, 381:364-366, 1996.

Wang et al., "Construction of a rice bacterial artificial chromosome library and identification of clones linked to the X-21 disease resistance locus," *Planta J.*, 7(3):525-533, 1995.

Wenzel et al., "Physical mapping of the mycoplasma pneumoniae genome," *Nucleic Acids Res.*, 16(17):8323-8336, 1988.

Yomo et al., "Histochemical studies on protease formation in the cotyledons of germinating bean seeds," *Planta*, 112(1):35-43, 1973.

Zhang et al., "Construction and characterization of two rice bacterial artificial chromosome libraries from the parents of a permanent recombinant inbred mapping population," *Mol. Breeding*, 2:11-24, 1996.

Zhang et al., "Physical mapping of the rice genome with BACs," *Plant Mol. Biol.*, 35:115-127, 1997.

Zwick et al., "Physical mapping of the liguleless linkage group in sorghum bicolor using rice RFLP-selected sorghum BACs," *Genetics*, 248:1983-1992, 1998.

Meinkoth et al. *Analyt. Biochem.* (1984) vol. 138, pp. 267-284.\*

Wing et al. NCBI accession number AZ134591, Jun. 2, 2000.\*

Chen et al., *Proc. Natl. Acad. Sci. USA*, 94:3431-3435 (1997).

\* cited by examiner

## GENOMIC PLANT SEQUENCES AND USES THEREOF

**Matter enclosed in heavy brackets [ ] appears in the original patent but forms no part of this reissue specification; matter printed in italics indicates the additions made by reissue; a claim printed with strikethrough indicates that the claim was canceled, disclaimed, or held invalid by a prior post-patent action or proceeding.**

### REFERENCES TO RELATED APPLICATIONS

This application is a continuation-in-part under 35 U.S.C. §120 of U.S. application Ser. No. 09/620,392, filed Jul. 19, 2000, *now abandoned, which claims the benefit of U.S. Provisional Application Ser. No. 60/144,351, filed Jul. 20, 1999; and is a continuation-in-part of U.S. application Ser. No. 09/702,134, filed Oct. 31, 2000, now abandoned*, the disclosures of which applications are incorporated herein by reference in their entirety.

### INCORPORATION OF SEQUENCE LISTING

Two copies of the sequence listing (*Copy 1 Replacement* and *Copy 2 Replacement*) and a computer readable form of the sequence listing (*Computer Readable Form (CRF) Replacement*), all on CD-ROMs, each containing the file named [Pa2\_00329.txt] *MONS232USRE seq.txt* which is [420,819,499] *420,834,519* bytes (measured in MS-DOS) and was created on [Mar. 23, 2001] *Sep. 26, 2014*, are herein incorporated by reference.

### INCORPORATION OF TABLES 1, 3, 4, 5 AND 6

Two copies of Tables 1, 3, 4, 5, and 6 on CD-ROMs, each containing 47,041,202 bytes (measured in MS-DOS) and each having the file name Pa\_00329.txt all created on [Mar. 16, 2001] *Apr. 15, 2009*, are herein incorporated by reference.

may be enhanced to have desirable agricultural, biosynthetic, commercial, chemical, insecticidal, industrial, nutritional, or pharmaceutical properties. Despite the availability of many molecular tools, however, the genetic modification of plants and seeds is often constrained by an insufficient or poorly localized expression of the engineered transgene.

Many intracellular processes may impact overall transgene expression, including transcription, translation, protein assembly and folding, methylation, phosphorylation, transport, and proteolysis. Intervention in one or more of these processes can increase the amount of transgene expression in genetically engineered plants and seeds. For example, raising the steady-state level of mRNA in the cytosol often yields an increased accumulation of transgene expression. Many factors may contribute to increasing the steady-state level of an mRNA in the cytosol, including the rate of transcription, promoter strength and other regulatory features of the promoter, efficiency of mRNA processing, and the overall stability of the mRNA.

Among these factors, the promoter plays a central role. Along the promoter, the transcription machinery is assembled and transcription is initiated. This early step is often rate-limiting relative to subsequent stages of protein production. Transcription initiation at the promoter may be regulated in several ways. For example, a promoter may be induced by the presence of a particular compound or external stimuli, express a gene only in a specific tissue, express a gene during a specific stage of development, or constitutively express a gene. Thus, transcription of a transgene may be regulated by operably linking the coding sequence to promoters with different regulatory characteristics. Accordingly, regulatory elements such as promoters, play a pivotal role in enhancing the agronomic, pharmaceutical or nutritional value of crops.

At least two types of information are useful in predicting promoter regions within a genomic DNA sequence. First, promoters may be identified on the basis of their sequence "content," such as transcription factor binding sites and various known promoter motifs. (Stormo, *Genome Research* 10:

### LENGTHY TABLES

The patent contains a lengthy table section. A copy of the table is available in electronic form from the USPTO web site ([http://seqdata.uspto.gov/?pageRequest=docDetail&DocID=\[US07365185B2\]USRE046033E1](http://seqdata.uspto.gov/?pageRequest=docDetail&DocID=[US07365185B2]USRE046033E1)). An electronic copy of the table will also be available from the USPTO upon request and payment of the fee set forth in 37 CFR 1.19(b)(3).

### FIELD OF THE INVENTION

The present invention relates to the field of plant biochemistry and genetics. Specifically, the invention relates to regulatory elements comprising genomic nucleic acid sequences from rice plants, and nucleic acid molecules containing the same. More specifically, the invention discloses nucleic acid sequences from *Oryza sativa* (rice) containing regulatory elements, such as promoters. The invention also discloses methods of modifying, producing, and using the regulatory elements.

### BACKGROUND OF THE INVENTION

#### Promoters

The genetic enhancement of plants and seeds provides significant benefits to society. For example, plants and seeds

394-397 (2000)). Such signals may be identified by computer programs that identify sites associated with promoters, such as TATA boxes, transcription factor (TF) binding sites, and CpG islands.

Second, promoters may be identified on the basis of their "location," i.e. their proximity to a known or suspected coding sequence. (Stormo, *Genome Research* 10: 394-397 (2000)). Promoters are typically contained within a region of DNA extending approximately 150-1500 basepairs in the 5' direction from the start codon of a coding sequence. Thus, promoter regions may be identified by locating the start codon of a coding sequence, and moving beyond the start codon in the 5' direction to locate the promoter region.

#### Rice

Approximately half a billion tons of rice is produced each year world-wide. More than 90% of this rice is for human consumption (Goff, *Curr. Opin. Plant Biol.* 2:86-89 (1999)).

Rice, however, is not only a commercially important crop; it is also a model for other cereal crops, such as sorghum, maize, barley and wheat.

Rice is a model crop for several reasons. First, the genes in rice are predicted to be generally arranged in the genome in an order that is similar to other cereal crops. In fact, comparisons of the physical and genetic maps of cereal genomes have suggested the existence of a colinearity of gene order among the various cereal genomes studied. (Goff, *Curr. Opin. Plant Biol.* 2:86-89 (1999)).

Second, studies of a number of individual genes indicate that there is considerable homology within gene families found among various cereal crops. This conservation of gene and protein sequences suggests that functional studies of genes or proteins from one cereal crop can help elucidate the function of similar genes or proteins in other cereal crops. Likewise, non-coding regulatory elements in rice, such as promoters, are predicted to display similar functions compared to related regulatory elements found in other cereal crops. Accordingly, a strong constitutive or tissue-specific promoters from one cereal is more likely to retain its function when introduced as a portion of a transgene into another cereal crop species (Goff, *Curr. Opin. Plant Biol.* 2:86-89 (1999)).

Third, rice can be used as a model for other cereal genomes because its genome is smaller than those of other major cereals. The size of the rice genome is estimated at 420 to 450 megabase pairs. Sorghum, maize, barley and wheat have larger genomes (1000, 3000, 5000 and 16000 Mbp, respectively). Despite such differences in genome size, however, the number of genes in each of these crops is on the same order of magnitude. Thus, the smaller genome size of rice results in a higher gene density relative to the other cereals. Based on estimates of 30,000 genes in a cereal genome, rice will have on average one gene approximately every 15 Kbp. In contrast, maize and wheat have one gene approximately every 100 and 500 Kbp, respectively. This higher gene density makes rice an attractive target for cereal gene discovery efforts, genomic sequence analysis, and identification of regulatory elements, such as promoters (Goff, *Curr. Opin. Plant Biol.* 2:86-89 (1999)).

For these reasons, rice is a model for other crops. Accordingly, discoveries in rice may be extended to other crops. Thus, the identification of new genes, regulatory elements (e.g., promoters), etc. that function in rice is useful not only in developing enhanced varieties of rice, but also in developing enhanced varieties of other crops. In particular, developments in rice are applicable to other cereal crops, such as sorghum, maize, barley and wheat.

Clearly, there exists a need in the art for new regulatory elements, such as promoters, that are capable of expressing heterologous nucleic acid sequences in important crop species.

#### SUMMARY OF THE INVENTION

The present invention includes and provides a substantially purified nucleic acid molecule comprising a nucleic acid sequence wherein the nucleic acid sequence: i) hybridizes under stringent conditions with a sequence selected from the group consisting of SEQ ID NO:1 through 57,467, and the complements thereof; or ii) exhibits an 85% or greater identity to a sequence selected from the group consisting of SEQ ID NO:1 through 57,467.

The present invention includes and provides a transgenic plant containing a nucleic acid molecule that comprises in the 5' to 3' direction: a nucleic acid sequence that: i) hybridizes

under stringent conditions with a sequence selected from the group consisting of SEQ ID NO:1 through 57,467, and the complements thereof; or ii) exhibits an 85% or greater identity to a sequence selected from the group consisting of SEQ ID NO:1 through 57,467; operably linked to a structural nucleic acid sequence; wherein the nucleic acid sequence is heterologous with respect to the structural nucleic acid sequence.

The present invention includes and provides a seed from a transgenic plant containing a nucleic acid molecule that comprises in the 5' to 3' direction: a nucleic acid sequence that: i) hybridizes under stringent conditions with a sequence selected from the group consisting of SEQ ID NO:1 through 57,467, and the complements thereof; or ii) exhibits an 85% or greater identity to a sequence selected from the group consisting of SEQ ID NO:1 through 57,467; operably linked to a structural nucleic acid sequence; wherein the nucleic acid sequence is heterologous with respect to the structural nucleic acid sequence.

The present invention includes and provides a fertile transgenic plant containing a nucleic acid molecule that comprises in the 5' to 3' direction: a nucleic acid sequence that: i) hybridizes under stringent conditions with a sequence selected from the group consisting of SEQ ID NO:1 through 57,467, and the complements thereof; or ii) exhibits an 85% or greater identity to a sequence selected from the group consisting of SEQ ID NO:1 through 57,467; operably linked to a structural nucleic acid sequence; wherein the nucleic acid sequence is heterologous with respect to the structural nucleic acid sequence.

The present invention includes and provides a method of transforming a host cell comprising: a) providing a nucleic acid molecule that comprises in the 5' to 3' direction: a nucleic acid sequence that: i) hybridizes under stringent conditions with a sequence selected from the group consisting of SEQ ID NO:1 through 57,467, and the complements thereof; or ii) exhibits an 85% or greater identity to a sequence selected from the group consisting of SEQ ID NO:1 through 57,467; operably linked to a structural nucleic acid sequence; and b) transforming said plant with the nucleic acid molecule.

#### DEFINITIONS

The following definitions are provided as an aid to understanding the detailed description of the present invention.

The phrases "coding sequence," "structural sequence," and "structural nucleic acid sequence" refer to a physical structure comprising an orderly arrangement of nucleic acids. The nucleic acids are arranged in a series of nucleic acid triplets that each form a codon. Each codon encodes for a specific amino acid. Thus the coding sequence, structural sequence, and structural nucleic acid sequence encode a series of amino acids forming a protein, polypeptide, or peptide sequence. The coding sequence, structural sequence, and structural nucleic acid sequence may be contained, without limitation, within a larger nucleic acid molecule, vector, etc. In addition, the orderly arrangement of nucleic acids in these sequences may be depicted, without limitation, in the form of a sequence listing, figure, table, electronic medium, etc.

The phrases "DNA sequence" and "nucleic acid sequence" refer to a physical structure comprising an orderly arrangement of nucleic acids. The DNA sequence or nucleic acid sequence may be contained within a larger nucleic acid molecule, vector, or the like. In addition, the orderly arrangement of nucleic acids in these sequences may be depicted in the form of a sequence listing, figure, table, electronic medium, or the like.

The term "expression" refers to the transcription of a gene to produce the corresponding mRNA and translation of this mRNA to produce the corresponding gene product (i.e., a peptide, polypeptide, or protein) and activity of the protein to confer a function.

The term "expression of antisense RNA" refers to the transcription of a DNA to produce a first RNA molecule capable of hybridizing to a second RNA molecule.

The term "gene" refers to chromosomal or genomic DNA, plasmid DNA, cDNA, synthetic DNA, or other DNA that encodes a peptide, polypeptide, protein, or RNA molecule.

"Homology" refers to the level of similarity between two or more nucleic acid or amino acid sequences in terms of percent of positional identity (i.e., sequence similarity or identity). Homology also refers to the concept of similar functional properties among different nucleic acids or proteins.

The phrase "heterologous" refers to the relationship between two or more nucleic acid or protein sequences that are derived from different sources. For example, a promoter is heterologous with respect to a coding sequence if such a combination is not normally found in nature. In addition, a particular sequence may be "heterologous" with respect to a cell or organism into which it is inserted (i.e. does not naturally occur in that particular cell or organism).

The term "hybridization" refers generally to the ability of nucleic acid molecules to join via complementary base strand pairing. Such hybridization may occur when nucleic acid molecules are contacted under appropriate conditions (see also, "specific hybridization," below).

The phrase "operably linked" refers to the functional spatial arrangement of two or more nucleic acid regions or nucleic acid sequences. For example, a promoter region may be positioned relative to a nucleic acid sequence such that transcription of the nucleic acid sequence is directed by the promoter region. Thus, the promoter region is "operably linked" to the nucleic acid sequence.

The term "promoter;" "promoter region;" or "promoter sequence" refer to a nucleic acid sequence, usually found upstream (5') to a coding sequence, that directs transcription of a nucleic acid sequence into mRNA. The promoter or promoter region typically provide a recognition site for RNA polymerase and the other factors necessary for proper initiation of transcription. As contemplated herein, a promoter or promoter region includes variations of promoters derived by inserting or deleting regulatory regions, subjecting the promoter to random or site-directed mutagenesis, etc. The activity or strength of a promoter may be measured in terms of the amounts of RNA it produces, or the amount of protein accumulation in a cell or tissue, relative to a promoter whose transcriptional activity has been previously assessed.

The term "recombinant vector" refers to any agent such as a plasmid, cosmid, virus, autonomously replicating sequence, phage, or linear or circular single-stranded or double-stranded DNA or RNA nucleotide sequence. The recombinant vector may be derived from any source; is capable of genomic integration or autonomous replication; and comprises a promoter nucleic acid sequence operably linked to one or more nucleic acid sequences. A recombinant vector is typically used to introduce such operably linked sequences into a suitable host.

"Regulatory sequence" refers to a nucleotide sequence located upstream (5'), within, or downstream (3') to a coding sequence. Transcription and expression of the coding sequence is typically impacted by the presence or absence of the regulatory sequence.

"Specifically hybridizes" refers to the ability of two nucleic acid molecules to form an anti-parallel, double-stranded nucleic acid structure. A nucleic acid molecule is said to be the "complement" of another nucleic acid molecule if they exhibit "complete complementarity;" i.e., each nucleotide in one sequence is complementary to its base pairing partner nucleotide in another sequence. Two molecules are said to be "minimally complementary" if they can hybridize to one another with sufficient stability to permit them to remain annealed to one another under at least conventional "low-stringency" conditions. Similarly, the molecules are said to be "complementary" if they can hybridize to one another with sufficient stability to permit them to remain annealed to one another under conventional "high-stringency" conditions. Nucleic acid molecules that hybridize to other nucleic acid molecules, e.g., at least under low stringency conditions are said to be "hybridizable cognates" of the other nucleic acid molecules. Conventional low stringency and high stringency conditions are described herein and by Sambrook et al., *Molecular Cloning, A Laboratory Manual*, 2nd Ed., Cold Spring Harbor Press, Cold Spring Harbor, N.Y. (1989) and by Haymes et al., *Nucleic Acid Hybridization, A Practical Approach*, IRL Press, Washington, D.C. (1985). Departures from complete complementarity are permissible, as long as such departures do not completely preclude the capacity of the molecules to form a double-stranded structure.

The term "substantially homologous" refers to two sequences which are at least 90% identical in sequence, as measured by the BestFit program described herein (Version 10; Genetics Computer Group, Inc., Madison, Wis.), using default parameters.

"Substantially purified" refers to a molecule separated from substantially all other molecules normally associated with it in its native state. More preferably a substantially purified molecule is the predominant species present in a preparation. A substantially purified molecule may be greater than 60% free, preferably 75% free, more preferably 90% free, and most preferably 95% free from the other molecules (exclusive of solvent) present in the natural mixture. The term "substantially purified" is not intended to encompass molecules present in their native state.

The term "transformation" refers to the introduction of nucleic acid into a recipient host. The term "host" refers to bacteria cells, fungi, animals and animal cells, plants and plant cells, or any plant parts or tissues including protoplasts, calli, roots, tubers, seeds, stems, leaves, seedlings, embryos, and pollen.

The term "transgenic" refers to an animal, plant, or other organism containing one or more heterologous nucleic acid sequences.

#### DETAILED DESCRIPTION OF THE INVENTION

The present invention includes nucleic acid molecules comprising promoter sequences useful for transcribing a heterologous structural nucleic acid sequence in plants, and methods of modifying, producing, and using the same. The invention also includes compositions, transformed host cells, transgenic plants, and seeds containing the promoters, and methods for preparing and using the same.

Nucleic Acid Molecules

The present invention includes a nucleic acid molecule having a nucleic acid sequence that hybridizes to SEQ ID NO:1 through SEQ ID NO:57,467, or any complements thereof; or any fragments thereof. The present invention also provides a nucleic acid molecule comprising a nucleic acid

sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, any complements thereof, and any fragments thereof.

Nucleic acid hybridization is a technique well known to those of skill in the art of DNA manipulation. The hybridization properties of a given pair of nucleic acids are an indication of their similarity or identity.

Low stringency conditions may be used to select nucleic acid sequences with lower sequence identities to a target nucleic acid sequence. One may wish to employ conditions such as about 0.15 M to about 0.9 M sodium chloride, at temperatures ranging from about 20° C. to about 55° C.

High stringency conditions may be used to select for nucleic acid sequences with higher degrees of identity to the disclosed nucleic acid sequences (Sambrook et al., 1989).

High stringency conditions typically involve nucleic acid hybridization in about 2× to about 10×SSC (diluted from a 20×SSC stock solution containing 3 M sodium chloride and 0.3 M sodium citrate, pH 7.0 in distilled water), about 2.5× to about 5×Denhardt's solution (diluted from a 50× stock solution containing 1% (w/v) bovine serum albumin, 1% (w/v) ficoll, and 1% (w/v) polyvinylpyrrolidone in distilled water), about 10 mg/mL to about 100 mg/mL fish sperm DNA, and about 0.02% (w/v) to about 0.1% (w/v) SDS, with an incubation at about 50° C. to about 70° C. for several hours to overnight. High stringency conditions are preferably provided by 6×SSC, 5×Denhardt's solution, 100 mg/mL fish sperm DNA, and 0.1% (w/v) SDS, with an incubation at 55° C. for several hours.

Hybridization is generally followed by several wash steps. The wash compositions generally comprise 0.5× to about 10×SSC, and 0.01% (w/v) to about 0.5% (w/v) SDS with a 15 minute incubation at about 20° C. to about 70° C. Preferably, the nucleic acid segments remain hybridized after washing at least one time in 0.1×SSC at 65° C.

A nucleic acid molecule preferably comprises a nucleic acid sequence that hybridizes, under low or high stringency conditions, with SEQ ID NO:1 through SEQ ID NO:57,467, any complements thereof, or any fragments thereof. A nucleic acid molecule most preferably comprises a nucleic acid sequence that hybridizes under high stringency conditions with SEQ ID NO:1 through SEQ ID NO:57,467, any complements thereof, or any fragments thereof.

In an alternative embodiment, the nucleic acid molecule comprises a nucleic acid sequence that exhibits 85% or greater identity, and more preferably at least 86 or greater, 87 or greater, 88 or greater, 89 or greater, 90 or greater, 91 or greater, 92 or greater, 93 or greater, 94 or greater, 95 or greater, 96 or greater, 97 or greater, 98 or greater, or 99% or greater identity to a nucleic acid molecule selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467 and complements thereof. The nucleic acid molecule most preferably comprises a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467 and complements thereof.

The percent of sequence identity is preferably determined using the "Best Fit" or "Gap" program of the Sequence Analysis Software Package™ (Version 10; Genetics Computer Group, Inc., Madison, Wis.). "Gap" utilizes the algorithm of Needleman and Wunsch (Needleman and Wunsch, *Journal of molecular Biology* 48:443-453, 1970) to find the alignment of two sequences that maximizes the number of matches and minimizes the number of gaps. "BestFit" performs an optimal alignment of the best segment of similarity between two sequences and inserts gaps to maximize the number of matches using the local homology algorithm of Smith and Waterman (Smith and Waterman, *Advances in*

*Applied Mathematics*, 2:482-489, 1981, Smith et al., *Nucleic Acids Research* 11:2205-2220, 1983). The percent identity is most preferably determined using the "Best Fit" program.

As used herein "sequence identity" refers to the extent to which two optimally aligned polynucleotide or peptide sequences are invariant throughout a window of alignment of components, e.g., nucleotides or amino acids. An "identity fraction" for aligned segments of a test sequence and a reference sequence is the number of identical components which are shared by the two aligned sequences divided by the total number of components in reference sequence segment, i.e., the entire reference sequence or a smaller defined part of the reference sequence. "Percent identity" is the identity fraction times 100.

Useful methods for determining sequence identity are also disclosed in *Guide to Huge Computers*, Martin J. Bishop, ed., Academic Press, San Diego, 1994, and Carillo, H., and Lip-ton, D., *Applied Math* (1988) 48:1073. More particularly, preferred computer programs for determining sequence identity include the Basic Local Alignment Search Tool (BLAST) programs which are publicly available from National Center Biotechnology Information (NCBI) at the National Library of Medicine, National Institute of Health, Bethesda, Md. 20894; see BLAST Manual, Altschul et al., NCBI, NLM, NIH; Altschul et al., *J. Mol. Biol.* 215:403-410 (1990); version 2.0 or higher of BLAST programs allows the introduction of gaps (deletions and insertions) into alignments; for peptide sequence BLASTX can be used to determine sequence identity; and, for polynucleotide sequence BLASTN can be used to determine sequence identity.

For purposes of this invention "percent identity" may also be determined using BLASTX version 2.0 for translated nucleotide sequences and BLASTN version 2.0 for polynucleotide sequences. In a preferred embodiment of the present invention, the presently disclosed rice genomic promoter sequences comprise nucleic acid molecules or fragments having a BLAST score of more than 200, preferably a BLAST score of more than 300, and even more preferably a BLAST score of more than 400 with their respective homologues.

Nucleic acid molecules of the present invention include nucleic acid sequences that are between about 0.01 Kb and about 50 Kb, more preferably between about 0.1 Kb and about 25 Kb, even more preferably between about 1 Kb and about 10 Kb, and most preferably between about 3 Kb and about 10 Kb, about 3 Kb and about 7 Kb, about 4 Kb and about 6 Kb, about 2 Kb and about 4 Kb, about 2 Kb and about 5 Kb, about 1 Kb and about 5 Kb, about 1 Kb and about 3 Kb, or about 1 Kb and about 2 Kb.

#### Promoters

Any of the nucleic acid molecules described herein may comprise nucleic acid sequences comprising promoters. Promoters of the present invention can include between about 300 bp upstream and about 10 kb upstream of the trinucleotide ATG sequence at the start site of a protein coding region. Promoters of the present invention can preferably include between about 300 bp upstream and about 5 kb upstream of the trinucleotide ATG sequence at the start site of a protein coding region. Promoters of the present invention can more preferably include between about 300 bp upstream and about 2 kb upstream of the trinucleotide ATG sequence at the start site of a protein coding region. Promoters of the present invention can include between about 300 bp upstream and about 1 kb upstream of the trinucleotide ATG sequence at the start site of a protein coding region. While in many circumstances a 300 bp promoter may be sufficient for expression,

additional sequences may act to further regulate expression, for example, in response to biochemical, developmental or environmental signals.

It is also preferred that the promoters of the present invention contain a CAAT and a TATA cis element. Moreover, the promoters of the present invention can contain one or more cis elements in addition to a CAAT and a TATA box.

By "regulatory element" it is intended a series of nucleotides that determines if, when, and at what level a particular gene is expressed. The regulatory DNA sequences specifically interact with regulatory or other proteins. Many regulatory elements act in cis ("cis elements") and are believed to affect DNA topology, producing local conformations that selectively allow or restrict access of RNA polymerase to the DNA template or that facilitate selective opening of the double helix at the site of transcriptional initiation. Cis elements occur within, but are not limited to promoters, and promoter modulating sequences (inducible elements). Cis elements can be identified using known cis elements as a target sequence or target motif in the BLAST programs of the present invention.

Promoters of the present invention include homologues of cis elements known to effect gene regulation that show homology with the promoter sequences of the present invention. These cis elements include, but are not limited to, oxygen responsive cis elements (Cowen et al., *J Biol. Chem.* 268(36):26904-26910 (1993)), light regulatory elements (Bruce and Quail, *Plant Cell* 2 (11):1081-1089 (1990); Bruce et al., *EMBO J.* 10:3015-3024 (1991); Rocholl et al., *Plant Sci.* 97:189-198 (1994); Block et al., *Proc. Natl. Acad. Sci. USA* 87:5387-5391 (1990); Giuliano et al., *Proc. Natl. Acad. Sci. USA* 85:7089-7093 (1988); Staiger et al., *Proc. Natl. Acad. Sci. USA* 86:6930-6934 (1989); Izawa et al., *Plant Cell* 6: 1277-1287 (1994); Menkens et al., *Trends in Biochemistry* 20:506-510 (1995); Foster et al., *FASEB J.* 8:192-200 (1994); Plesse et al., *Mol Gen Gene* 254:258-266 (1997); Green et al., *EMBO J.* 6:2543-2549(1987); Kuhlemeier et al., *Ann. Rev. Plant Physiol.* 38:221-257 (1987); Villain et al., *J. Biol. Chem.* 271:32593-32598 (1996); Lam et al., *Plant Cell* 2:857-866 (1990); Gilmartin et al., *Plant Cell* 2:369-378 (1990); Datta et al., *Plant Cell* 1:1069-1077 (1989); Gilmartin et al., *Plant Cell* 2:369-378 (1990); Castresana et al., *EMBO J.* 7:1929-1936 (1988); Ueda et al., *Plant Cell* 1:217-227 (1989); Terzaghi et al., *Annu. Rev. Plant Physiol. Plant Mol. Biol.* 46:445-474 (1995); Green et al., *EMBO J.* 6:2543-2549 (1987); Villain et al., *J. Biol. Chem.* 271:32593-32598 (1996); Tjaden et al., *Plant Cell* 6: 107-118 (1994); Tjaden et al., *Plant Physiol.* 108:1109-1117 (1995); Ngai et al., *Plant J.* 12:1021-1234 (1997); Bruce et al., *EMBO J* 10:3015-3024 (1991); Ngai et al., *Plant J.* 12:1021-1034 (1997)), elements responsive to gibberellin, (Muller et al., *J. Plant Physiol.* 145:606-613 (1995); Croissant et al., *Plant Science* 116:27-35 (1996); Lohmer et al., *EMBO J.* 10:617-624 (1991); Rogers et al., *Plant Cell* 4:1443-1451 (1992); Lanahan et al., *Plant Cell* 4:203-211 (1992); Skriver et al., *Proc. Natl. Acad. Sci. USA* 88:7266-7270 (1991); Gilmartin et al., *Plant Cell* 2:369-378 (1990); Huang et al., *Plant Mol. Biol.* 14:655-668 (1990), Gubler et al., *Plant Cell* 7:1879-1891 (1995)), elements responsive to abscisic acid, (Busk et al., *Plant Cell* 9:2261-2270 (1997); Guiltinan et al., *Science* 250:267-270 (1990); Shen et al., *Plant Cell* 7:295-307 (1995); Shen et al., *Plant Cell* 8:1107-1119 (1996); Seo et al., *Plant Mol. Biol.* 27:1119-1131 (1995); Marcotte et al., *Plant Cell* 1:969-976 (1989); Shen et al., *Plant Cell* 7:295-307 (1995); Iwasaki et al., *Mol Gen Genet* 247:391-398 (1995); Hattori et al., *Genes Dev.* 6:609-618 (1992); Thomas et al., *Plant Cell* 5:1401-1410 (1993)), elements similar to abscisic acid responsive

elements, (Ellerstrom et al., *Plant Mol. Biol.* 32:1019-1027 (1996)), auxin responsive elements (Liu et al., *Plant Cell* 6:645-657 (1994); Liu et al., *Plant Physiol.* 115:397-407 (1997); Kosugi et al., *Plant J.* 7:877-886 (1995); Kosugi et al., *Plant Cell* 9:1607-1619 (1997); Ballas et al., *J. Mol. Biol.* 233.580-596 (1993)), a cis element responsive to methyl jasmonate treatment (Beaudoin and Rotstein *Plant Mol. Biol.* 33:835-846 (1997)), a cis element responsive to abscisic acid and stress response (Straub et al., *Plant Mol. Biol.* 26:617-630 (1994)), ethylene responsive cis elements (Itzhaki et al., *Proc. Natl. Acad. Sci. USA* 91:8925-8929 (1994); Montgomery et al., *Proc. Natl. Acad. Sci. USA* 90:5939-5943 (1993); Sessa et al., *Plant Mol. Biol.* 28:145-153 (1995); Shinsai et al., *Plant Mol. Biol.* 27:923-932 (1995)), salicylic acid cis responsive elements, (Strange et al., *Plant J.* 11:1315-1324 (1997); Qin et al., *Plant Cell* 6:863-874 (1994)), a cis element that responds to water stress and abscisic acid (Lam et al., *J. Biol. Chem.* 266:17131-17135 (1991); Thomas et al., *Plant Cell* 5:1401-1410 (1993); Pla et al., *Plant Mol Biol* 21:259-266 (1993)), a cis element essential for M phase-specific expression (Ito et al., *Plant Cell* 10:331-341 (1998)), sucrose responsive elements (Huang et al., *Plant Mol. Biol.* 14:655-668 (1990); Hwang et al., *Plant Mol Biol* 36:331-341 (1998); Grierson et al., *Plant J.* 5:815-826 (1994)), heat shock response elements (Pelham et al., *Trends Genet.* 1:31-35 (1985)), elements responsive to auxin and/or salicylic acid and also reported for light regulation (Lam et al., *Proc. Natl. Acad. Sci. USA* 86:7890-7897 (1989); Benfey et al., *Science* 250:959-966 (1990)), elements responsive to ethylene and salicylic acid (Ohme-Takagi et al., *Plant Mol. Biol.* 15:941-946 (1990)), elements responsive to wounding and abiotic stress (Loake et al., *Proc. Natl. Acad. Sci. USA* 89:9230-9234 (1992); Mhiri et al., *Plant Mol. Biol.* 33:257-266 (1997)), antioxidant response elements (Rushmore et al., *J. Biol. Chem.* 266: 11632-11639; Dalton et al., *Nucleic Acids Res.* 22:5016-5023 (1994)), Sph elements (Suzuki et al., *Plant Cell* 9:799-807 (1997)), elicitor responsive elements, (Fukuda et al., *Plant Mol. Biol.* 34:81-87 (1997); Rushton et al., *EMBO J.* 15:5690-5700 (1996)), metal responsive elements (Stuart et al., *Nature* 317:828-831 (1985); Westin et al., *EMBO J.* 7:3763-3770 (1988); Thiele et al., *Nucleic Acids Res.* 20:1183-1191 (1992); Faisst et al., *Nucleic Acids Res.* 20:3-26 (1992)), low temperature responsive elements, (Baker et al., *Plant Mol. Biol.* 24:701-713 (1994); Jiang et al., *Plant Mol. Biol.* 30:679-684 (1996); Nordin et al., *Plant Mol. Biol.* 21:641-653 (1993); Zhou et al., *J. Biol. Chem.* 267:23515-23519 (1992)), drought responsive elements, (Yamaguchi et al., *Plant Cell* 6:251-264 (1994); Wang et al., *Plant Mol. Biol.* 28:605-617 (1995); Bray E A, *Trends in Plant Science* 2:48-54 (1997)) enhancer elements for glutenin, (Colot et al., *EMBO J.* 6:3559-3564 (1987); Thomas et al., *Plant Cell* 2:1171-1180 (1990); Kreis et al., *Philos. Trans. R. Soc. Lond.*, B314:355-365 (1986)), light-independent regulatory elements, (Lagrange et al., *Plant Cell* 9:1469-1479 (1997); Villain et al., *J. Biol. Chem.* 271:32593-32598(1996)), OCS enhancer elements, (Bouchez et al., *EMBO J.* 8:4197-4204 (1989); Foley et al., *Plant J.* 3:669-679 (1993)), ACGT elements, (Foster et al., *FASEB J.* 8:192-200 (1994); Izawa et al., *Plant Cell* 6:1277-1287 (1994); Izawa et al., *J. Mol. Biol.* 230:1131-1144 (1993)), negative cis elements in plastid related genes, (Zhou et al., *J. Biol. Chem.* 267:23515-23519 (1992); Lagrange et al., *Mol. Cell Biol.* 13:2614-2622 (1993); Lagrange et al., *Plant Cell* 9:1469-1479 (1997); Zhou et al., *J. Biol. Chem.* 267:23515-23519 (1992)), prolamin box elements, (Forde et al., *Nucleic Acids Res.* 13:7327-7339 (1985); Colot et al., *EMBO J.* 6:3559-3564 (1987); Thomas et al., *Plant Cell* 2:1171-1180 (1990); Thompson et al., *Plant*

Mol. Biol. 15:755-764 (1990); Vicente et al., Proc. Natl. Acad. Sci. USA 94:7685-7690 (1997)), elements in enhancers from the IgM heavy chain gene (Gillies et al., Cell 33:717-728 (1983); Whittier et al., Nucleic Acids Res. 15:2515-2535 (1987)).

#### Promoter Activity

The activity or strength of a promoter may be measured in terms of the amount of mRNA or protein accumulation it specifically produces, relative to the total amount of mRNA or protein. The promoter preferably expresses an operably linked nucleic acid sequence at a level greater than 0.01%; more preferably greater than 0.05, 0.1, 0.25, 0.5, 0.75, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, or 20% (w/w) of the total cellular RNA or protein.

As used herein, an "expression pattern" is any pattern of differential gene expression. In a preferred embodiment, an expression pattern is selected from the group consisting of tissue, temporal, spatial, developmental, stress, environmental, physiological, pathological, cell cycle, and chemically responsive expression patterns.

As used herein, an "enhanced expression pattern" is any expression pattern for which an operably linked nucleic acid sequence is expressed at a level greater than 0.01%; more preferably greater than 0.05, 0.1, 0.25, 0.5, 0.75, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, or 20%(w/w) of the total cellular RNA or protein.

Alternatively, the activity or strength of a promoter may be expressed relative to a well-characterized promoter (for which transcriptional activity was previously assessed). For example, a less-characterized promoter may be operably linked to a reporter sequence (e.g., GUS) and introduced into a specific cell type. A well-characterized promoter (e.g. the 35S promoter) is similarly prepared and introduced into the same cellular context. Transcriptional activity of the unknown promoter is determined by comparing the amount of reporter expression, relative to the well characterized promoter. In one embodiment, the activity of the present promoter is as strong as the 35S promoter when compared in the same cellular context. The cellular context is preferably rice, sorghum, maize, barley, wheat, canola, soybean, or maize; and more preferably is rice, sorghum, maize, barley, or wheat; and most preferably is rice.

#### Structural Nucleic Acid Sequences

The promoter of the present invention may be operably linked to a structural nucleic acid sequence that is heterologous with respect to the promoter. The structural nucleic acid sequence may generally be any nucleic acid sequence for which an increased level of transcription is desired. The structural nucleic acid sequence preferably encodes a polypeptide that is suitable for incorporation into the diet of a human or an animal. Suitable structural nucleic acid sequences include those encoding a yield protein, a stress resistance protein, a developmental control protein, a tissue differentiation protein, a meristem protein, an environmentally responsive protein, a senescence protein, a hormone responsive protein, an abscission protein, a source protein, a sink protein, a flower control protein, a seed protein, an herbicide resistance protein, a disease resistance protein, a fatty acid biosynthetic enzyme, a tocopherol biosynthetic enzyme, an amino acid biosynthetic enzyme, and an insecticidal protein.

Alternatively, the promoter and structural nucleic acid sequence may be designed to down-regulate a specific nucleic acid sequence. This is typically accomplished by linking the promoter to a structural nucleic acid sequence that is oriented in the antisense direction. One of ordinary skill in the art is familiar with such antisense technology. Briefly, as the antisense nucleic acid sequence is transcribed, it hybridizes to

and sequesters a complimentary nucleic acid sequence inside the cell. This duplex RNA molecule cannot be translated into a protein by the cell's translational machinery. Any nucleic acid sequence may be negatively regulated in this manner.

#### 5 Modified Structural Nucleic Acid Sequences

The promoter of the present invention may also be operably linked to a modified structural nucleic acid sequence that is heterologous with respect to the promoter. The structural nucleic acid sequence may be modified to provide various desirable features. For example, a structural nucleic acid sequence may be modified to increase the content of essential amino acids, enhance translation of the amino acid sequence, alter post-translational modifications (e.g., phosphorylation sites), transport a translated product to a compartment inside or outside of the cell, improve protein stability, insert or delete cell signaling motifs, etc.

#### Codon Usage in Structural Nucleic Acid Sequences

Due to the degeneracy of the genetic code, different nucleotide codons may be used to code for a particular amino acid. A host cell often displays a preferred pattern of codon usage. Structural nucleic acid sequences are preferably constructed to utilize the codon usage pattern of the particular host cell. This generally enhances the expression of the structural nucleic acid sequence in a transformed host cell. Any of the above described nucleic acid and amino acid sequences may be modified to reflect the preferred codon usage of a host cell or organism in which they are contained. Modification of a structural nucleic acid sequence for optimal codon usage in plants is described in U.S. Pat. No. 5,689,052.

#### Other Modifications of Structural Nucleic Acid Sequences

Additional variations in the structural nucleic acid sequences described above may encode proteins having equivalent or superior characteristics when compared to the proteins from which they are engineered. Mutations may include deletions, insertions, truncations, substitutions, fusions, shuffling of motif sequences, and the like.

Mutations to a structural nucleic acid sequence may be introduced in either a specific or random manner, both of which are well known to those of skill in the art of molecular biology. A myriad of site-directed mutagenesis techniques exist, typically using oligonucleotides to introduce mutations at specific locations in a structural nucleic acid sequence. Examples include single strand rescue (Kunkel et al., Proc. Natl. Acad. Sci. U.S.A., 82: 488-492, 1985), unique site elimination (Deng and Nickloff, Anal. Biochem. 200:81, 1992), nick protection (Vandeyar, et al. Gene 65: 129-133, 1988), and PCR (Costa et al., Methods Mol. Biol. 57: 31-44, 1996). Random or non-specific mutations may be generated by chemical agents (for a general review, see Singer and Kusmierk, Ann. Rev. Biochem. 52: 655-693, 1982) such as nitrosoguanidine (Cerdeira-Olmedo et al., J. Mol. Biol. 33: 705-719, 1968; Guerola, et al. Nature New Biol. 230: 122-125, 1971) and 2-aminopurine (Rogan and Bessman, J. Bacteriol. 103: 622-633, 1970); or by biological methods such as passage through mutator strains (Greener, et al. Mol. Biotechnol. 7:189-195,1997).

The modifications may result in either conservative or non-conservative changes in the amino acid sequence. Conservative changes result from additions, deletions, substitutions, etc. in the structural nucleic acid sequence which do not alter the final amino acid sequence of the protein. In a preferred embodiment, the protein has between 20 and 500 conservative changes, more preferably between 15 and 300 conservative changes, even more preferably between 10 and 150 conservative changes, and most preferably between 5 and 75 conservative changes.

Non-conservative changes include additions, deletions, and substitutions which result in an altered amino acid sequence. In a preferred embodiment, the protein has between 10 and 250 non-conservative amino acid changes, more preferably between 5 and 100 non-conservative amino acid changes, even more preferably between 2 and 50 non-conservative amino acid changes, and most preferably between 1 and 30 non-conservative amino acid changes.

Additional methods of making the alterations described above are described by Ausubel et al., *Current Protocols in Molecular Biology*, John Wiley and Sons, Inc., 1995, Bauer et al., *Gene*, 37:73, 1985; Craik, *BioTechniques*, 3: 12-19, 1985; Frits Eckstein et al., *Nucleic Acids Research*, 10: 6487-6497, 1982; Sambrook, et al., *Molecular Cloning: A Laboratory Manual*, Second Edition, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y., 1989, Smith, et al., In: *Genetic Engineering: Principles and Methods*, Setlow et al., Eds., Plenum Press, N.Y., 1-32, 1981, and Osuna, et al., *Critical Reviews In Microbiology*, 20: 107-116, 1994.

Modifications may be made to the protein sequences of the present invention and the nucleic acid segments which encode them that maintain the desired properties of the molecule. The following is a discussion based upon changing the amino acid sequence of a protein to create an equivalent, or possibly an improved, second-generation molecule. The amino acid changes may be achieved by changing the codons of the structural nucleic acid sequence, according to the codons given in Table A.

TABLE A

Codon degeneracy of amino acids			
Amino acid	One letter	Three letter	Codons
Alanine	A	Ala	GCA GCC GCG GCT
Cysteine	C	Cys	TGC TGT
Aspartic acid	D	Asp	GAC GAT
Glutamic acid	E	Glu	GAA GAG
Phenylalanine	F	Phe	ITC TTT
Glycine	G	Gly	GGA GGC GGG GGT
Histidine	H	His	CAC CAT
Isoleucine	I	Ile	ATA ATC ATT
Lysine	K	Lys	AAA AAG
Leucine	L	Leu	TTA TTG CTA CTC CTG CTT
Methionine	M	Met	ATG
Asparagine	N	Asn	AAC AAT
Proline	P	Pro	CCA CCC CCG CCT
Glutamine	Q	Gln	CAA CAG
Arginine	R	Arg	AGA AGG CGA CGC CGG CGT
Serine	S	Ser	AGC AGT TCA TCC TCG TCT
Threonine	T	Thr	ACA ACC ACG ACT
Valine	V	Val	GTA GTC GTG GTT
Tryptophan	W	Trp	TGG
Tyrosine	Y	Tyr	TAC TAT

In making such changes, the hydropathic index of amino acids may be considered. The importance of the hydropathic amino acid index in conferring interactive biological function on a protein is generally understood in the art (Kyte and Doolittle, *J. Mol. Biol.*, 157: 105-132, 1982). It is accepted that the relative hydropathic character of the amino acid contributes to the secondary structure of the resultant protein, which in turn defines the interaction of the protein with other molecules, for example, enzymes, substrates, receptors, DNA, antibodies, antigens, and the like.

Each amino acid has been assigned a hydropathic index on the basis of their hydrophobicity and charge characteristics. These are: isoleucine (+4.5); valine (+4.2); leucine (+3.8); phenylalanine (+2.8); cysteine/cysteine (+2.5); methionine

(+1.9); alanine (+1.8); glycine (-0.4); threonine (-0.7); serine (-0.8); tryptophan (-0.9); tyrosine (-1.3); proline (-1.6); histidine (-3.2); glutamate/glutamine/aspartate/asparagine (-3.5); lysine (-3.9); and arginine (4.5).

It is known in the art that certain amino acids may be substituted by other amino acids having a similar hydropathic index or score and still result in a protein with similar biological activity, i.e., still obtain a biologically functional protein. In making such changes, the substitution of amino acids whose hydropathic indices are within  $\pm 2$  is preferred, those within  $\pm 1$  are more preferred, and those within  $\pm 0.5$  are most preferred.

It is also understood in the art that the substitution of like amino acids may be made effectively on the basis of hydrophobicity. U.S. Pat. No. 4,554,101 (Hopp) states that the greatest local average hydrophobicity of a protein, as governed by the hydrophobicity of its adjacent amino acids, correlates with a biological property of the protein. The following hydrophobicity values have been assigned to amino acids: arginine/lysine (+3.0); aspartate/glutamate (+3.0 $\pm$ 1); serine (+0.3); asparagine/glutamine (+0.2); glycine (0); threonine (-0.4); proline (-0.5 $\pm$ 1); alanine/histidine (-0.5); cysteine (-1.0); methionine (-1.3); valine (-1.5); leucine/isoleucine (-1.8); tyrosine (-2.3); phenylalanine (-2.5); and tryptophan (-3.4).

It is understood that an amino acid may be substituted by another amino acid having a similar hydrophobicity score and still result in a protein with similar biological activity, i.e., still obtain a biologically functional protein. In making such changes, the substitution of amino acids whose hydropathic indices are within  $\pm 2$  is preferred, those within  $\pm 1$  are more preferred, and those within  $\pm 0.5$  are most preferred.

As outlined above, amino acid substitutions are therefore based on the relative similarity of the amino acid side-chain substituents, for example, their hydrophobicity, hydrophobicity, charge, size, and the like. Exemplary substitutions which take various of the foregoing characteristics into consideration are well known to those of skill in the art and include: arginine and lysine; glutamate and aspartate; serine and threonine; glutamine and asparagine; and valine, leucine, and isoleucine. Changes which are not expected to be advantageous may also be used if these resulted proteins having improved rumen resistance, increased resistance to proteolytic degradation, or both improved rumen resistance and increased resistance to proteolytic degradation, relative to the unmodified polypeptide from which they are engineered.

#### Recombinant Vectors

Any of the promoters and structural nucleic acid sequences described above may be provided in a recombinant vector. A recombinant vector typically comprises, in a 5' to 3' orientation: a promoter to direct the transcription of a structural nucleic acid sequence and a structural nucleic acid sequence. The recombinant vector may further comprise a 3' transcriptional terminator, a 3' polyadenylation signal, other untranslated nucleic acid sequences, transit and targeting nucleic acid sequences, selectable markers, enhancers, and operators, as desired.

Means for preparing recombinant vectors are well known in the art. Methods for making recombinant vectors particularly suited to plant transformation include, without limitation, those described in U.S. Pat. Nos. 4,971,908, 4,940,835, 4,769,061 and 4,757,011. These type of vectors have also been reviewed (Rodriguez, et al. *Vectors: A Survey of Molecular Cloning Vectors and Their Uses*, Butterworths, Boston, 1988; Glick et al., *Methods in Plant Molecular Biology and Biotechnology*, CRC Press, Boca Raton, Fla., 1993).

Typical vectors useful for expression of nucleic acids in higher plants are well known in the art and include vectors

derived from the tumor-inducing (Ti) plasmid of *Agrobacterium tumefaciens* (Rogers, et al., Meth. In Enzymol, 153: 253-277, 1987). Other recombinant vectors useful for plant transformation, including the pCaMVCN transfer control vector, have also been described (Fromm et al., Proc. Natl. Acad. Sci. USA, 82(17): 5824-5828, 1985).

#### Promoters in the Recombinant Vectors

The promoter used in the recombinant vector preferably transcribes a heterologous structural nucleic acid sequence at a high level in a plant. More preferably, the promoter hybridizes to a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, or any complements thereof; or any fragments thereof. Suitable hybridization conditions include those described above. A nucleic acid sequence of the promoter preferably hybridizes, under low or high stringency conditions, with SEQ ID NO:1 through SEQ ID NO:57,467, and any complements thereof. The promoter most preferably hybridizes under high stringency conditions to a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, and any complements thereof.

In an alternative embodiment, the promoter comprises a nucleic acid sequence that exhibits 85% or greater identity, and more preferably at least 86 or greater, 87 or greater, 88 or greater, 89 or greater, 90 or greater, 91 or greater, 92 or greater, 93 or greater, 94 or greater, 95 or greater, 96 or greater, 97 or greater, 98 or greater, or 99% or greater identity to a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, and complements thereof. The promoter most preferably comprises a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, any complements thereof, and any fragments thereof.

#### Additional Promoters in the Recombinant Vector

One or more additional promoters may also be provided in the recombinant vector. These promoters may be operably linked to any of the structural nucleic acid sequences described above. Alternatively, the promoters may be operably linked to other nucleic acid sequences, such as those encoding transit peptides, selectable marker proteins, or antisense sequences.

These additional promoters may be selected on the basis of the cell type into which the vector will be inserted. Promoters which function in bacteria, yeast, and plants are all well taught in the art. The additional promoters may also be selected on the basis of their regulatory features. Examples of such features include enhancement of transcriptional activity, inducibility, tissue-specificity, and developmental stage-specificity. In plants, promoters that are inducible, of viral or synthetic origin, constitutively active, temporally regulated, and spatially regulated have been described (Poszkowski, et al., EMBO J., 3: 2719, 1989; Odell, et al., Nature, 313:810, 1985; Chau et al., Science, 244:174-181, 1989).

Often-used constitutive promoters include the CaMV 35S promoter (Odell, et al., Nature, 313: 810, 1985), the enhanced CaMV 35S promoter, the Figwort Mosaic Virus (FMV) promoter (Richins, et al., Nucleic Acids Res. 20: 8451, 1987), the mannopine synthase (mas) promoter, the nopaline synthase (nos) promoter, and the octopine synthase (ocs) promoter.

Useful inducible promoters include promoters induced by salicylic acid or polyacrylic acids (PR-1; Williams, et al., Biotechnology 10:540-543, 1992), induced by application of safeners (substituted benzenesulfonamide herbicides; Hershey and Stoner, Plant Mol. Biol. 17: 679-690, 1991), heat-shock promoters (Ou-Lee et al., Proc. Natl. Acad. Sci. U.S.A. 83: 6815, 1986; Ainley et al., Plant Mol. Biol. 14: 949, 1990), a nitrate-inducible promoter derived from the spinach nitrite

reductase structural nucleic acid sequence (Back et al., Plant Mol. Biol. 17: 9, 1991), hormone-inducible promoters (Yamaguchi-Shinozaki et al., Plant Mol. Biol. 15: 905, 1990), and light-inducible promoters associated with the small subunit of RuBP carboxylase and LHCP families (Kuhlemeier et al., Plant Cell 1: 471, 1989; Feinbaum et al., Mol. Gen. Genet. 226: 449-456, 1991; Weisshaar, et al., EMBO J. 10: 1777-1786, 1991; Lam and Chua, J. Biol. Chem. 266: 17131-17135, 1990; Castresana et al., EMBO J. 7: 1929-1936, 1988; Schulze-Lefert, et al., EMBO J. 8: 651, 1989).

Examples of useful tissue-specific, developmentally-regulated promoters include the  $\beta$ -conglycinin 7S $\alpha$  promoter (Doyle et al., J. Biol. Chem. 261: 9228-9238, 1986; Slighton and Beachy, Planta 172: 356, 1987), and seed-specific promoters (Knutzon, et al., Proc. Natl. Acad. Sci. U.S.A. 89: 2624-2628, 1992; Bustos, et al., EMBO J 10: 1469-1479, 1991; Lam and Chua, Science 248: 471, 1991). Plant functional promoters useful for preferential expression in seed plastid include those from plant storage proteins and from proteins involved in fatty acid biosynthesis in oilseeds. Examples of such promoters include the 5' regulatory regions from such structural nucleic acid sequences as napin (Kridl et al., Seed Sci. Res. 1: 209, 1991), phaseolin, zein, soybean trypsin inhibitor, ACP, stearyl-ACP desaturase, and oleosin. Seed-specific regulation is discussed in EP 0 255 378.

Another exemplary tissue-specific promoter is the lectin promoter, which is specific for seed tissue. The Lectin protein in soybean seeds is encoded by a single structural nucleic acid sequence (Lel) that is only expressed during seed maturation and accounts for about 2 to about 5% of total seed mRNA. The lectin structural nucleic acid sequence and seed-specific promoter have been fully characterized and used to direct seed specific expression in transgenic tobacco plants (Vodkin, et al., Cell, 34: 1023, 1983; Lindstrom, et al., Developmental Genetics, 11: 160, 1990).

Particularly preferred additional promoters in the recombinant vector include the nopaline synthase (nos), mannopine synthase (mas), and octopine synthase (ocs) promoters, which are carried on tumor-inducing plasmids of *Agrobacterium tumefaciens*; the cauliflower mosaic virus (CaMV) 19S and 35S promoters; the enhanced CaMV 35S promoter, the Figwort Mosaic Virus (FMV) 35S promoter; the light-inducible promoter from the small subunit of ribulose-1,5-bisphosphate carboxylase (ssRUBISCO); the EIF-4A promoter from tobacco (Mandel, et al., Plant Mol. Biol. 29: 995-1004, 1995); corn sucrose synthetase 1 (Yang, et al., Proc. Natl. Acad. Sci. USA, 87: 4144-48, 1990); corn alcohol dehydrogenase 1 (Vogel, et al., J. Cell Biochem., (Suppl) 13D: 312, 1989); corn light harvesting complex (Simpson, Science, 233: 34, 1986); corn heat shock protein (Odell, et al., Nature, 313: 810, 1985); the chitinase promoter from *Arabidopsis* (Samac, et al., Plant Cell, 3:1063-1072, 1991); the LTP (Lipid Transfer Protein) promoters from broccoli (Pyee, et al., Plant J., 7: 49-59, 1995); petunia chalcone isomerase (Van Tunen, et al., EMBO J. 7: 1257, 1988); bean glycine rich protein 1 (Keller, et al., EMBO L., 8: 1309-1314, 1989); Potato patatin (Wenzler, et al., Plant Mol. Biol., 12: 41-50, 1989); the ubiquitin promoter from maize (Christensen et al., Plant Mol. Biol., 18: 675,689, 1992); and the actin promoter from rice (McEloy, et al., Plant Cell, 2:163-171, 1990).

The additional promoter is preferably seed selective, tissue selective, constitutive, or inducible. The promoter is most preferably the nopaline synthase (NOS), octopine synthase (OCS), mannopine synthase (MAS), cauliflower mosaic virus 19S and 35S (CaMV19S, CaMV35S), enhanced CaMV (eCaMV), ribulose 1,5-bisphosphate carboxylase (ss-RUBISCO), figwort mosaic virus (FMV), CaMV derived

AS4, tobacco RB7, wheat PDX1, tobacco EIF-4, lectin protein (Lel), or rice RC2 promoter.

#### Structural Nucleic Acid Sequences in the Recombinant Nucleic Acid Vector

The promoter in the recombinant vector is preferably operably linked to a structural nucleic acid sequence. Exemplary structural nucleic acid sequences, and modified forms thereof, are described in detail above. The promoter of the present invention may be operably linked to a structural nucleic acid sequence that is heterologous with respect to the promoter. In one aspect, the structural nucleic acid sequence may generally be any nucleic acid sequence for which an increased level of transcription is desired. The structural nucleic acid sequence preferably encodes a polypeptide that is suitable for incorporation into the diet of a human or an animal. Suitable structural nucleic acid sequences include those encoding a yield protein, a stress resistance protein, a developmental control protein, a tissue differentiation protein, a meristem protein, an environmentally responsive protein, a senescence protein, a hormone responsive protein, an abscission protein, a source protein, a sink protein, a flower control protein, a seed protein, an herbicide resistance protein, a disease resistance protein, a fatty acid biosynthetic enzyme, a tocopherol biosynthetic enzyme, an amino acid biosynthetic enzyme, and an insecticidal protein.

Alternatively, the promoter and structural nucleic acid sequence may be designed to down-regulate a specific nucleic acid sequence. This is typically accomplished by linking the promoter to a structural nucleic acid sequence that is oriented in the antisense direction. One of ordinary skill in the art is familiar with such antisense technology. Using such an approach, a cellular nucleic acid sequence is effectively down regulated as the subsequent steps of translation are disrupted. Nucleic acid sequences may be negatively regulated in this manner.

#### Recombinant Vectors Having Additional Structural Nucleic Acid Sequences

The recombinant vector may also contain one or more additional structural nucleic acid sequences. These additional structural nucleic acid sequences may generally be any sequences suitable for use in a recombinant vector. Such structural nucleic acid sequences include any of the structural nucleic acid sequences, and modified forms thereof, described above. The additional structural nucleic acid sequences may also be operably linked to any of the above-described promoters. The one or more structural nucleic acid sequences may each be operably linked to separate promoters. Alternatively, the structural nucleic acid sequences may be operably linked to a single promoter (i.e. a single operon).

The additional structural nucleic acid sequences preferably encode a yield protein, a stress resistance protein, a developmental control protein, a tissue differentiation protein, a meristem protein, an environmentally responsive protein, a senescence protein, a hormone responsive protein, an abscission protein, a source protein, a sink protein, a flower control protein, a seed protein, an herbicide resistance protein, a disease resistance protein, a fatty acid biosynthetic enzyme, a tocopherol biosynthetic enzyme, an amino acid biosynthetic enzyme, and an insecticidal protein.

Alternatively, the second structural nucleic acid sequence may be designed to down-regulate a specific nucleic acid sequence. This is typically accomplished by operably linking the second structural amino acid, in an antisense orientation, with a promoter. One of ordinary skill in the art is familiar with such antisense technology. The process is also briefly described above. Any nucleic acid sequence may be negatively regulated in this manner.

#### Selectable Markers

The recombinant vector may further comprise a selectable marker. The nucleic acid sequence serving as the selectable marker functions to produce a phenotype in cells which facilitates their identification relative to cells not containing the marker.

Examples of selectable markers include, but are not limited to, a neo gene (Potrykus, et al., *Ann. Rev. Plant Physiol. Plant Mol. Biol.*, 42: 205, 1991), which codes for kanamycin resistance and can be selected for using kanamycin, G418, etc.; a bar gene which codes for bialaphos resistance; a mutant EPSP synthase gene (Hinchee, et al., *Bio/Technology* 6:915-922, 1988) which encodes glyphosate resistance; a nitrilase gene which confers resistance to bromoxynil (Stalker, et al., *J. Biol. Chem.* 263:6310-6314, 1988); a mutant acetolactate synthase gene (ALS) which confers imidazolinone or sulphonylurea resistance (European Patent Application No. 0154204); green fluorescent protein (GFP); and a methotrexate resistant DHFR gene. (Thillet, et al., *J. Biol. Chem.* 263: 12500-12508, 1988).

Other exemplary selectable markers include: a  $\beta$ -glucuronidase or uidA gene, (GUS), which encodes an enzyme for which various chromogenic substrates are known (Jefferson, *Plant Mol. Biol. Rep.* 5:387-405, 1987; Jefferson, et al., *EMBO J.* 6:3901-3907, 1987); an R-locus gene, which encodes a product that regulates the production of anthocyanin pigments (red color) in plant tissues (Dellaporta et al., *Stadler Symposium* 11:263-282, 1988); a  $\beta$ -lactamase gene (Sutcliffe et al., *Proc. Natl. Acad. Sci. (U.S.A.)* 75:3737-3741, 1978), which encodes an enzyme for which various chromogenic substrates are known (e.g., PADAC, a chromogenic cephalosporin); a luciferase gene (Ow, et al., *Science* 234:856-859, 1986); a xylE gene (Zukowsky, et al., *Proc. Natl. Acad. Sci. (U.S.A.)* 80:1101-1105, 1983) which encodes a catechol dioxygenase that can convert chromogenic catechols; an  $\alpha$ -amylase gene (Ikata et al., *Bio/Technol.* 8:241-242, 1990); a tyrosinase gene (Katz et al., *J. Gen. Microbiol.* 129:2703-2714, 1983), which encodes an enzyme capable of oxidizing tyrosine to DOPA and dopaquinone (which in turn condenses to melanin); and an  $\alpha$ -galactosidase, which will turn a chromogenic  $\alpha$ -galactose substrate.

Included within the term "selectable markers" are also genes which encode a secretable marker whose secretion can be detected as a means of identifying or selecting for transformed cells. Examples include markers that encode a secretable antigen that can be identified by antibody interaction, or even secretable enzymes which can be detected catalytically. Selectable secreted marker proteins fall into a number of classes, including small, diffusible proteins which are detectable, (e.g., by ELISA), small active enzymes which are detectable in extracellular solution (e.g.,  $\alpha$ -amylase,  $\beta$ -lactamase, phosphinothricin transferase), or proteins which are inserted or trapped in the cell wall (such as proteins which include a leader sequence such as that found in the expression unit of extension or tobacco PR-S). Other possible selectable marker genes will be apparent to those of skill in the art.

The selectable marker is preferably GUS, green fluorescent protein (GFP), neomycin phosphotransferase II (nptII), luciferase (LUX), an antibiotic resistance coding sequence, or an herbicide (e.g., glyphosate) resistance coding sequence. The selectable marker is most preferably a kanamycin, hygromycin, or herbicide resistance marker.

#### Other Elements in the Recombinant Vector

Various cis-acting untranslated 5' and 3' regulatory sequences may be included in the recombinant nucleic acid vector. Any such regulatory sequences may be provided in a

recombinant vector with other regulatory sequences. Such combinations can be designed or modified to produce desirable regulatory features.

A 3' non-translated region typically provides a transcriptional termination signal, and a polyadenylation signal which functions in plants to cause the addition of adenylate nucleotides to the 3' end of the mRNA. These may be obtained from the 3' regions to the nopaline synthase (nos) coding sequence, the soybean 7S $\alpha$  storage protein coding sequence, the albumin coding sequence, and the pea ssRUBISCO E9 coding sequence. Particularly preferred 3' nucleic acid sequences include nos 3', E93', ADR12 3', 7S $\alpha$  3', 1,1S 3', and albumin 3'.

Typically, nucleic acid sequences located a few hundred base pairs downstream of the polyadenylation site serve to terminate transcription. These regions are required for efficient polyadenylation of transcribed mRNA.

Translational enhancers may also be incorporated as part of the recombinant vector. Thus the recombinant vector may preferably contain one or more 5' non-translated leader sequences which serve to enhance expression of the nucleic acid sequence. Such enhancer sequences may be desirable to increase or alter the translational efficiency of the resultant mRNA. Preferred 5' nucleic acid sequences include dSSU 5', PetHSP70 5', and GmHSP17.9 5'.

The recombinant vector may further comprise a nucleic acid sequence encoding a transit peptide. This peptide may be useful for directing a protein to the extracellular space, a chloroplast, or to some other compartment inside or outside of the cell (see, e.g., European Patent Application Publication Number 0218571).

The structural nucleic acid sequence in the recombinant vector may comprise introns. The introns may be heterologous with respect to the structural nucleic acid sequence. Preferred introns include the rice actin intron and the corn HSP70 intron.

#### Fusion Proteins

Any of the above described structural nucleic acid sequences, and modified forms thereof, may be linked with additional nucleic acid sequences to encode fusion proteins. The additional nucleic acid sequence preferably encodes at least 1 amino acid, peptide, or protein. Production of fusion proteins is routine in the art and many possible fusion combinations exist.

For instance, the fusion protein may provide a "tagged" epitope to facilitate detection of the fusion protein, such as GST, GFP, FLAG, or polyHIS. Such fusions preferably encode between 1 and 50 amino acids, more preferably between 5 and 30 additional amino acids, and even more preferably between 5 and 20 amino acids.

Alternatively, the fusion may provide regulatory, enzymatic, cell signaling, or intercellular transport functions. For example, a sequence encoding a chloroplast transit peptide may be added to direct a fusion protein to the chloroplasts within a plant cell. Such fusion partners preferably encode between 1 and 1000 additional amino acids, more preferably between 5 and 500 additional amino acids, and even more preferably between 10 and 250 amino acids.

#### Probes and Primers

Short nucleic acid sequences having the ability to specifically hybridize to complementary nucleic acid sequences may be produced and utilized in the present invention. These short nucleic acid molecules may be used as probes to identify the presence of a complementary nucleic acid sequence in a given sample. Thus, by constructing a nucleic acid probe which is complementary to a small portion of a particular

nucleic acid sequence, the presence of that nucleic acid sequence may be detected and assessed.

Use of these probes may greatly facilitate the identification of transgenic plants which contain the presently disclosed nucleic acid molecules. The probes may also be used to screen cDNA or genomic libraries for additional nucleic acid sequences related or sharing homology to the presently disclosed promoters and structural nucleic acid sequences.

Alternatively, the short nucleic acid sequences may be used as oligonucleotide primers to amplify or mutate a complementary nucleic acid sequence using PCR technology. These primers may also facilitate the amplification of related complementary nucleic acid sequences (e.g. related nucleic acid sequences from other species).

The short nucleic acid sequences may be used as probes and specifically as PCR probes. A PCR probe is a nucleic acid molecule capable of initiating a polymerase activity while in a double-stranded structure with another nucleic acid. Various methods for determining the structure of PCR probes and PCR techniques exist in the art. Computer generated searches using programs such as Primer3 ([www-genome.wi.mit.edu/cgi-bin/primer/primer2.cgi](http://www-genome.wi.mit.edu/cgi-bin/primer/primer2.cgi)), STSPipeline ([www-genome.wi.mit.edu/cgi-bin/www.STS\\_Pipeline](http://www-genome.wi.mit.edu/cgi-bin/www.STS_Pipeline)), or GeneUp (Pesole, et al., *BioTechniques* 25:112-123, 1998), for example, can be used to identify potential PCR primers.

The primer or probe is generally complementary to a portion of a nucleic acid sequence that is to be identified, amplified, or mutated. The primer or probe should be of sufficient length to form a stable and sequence-specific duplex molecule with its complement. The primer or probe preferably is about 10 to about 200 nucleotides long, more preferably is about 10 to about 100 nucleotides long, even more preferably is about 10 to about 50 nucleotides long, and most preferably is about 14 to about 30 nucleotides long.

The primer or probe may be prepared by direct chemical synthesis, by PCR (See, for example, U.S. Pat. No. 4,683,195, and 4,683,202), or by excising the nucleic acid specific fragment from a larger nucleic acid molecule.

#### Sequence Analysis

In the present invention, sequence similarity or identity is preferably determined using the "Best Fit" or "Gap" programs of the Sequence Analysis Software Package™ (Version 10; Genetics Computer Group, Inc., Center, Madison, Wis.). "Gap" utilizes the algorithm of Needleman and Wunsch (Needleman and Wunsch, *Journal of Molecular Biology* 48:443453, 1970) to find the alignment of two sequences that maximizes the number of matches and minimizes the number of gaps. "BestFit" performs an optimal alignment of the best segment of similarity between two sequences. Optimal alignments are found by inserting gaps to maximize the number of matches using the local homology algorithm of Smith and Waterman (Smith, et al. In: *Genetic Engineering: Principles and Methods*, Setlow et al., Eds., Plenum Press, N.Y., 1-32, 1981; Smith and Waterman *Advances in Applied Mathematics*, 2:482-489, 1981).

The Sequence Analysis Software Package described above contains a number of other useful sequence analysis tools for identifying homologues of the presently disclosed nucleotide and amino acid sequences. For example, the "BLAST" program (Altschul, et al., *Journal of Molecular Biology* 215: 403-410, 1990) searches for sequences similar to a query sequence (either peptide or nucleic acid) in a specified database (e.g., sequence databases maintained at the National Center for Biotechnology Information (NCBI) in Bethesda, Md., USA); "FastA" (Lipman and Pearson, *Science*, 227: 1435-1441, 1985; see also Pearson and Lipman, *Proceedings of the National Academy of Sciences USA* 85, 2444-2448,

1988; Pearson, *Methods in Enzymology*, (R. Doolittle, ed.), 183, 63-98, Academic Press, San Diego, Calif., USA, 1990) performs a Pearson and Lipman search for similarity between a query sequence and a group of sequences of the same type (nucleic acid or protein); "TfastA" performs a Pearson and Lipman search for similarity between a protein query sequence and any group of nucleotide sequences (it translates the nucleotide sequences in all six reading frames before performing the comparison); "FastX" performs a Pearson and Lipman search for similarity between a nucleotide query sequence and a group of protein sequences, taking frameshifts into account. "TfastX" performs a Pearson and Lipman search for similarity between a protein query sequence and any group of nucleotide sequences, taking frameshifts into account (it translates both strands of the nucleic acid sequence before performing the comparison).

#### Transgenic Plants and Plant Cells

The invention also includes and provides transformed plant cells which comprise a nucleic acid molecule of the present invention.

Preferred nucleic acid sequences of the present invention include, without limitation, recombinant vectors, structural nucleic acid sequences, promoters, and other regulatory elements, are described above. A promoter preferably comprises a nucleic acid sequence that hybridizes under stringent conditions with a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, any complement thereof; and any fragment thereof; or exhibits 85% or greater identity, and more preferably at least 86 or greater, 87 or greater, 88 or greater, 89 or greater, 90 or greater, 91 or greater, 92 or greater, 93 or greater, 94 or greater, 95 or greater, 96 or greater, 97 or greater, 98 or greater, or 99% or greater identity to a nucleic acid sequence selected from the group consisting of SEQ ID NO:1 through SEQ ID NO:57,467, any complement thereof; and any fragment thereof. A promoter most preferably comprises SEQ ID NO:1 through SEQ ID NO:57,467.

Methods for preparing such recombinant vectors are well known in the art. For example, methods for making recombinant vectors particularly suited to plant transformation are described in U.S. Pat. Nos. 4,971,908, 4,940,835, 4,769,061 and 4,757,011. These vectors have also been reviewed (Rodriguez et al., *Vectors: A Survey of Molecular Cloning Vectors and Their Uses*, Butterworths, Boston, 1988; Glick, et al., *Methods in Plant Molecular Biology and Biotechnology*, CRC Press, Boca Raton, Fla., 1993) and are described above.

Typical vectors useful for expression of nucleic acids in plant cells are well known in the art and include vectors derived from the tumor-inducing (Ti) plasmid of *Agrobacterium tumefaciens* (Rogers, et al., *Meth In Enzymol*, 153: 253-277, 1987). Other recombinant vectors useful for plant transformation, have also been described (Fromm et al., *Proc. Natl. Acad. Sci. USA*, 82(17): 5824-5828, 1985). Elements of such recombinant vectors are discussed above.

The transformed plant or cell may generally be any plant or cell that is compatible with the present invention. The plant or cell preferably is alfalfa, apple, banana, barley, bean, broccoli, cabbage, carrot, castorbean, celery, citrus, clover, coconut, coffee, corn, cotton, cucumber, garlic, grape, linseed, melon, oat, olive, onion, palm, parsnip, pea, peanut, pepper, potato, radish, rapeseed, rice, rye, sorghum, soybean, spinach, strawberry, sugarbeet, sugarcane, sunflower, tobacco, tomato, or wheat. The transformed plant or cell is more preferably rice, sorghum, maize, barley, wheat, canola, soybean, or maize; even more preferably a rice, sorghum, maize, barley, or wheat; and most preferably is rice. The rice plant or cell

is preferably *Oryza sativa* L. (japonica type), and more preferably *Oryza sativa* L. to (japonica type), cv. Nipponbare.

Typical vectors useful for expression of nucleic acids in cells and higher plants are well known in the art and include vectors derived from the tumor-inducing (Ti) plasmid of *Agrobacterium tumefaciens* (Rogers, et al., *Meth. In Enzymol*, 153: 253-277, 1987). Other recombinant vectors useful for plant transformation, have also been described (Fromm et al., *Proc. Natl. Acad. Sci. USA*, 82(17): 5824-5828, 1985). Elements of such recombinant vectors are discussed above.

#### Method for Preparing Transformed Cells

The invention is also directed to a method of producing transformed cells which comprise, in a 5' to 3' orientation, a promoter operably linked to a heterologous structural nucleic acid sequence. Other sequences may also be introduced into the cell along with the promoter and structural nucleic acid sequence. These other sequences may include 3' transcriptional terminators, 3' polyadenylation signals, other untranslated sequences, transit or targeting sequences, selectable markers, enhancers, and operators.

Preferred recombinant vectors, structural nucleic acid sequences, promoters, and other regulatory elements are described above. The promoter preferably has a nucleic acid sequence that hybridizes under stringent conditions with SEQ ID NO:1 through SEQ ID NO:57,467, or any complement thereof; or exhibits 85% or greater identity, and more preferably at least 86 or greater, 87 or greater, 88 or greater, 89 or greater, 90 or greater, 91 or greater, 92 or greater, 93 or greater, 94 or greater, 95 or greater, 96 or greater, 97 or greater, 98 or greater, or 99% or greater identity to SEQ ID NO:1 through SEQ ID NO:57,467.

The method generally comprises the steps of selecting a suitable host cell, transforming the host cell with a recombinant vector, and obtaining the transformed host cell.

There are many methods for introducing nucleic acids into plant cells. Suitable methods include bacterial infection (e.g. *Agrobacterium*), binary bacterial artificial chromosome vectors, direct delivery of DNA (e.g. via PEG-mediated transformation, desiccation/inhibition-mediated DNA uptake, electroporation, agitation with silicon carbide fibers, and acceleration of DNA coated particles, etc. (reviewed in Potrykus, et al., *Ann. Rev. Plant Physiol. Plant Mol. Biol.*, 42: 205, 1991).

Technology for introduction of DNA into cells is well known to those of skill in the art. These methods can generally be classified into four categories: (1) chemical methods (Graham and Van der Eb, *Virology*, 54(2): 536-539, 1973; Zatloukal, et al., *Ann. N.Y. Acad. Sci.*, 660: 136-153, 1992); (2) physical methods such as microinjection (Capocchi, *Cell*, 22(2): 479-488, 1980), electroporation (Wong and Neumann, *Biochim. Biophys. Res. Commun.*, 107(2): 584-587, 1982; Fromm et al., *Proc. Natl. Acad. Sci. USA*, 82(17): 5824-5828, 1985; U.S. Pat. No. 5,384,253) and particle acceleration (Johnston and Tang, *Methods Cell Biol.*, 43(A): 353-365, 1994; Fynan et al., *Proc. Natl. Acad. Sci. USA*, 90(24): 11478-11482, 1993); (3) viral vectors (Clapp, *Clin. Perinatol.*, 20(1): 155-168, 1993; Lu, et al., *J. Exp. Med.*, 178(6): 2089-2096, 1993; Eglitis and Anderson, *Biotechniques*, 6(7): 608-614, 1988); and (4) receptor-mediated mechanisms (Curiel et al., *Hum. Gen. Ther.*, 3(2):147-154, 1992; Wagner, et al., *Proc. Natl. Acad. Sci. USA*, 89(13): 6099-6103, 1992). Alternatively, nucleic acids can be directly introduced into pollen by directly injecting a plant's reproductive organs (Zhou, et al., *Methods in Enzymology*, 101: 433, 1983; Hess, *Intern Rev. Cytol.*, 107: 367, 1987; Luo, et al., *Plant Mol Biol. Reporter* 6: 165, 1988; Pena, et al., *Nature*, 325: 274, 1987).

The nucleic acids may also be injected into immature embryos (Neuhaus, et al., *Theor. Appl. Genet.*, 75: 30, 1987).

The recombinant vector used to transform the host cell typically comprises, in a 5' to 3' orientation: a promoter to direct the transcription of a structural nucleic acid sequence, a structural nucleic acid sequence, a 3' transcriptional terminator, and a 3' polyadenylation signal. The recombinant vector may further comprise untranslated nucleic acid sequences, transit and targeting nucleic acid sequences, selectable markers, enhancers, or operators.

Suitable recombinant vectors, structural nucleic acid sequences, promoters, and other regulatory elements include, without limitation, those described above.

The regeneration, development, and cultivation of plants from transformed plant protoplast or explants is well taught in the art (Weissbach and Weissbach, *Methods for Plant Molecular Biology*, (Eds.), Academic Press, Inc., San Diego, Calif., 1988; Horsch et al., *Science*, 227: 1229-1231, 1985). In this method, transformants are generally cultured in the presence of a selective media which selects for the successfully transformed cells and induces the regeneration of plant shoots (Fraleley et al., *Proc. Natl. Acad. Sci. U.S.A.*, 80: 4803, 1983). These shoots are typically obtained within two to four months.

The shoots are then transferred to an appropriate root-inducing medium containing the selective agent and an antibiotic to prevent bacterial growth. Many of the shoots will develop roots. These are then transplanted to soil or other media to allow the continued development of roots. The method, as outlined, will generally vary depending on the particular plant strain employed.

The regenerated transgenic plants are self-pollinated to provide homozygous transgenic plants. Alternatively, pollen obtained from the regenerated transgenic plants may be crossed with non-transgenic plants, preferably inbred lines of agronomically important species. Conversely, pollen as from non-transgenic plants may be used to pollinate the regenerated transgenic plants.

The transgenic plant may pass along the transformed nucleic acid sequence to its progeny. The transgenic plant is preferably homozygous for the transformed nucleic acid sequence and transmits that sequence to all of its offspring upon as a result of sexual reproduction. Progeny may be grown from seeds produced by the transgenic plant. These additional plants may then be self-pollinated to generate a true breeding line of plants.

The progeny from these plants are evaluated, among other things, for gene expression. The gene expression may be detected by several common methods such as western blotting, northern blotting, immunoprecipitation, and ELISA.

Methods for transforming dicots, primarily by use of *Agrobacterium tumefaciens* and obtaining transgenic plants have been published for cotton (U.S. Pat. Nos. 5,004,863; 5,159,135; 5,518,908); soybean (U.S. Pat. No. 5,569,834; 5,416,011; McCabe, et al., *Biotechnology*, 6: 923, 1988; Christou et al., *Plant Physiol.* 87:671-674 (1988)); Brassica (U.S. Pat. No. 5,463,174); peanut (Cheng et al., *Plant Cell Rep.* 15:653-657 (1996), McKenty et al., *Plant Cell Rep.* 14:699-703 (1995)); papaya; and pea (Grant et al., *Plant Cell Rep.* 15:254-258 (1995)).

Transformation of monocotyledons using electroporation, particle bombardment and *Agrobacterium* have also been reported. Transformation and plant regeneration have been achieved in asparagus (Bytebier et al., *Proc. Natl. Acad. Sci. (USA)* 84:5354 (1987)); barley (Wan and Lemaux, *Plant Physiol* 104:37 (1994)); maize (Rhodes et al., *Science* 240: 204 (1988); Gordon-Kamm et al., *Plant Cell* 2:603-618

(1990); Fromm et al., *Bio/Technology* 8:833 (1990); Koziel et al., *Bio/Technology* 11: 194 (1993); Armstrong et al., *Crop Science* 35:550-557 (1995)); oat (Somers et al., *Bio/Technology* 10:1589 (1992)); orchard grass (Horn et al., *Plant Cell Rep.* 7:469 (1988)); rice (Toriyama et al., *Theor Appl. Genet.* 205:34 (1986); Part et al., *Plant Mol.* 32:1135-1148 (1996); Abedinia et al., *Aust. J. Plant Physiol.* 24:133-141 (1997); Zhang and Wu, *Theor. Appl. Genet.* 76:835 (1988); Zhang et al., *Plant Cell Rep.* 7:379 (1988); Battraw and Hall, *Plant Sci.* 86:191-202 (1992); Christou et al., *Bio/Technology* 9:957 (1991)); rye (De la Pena et al., *Nature* 325:274 (1987)); sugarcane (Bower and Birch, *Plant J.* 2:409 (1992)); tall fescue (Wang et al., *Bio/Technology* 10:691 (1992)) and wheat (Vasil et al., *Bio/Technology* 10:667 (1992); U.S. Pat. No. 5,631,152).

#### Other Transformed Organisms

Any of the above described promoters and structural nucleic acid sequences may be introduced into any cell or organism such as a mammalian cell, mammal, fish cell, fish, bird cell, bird, algae cell, algae, fungal cell, fungi, or bacterial cell. Preferred hosts and transformants include: fungal cells such as *Aspergillus*, yeasts, mammals (particularly bovine and porcine), insects, bacteria and algae.

The transformed cell or organism is preferably prokaryotic, more preferably a bacterial cell, even more preferably a *Agrobacterium*, *Bacillus*, *Escherichia*, *Pseudomonas* cell, and most preferably is an *Escherichia coli* cell. Alternatively, the transformed organism is preferably a yeast or fungal cell. The yeast cell is preferably a *Saccharomyces cerevisiae*, *Schizosaccharomyces pombe*, or *Pichia pastoris*.

Methods to transform such cells or organisms are known in the art (EP 0238023; Yelton et al., *Proc. Natl. Acad. Sci. (U.S.A.)*, 81:1470-1474 (1984); Malardier et al., *Gene*, 78:147-156 (1989); Becker and Guarente, In: Abelson and Simon (eds.), *Guide to Yeast Genetics and Molecular Biology*, *Methods Enzymol.*, Vol. 194, pp. 182-187, Academic Press, Inc., New York; Ito et al., *J. Bacteriology*, 153:163 (1983); Hinnen et al., *Proc. Natl. Acad. Sci. (U.S.A.)*, 75:1920 (1978); Bennett and LaSure (eds.), *More Gene Manipulations in Fungi*, Academic Press, Calif. (1991)). Methods to produce proteins of the present invention from such organisms are also known (Kudla et al., *EMBO*, 9:1355-1364 (1990); Jarai and Buxton, *Current Genetics*, 26:2238-2244 (1994); Verdier, *Yeast*, 6:271-297 (1990); MacKenzie et al., *Journal of Gen. Microbiol.*, 139:2295-2307 (1993); Hartl et al., *TIBS*, 19:20-25 (1994); Bergeron et al., *TIBS*, 19:124-128 (1994); Demolder et al., *J. Biotechnology*, 32:179-189 (1994); Craig, *Science*, 260:1902-1903 (1993); Gething and Sambrook, *Nature*, 355:33-45 (1992); Puig and Gilbert, *J. Biol. Chem.*, 269:7764-7771 (1994); Wang and Tsou, *FASEB Journal*, 7:1515-1517 (1993); Robinson et al., *Bio/Technology*, 1:381-384 (1994); Enderlin and Ogrzydzki, *Yeast*, 10:67-79 (1994); Fuller et al., *Proc. Natl. Acad. Sci. (U.S.A.)*, 86:1434-1438 (1989); Julius et al., *Cell*, 37:1075-1089 (1984); Julius et al., *Cell*, 32:839-852 (1983)).

#### Exemplary Uses

The presently disclosed promoter sequences may be used as genetic markers and employed in genetic mapping studies using linkage analysis. A genetic linkage map shows the relative locations of specific DNA markers along a chromosome. Maps are used for the identification of genes associated with genetic diseases or phenotypic traits, comparative genomics, and as a guide for physical mapping. Through genetic mapping, a fine scale linkage map can be developed using DNA markers, and, then, a genomic DNA library of large-sized fragments can be screened with molecular markers linked to the desired trait. In a preferred embodiment of

the present invention, a genomic library is screened with the promoter sequences of the present invention.

Mapping marker locations is based on the observation that two markers located near each other on the same chromosome will tend to be passed together from parent to offspring. During gamete production, DNA strands occasionally break and rejoin in different places on the same chromosome or on the homologous chromosome. The closer the markers are to each other, the more tightly linked and the less likely a recombination event will fall between and separate them. Recombination frequency thus provides an estimate of the distance between two markers.

In segregating populations, target genes have been reported to have been placed within an interval of 5-10 cM with a high degree of certainty (Tanksley et al., *Trends in Genetics* 11(2): 63-68 (1995)). The markers defining this interval are used to screen a larger segregating population to identify individuals derived from one or more gametes containing a crossover in the given interval. Such individuals are useful in orienting other markers closer to the target gene. Once identified, these individuals can be analyzed in relation to all molecular markers within the region to identify those closest to the target.

Markers of the present invention can be employed to construct linkage maps and to locate genes with qualitative and quantitative effects. The genetic linkage of additional marker molecules can be established by a genetic mapping model such as, without limitation, the flanking marker model reported by Lander and Botstein, *Genetics*, 121:185-199 (1989), and the interval mapping, based on maximum likelihood methods described by Lander and Botstein, *Genetics*, 121:185-199 (1989), and implemented in the software package MAPMAKER/QTL (Lincoln and Lander, *Mapping Genes Controlling Quantitative Traits Using MAPMAKER/QTL*, Whitehead Institute for Biomedical Research, Massachusetts, (1990)). Additional software includes Qgene, Version 2.23 (1996), Department of Plant Breeding and Biometry, 266 Emerson Hall, Cornell University, Ithaca, N.Y.). Use of the Qgene software is a particularly preferred approach.

A maximum likelihood estimate (MLE) for the presence of a marker is calculated, together with an MLE assuming no QTL effect, to avoid false positives. A log<sub>10</sub> of an odds ratio (LOD) is then calculated as:  $LOD = \log_{10} (MLE \text{ for the presence of a QTL} / MLE \text{ given no linked QTL})$ .

The LOD score essentially indicates how much more likely the data are to have arisen assuming the presence of a QTL than in its absence. The LOD threshold value for avoiding a false positive with a given confidence, say 95%, depends on the number of markers and the length of the genome. Graphs indicating LOD thresholds are set forth in Lander and Botstein, *Genetics*, 121:185-199 (1989), and further described by Arús and Moreno-González, *Plant Breeding*, Hayward, Rosemark, Romagosa (eds.) Chapman & Hall, London, pp. 314-331 (1993).

Additional models can be used. Many modifications and alternative approaches to interval mapping have been reported, including the use of non-parametric methods (Kruglyak and Lander, *Genetics*, 139:1421-1428 (1995)). Multiple regression methods or models can be also be used, in which the trait is regressed on a large number of markers (Jansen, *Biometrics in Plant Breed*, van Oijen, Jansen (eds.) *Proceedings of the Ninth Meeting of the Eucarpia Section Biometrics in Plant Breeding*, The Netherlands, pp. 116-124 (1994); Weber and Wricke, *Advances in Plant Breeding*, Blackwell, Berlin, 16 (1994). Procedures combining interval mapping with regression analysis, whereby the phenotype is regressed onto a single putative QTL at a given marker inter-

val, and at the same time onto a number of markers that serve as 'cofactors,' have been reported by Jansen and Stam, *Genetics*, 136:1447-1455 (1994) and Zeng, *Genetics*, 136:1457-1468 (1994). Generally, the use of cofactors reduces the bias and sampling error of the estimated QTL positions (Utz and Melchinger, *Biometrics in Plant Breeding*, van Oijen, Jansen (eds.) *Proceedings of the Ninth Meeting of the Eucarpia Section Biometrics in Plant Breeding*, The Netherlands, pp.195-204 (1994), thereby improving the precision and efficiency of QTL mapping (Zeng, *Genetics*, 136:1457-1468 (1994)). These models can be extended to multi-environment experiments to analysis genotype-environment interactions (Jansen et al., *Theo. Appl. Genet.* 91:33-37 (1995)).

Selection of an appropriate mapping population is important to map construction. The choice of appropriate mapping population depends on the type of marker systems employed (Tanksley et al., J. P. Gustafson and R. Appels (eds.), Plenum Press, New York, pp. 157-173 (1988)). Consideration must be given to the source of parents (adapted vs. exotic) used in the mapping population. Chromosome pairing and recombination rates can be severely disturbed (suppressed) in wide crosses (adapted x exotic) and generally yield greatly reduced linkage distances. Wide crosses will usually provide segregating populations with a relatively large array of polymorphisms when compared to progeny in a narrow cross (adapted x adapted).

An F<sub>2</sub> population is the first generation of selfing after the hybrid seed is produced. Usually a single F<sub>1</sub> plant is selfed to generate a population segregating for all the genes in Mendelian (1:2:1) fashion. Maximum genetic information is obtained from a completely classified F<sub>2</sub> population using a codominant marker system (Mather, *Measurement of Linkage in Heredity*: Methuen and Co., (1938)). In the case of dominant markers, progeny tests (e.g., F<sub>3</sub>, BCF<sub>2</sub>) are required to identify the heterozygotes, thus making it equivalent to a completely classified F<sub>2</sub> population. However, this procedure is often prohibitive because of the cost and time involved in progeny testing. Progeny testing of F<sub>2</sub> individuals is often used in map construction where phenotypes do not consistently reflect genotype (e.g., disease resistance) or where trait expression is controlled by a QTL. Segregation data from progeny test populations (e.g., F<sub>3</sub> or BCF<sub>2</sub>) can be used in map construction. Marker-assisted selection can then be applied to cross progeny based on marker-trait map associations (F<sub>2</sub>, F<sub>3</sub>), where linkage groups have not been completely disassociated by recombination events (i.e., maximum disequilibrium).

Recombinant inbred lines (RIL) (genetically related lines; usually >F<sub>5</sub>, developed from continuously selfing F<sub>2</sub> lines towards homozygosity) can be used as a mapping population. Information obtained from dominant markers can be maximized by using RIL because all loci are homozygous or nearly so. Under conditions of tight linkage (i.e., about <10% recombination), dominant and co-dominant markers evaluated in RIL populations provide more information per individual than either marker type in backcross populations (Reiter, *Proc. Natl. Acad. Sci. USA* 89:1477-1481 (1992)). However, as the distance between markers becomes larger (i.e., loci become more independent), the information in RIL populations decreases dramatically when compared to codominant markers.

Backcross populations (e.g., generated from a cross between a successful variety (recurrent parent) and another variety (donor parent) carrying a trait not present in the former) can be utilized as a mapping population. A series of backcrosses to the recurrent parent can be made to recover most of its desirable traits. Thus a population is created con-

sisting of individuals nearly like the recurrent parent but each individual carries varying amounts or mosaic of genomic regions from the donor parent. Backcross populations can be useful for mapping dominant markers if all loci in the recurrent parent are homozygous and the donor and recurrent parent have contrasting polymorphic marker alleles (Reiter et al., Proc. Natl. Acad. Sci. USA 89:1477-1481 (1992)). Information obtained from backcross populations using either codominant or dominant markers is less than that obtained from F2 populations because one, rather than two, recombinant gametes are sampled per plant. Backcross populations, however, are more informative (at low marker saturation) when compared to RILs as the distance between linked loci increases in RIL populations (i.e., about 0.15% recombination). Increased recombination can be beneficial for resolution of tight linkages, but may be undesirable in the construction of maps with low marker saturation.

Near-isogenic lines (NIL)(created by many backcrosses to produce an array of individuals that are nearly identical in genetic composition except for the trait or genomic region under interrogation) can be used as a mapping population. In mapping with NILs, only a portion of the polymorphic loci are expected to map to a selected region.

Bulk segregant analysis (BSA) is a method developed for the rapid identification of linkage between markers and traits of interest (Michelmore et al., Proc. Natl. Acad. Sci. USA 88:9828-9832 (1991)). In BSA, two bulked DNA samples are drawn from a segregating population originating from a single cross. These bulks contain individuals that are identical for a particular trait (resistant or susceptible to particular disease) or genomic region but arbitrary at unlinked regions (i.e., heterozygous). Regions unlinked to the target region will not differ between the bulked samples of many individuals in BSA.

It is understood that one or more of the nucleic acid molecules of the present invention may in one embodiment be used as markers in genetic mapping. In a preferred embodiment, nucleic acid molecules of the present invention may in one embodiment be used as markers with rice.

Nucleic acid molecules of the present invention can be used in comparative mapping (physical and genetic) and to isolate molecules from other cereals based on the syntenic relationship between cereals. Comparative mapping within families provides a method to the degree of sequence conservation, gene order, ploidy of species, ancestral relationships and the rates at which individual genomes are evolving. Comparative mapping has been carried out by cross-hybridizing molecular markers across species within a given family.

In a preferred embodiment, the nucleic acid molecules of the present invention can be utilized to isolate corresponding syntenic regions in non-rice plants (Bennetzen and Freeling, Trends in Genet., 9(8):259-261 (1993); Ahn et al., Mol. Gen. Genet., 241(5-6):483-490 (1993); Schwarzacher, Cur. Opin. Genet. & Devel., 4(6): 868-874 (1994); Kurata et al., Bio/Technology, 12:276-278 (1994); Kilian et al., Nucl. Acids Res., 23(14):2729-2733 (1995); Bennett, Symp. Soc. Exp. Biol., 50:45-52 (1996); Hu et al., Genetics, 142(3):1021-1031 (1996); Kilian, Plant Mol. Biol., 35:187-195 (1997); Bennetzen and Freeling, Genome Res., 7(4):301-306 (1997); Foote et al., Genetics, 147(2):801-807 (1997); Gallego et al., Genome, 41(3):328-336 (1998). Gale and Devos Proc. Natl. Acad. Sci. USA 95:1971-1974 (1998); Bennetzen et al., Proc. Natl. Acad. Sci. USA, 95:1975-1978 (1998); Messing and Llaça, Proc. Natl. Acad. Sci. USA 95:2017-2020 (1998); McCouch, Proc. Natl. Acad. Sci. USA, 95:1983-1985 (1998); Goff, Cur. Opin. Plant Biol. 2:85-89 (1999); Bailey et al., Theor. Appl. Genet., 98:281-284 (1999); Zhang et al., Proc.

Natl. Acad. Sci. USA, 91:8675-8679 (1994); Yano and Sasaki, Plant Mol. Biol., 35:145-153 (1997); Leister et al., Proc. Natl. Acad. Sci. USA, 95:370-375 (1998); Lin et al., Phytopathology 86(11):1156-1159 (1996); Havukkala, Cur. Opin. Genet. Dev., 96:711-713 (1996); and Lee, The Society for Experimental Biology, pp.31-38 (1996). Synteny between rice and barley has recently been reported in the genomic region carrying malting quality Quantitative Trait Loci (QTL) (Kleinohfs et al., Genome 41:373-380 (1998)). Likewise, mapping of the liguleless region of sorghum, a region containing a developmental control gene, was facilitated using molecular markers from a syntenic region of the rice genome (Christou et al., Genetics 148:1983-1992 (1998)).

In a particularly preferred embodiment, the nucleic acid molecules of the present invention that define a genomic region in rice plants associated with a desirable phenotype are utilized to obtain corresponding syntenic regions in non-rice plants. A region can be defined either physically or genetically. In an even more preferred embodiment, the nucleic acid molecules of the present invention that define a genomic region in rice plants associated with a desirable phenotype are utilized to obtain corresponding syntenic regions in corn plants. A region can be defined either physically or genetically.

One or more of the nucleic acids molecules may be used to define a physical genomic region. For example, two nucleic acid molecules of the present invention can act to define a physical genomic region that lies between them. Moreover, for example, a physical genomic region may be defined by a distance relative to a nucleic acid molecule. In a preferred embodiment of the present invention, the defined physical genomic region is less than about 1,000 kb, more preferably less than about 500 kb, even more preferably less than about 100 kb or less than about 50 kb.

One or more of the nucleic acids molecules may be used to define a genomic region by its genetic distance from one or more of the nucleic acid molecules of the present invention. In a preferred embodiment of the present invention, the genomic region is defined by its linkage to a nucleic acid molecule of the present invention. In such a preferred embodiment, the genomic region that is defined by one or more nucleic acid molecules of the present invention is located within about 50 centimorgans, more preferably within about 20 centimorgans, even more preferably with about 10, about 5 or about 2 centimorgans of the trait or marker at issue.

In another particularly preferred embodiment, two or more nucleic acid molecules of the present invention derived from rice plants that flank a genomic region of interest in rice plants are used to isolate the syntenic region in another cereal, more preferably maize, sorghum, barley, or wheat. Regions of interest in rice include, without limitation, those regions that are associated with a commercially desirable phenotype in rice. In another particularly preferred embodiment the desirable phenotype in rice is the result of a quantitative trait locus (QTL) present in the region.

One exemplary approach to isolate syntenic genomic regions is as follows. Nucleic acid sequences derived from rice of the present invention can be used to select large insert clones from a total genomic DNA library of a related species such as maize, sorghum, barley, or wheat. Any appropriate method to screen the genomic library with a nucleic acid molecule of the present invention may be used to select the required clones (See, for example, Birren et al., Detecting Genes: A Laboratory Manual, Cold Spring Harbor, New York, N.Y. (1998)). For example, direct hybridization of a nucleic acid molecule of the present invention to mapping filters comprising the genomic DNA of the syntenic species

can be used to select large insert clones from a total genomic DNA library of a related species. The selected clones can then be used to physically map the region in the target species. An advantage of this method for comparative mapping is that no mapping population or linkage map of the target species is needed and the clones may also be used in other closely related species. By comparing the results obtained by genetic mapping in model plants, with those from other species, similarities of genomic structure among plants species can be established. Cross-hybridization of RFLP markers have been reported and conserved gene order has been established in many studies. Such macroscopic synteny is utilized for the estimation of correspondence of loci among these crops. These loci include not only Mendelian genes but also Quantitative Trait Loci (QTL) (Mohan et al., *Molecular Breeding* 3:87-103 (1997)). Other methods to isolate syntenic nucleic acid molecules may be used.

It is understood that markers of the present invention may be used in comparative mapping. In a preferred embodiment the markers of present invention may be used in the comparative mapping of cereals, more preferably maize, barley, sorgham, and wheat.

It is understood that markers of the present invention may be used to isolate promoters and other nucleic acid sequences from other cereals based on the syntenic relationship between such cereals. In a preferred embodiment the cereal is selected from the group of maize, sorgham, barley, and wheat.

The nucleic acid molecules of the present invention can be used to identify polymorphisms. In one embodiment, one or more of the nucleic acid molecules may be employed as a marker nucleic acid molecule to identify such polymorphism (s). Alternatively, such polymorphisms can be detected through the use of a marker nucleic acid molecule or a marker protein that is genetically linked to (i.e., a polynucleotide that co-segregates with) such polymorphism(s). In a preferred embodiment, the plant is selected from the group consisting of cereals, and more preferably rice, maize, barley, sorgham, and wheat.

In an alternative embodiment, such polymorphisms can be detected through the use of a marker nucleic acid molecule that is physically linked to such polymorphism(s). For this purpose, marker nucleic acid sequences located within 1 mb of the polymorphism(s), and more preferably within 100 kb of the polymorphism(s), and most preferably within 10 kb of the polymorphism(s) can be employed.

The genomes of animals and plants naturally undergo spontaneous mutation in the course of their continuing evolution (Gusella, *Ann. Rev. Biochem.* 55:831-854 (1986)). A "polymorphism" is a variation or difference in the sequence of the gene or its flanking regions that arises in some of the members of a species. The variant sequence and the "original" sequence co-exist in the species' population. In some instances, such co-existence is in stable or quasi-stable equilibrium.

A polymorphism is thus said to be "allelic," in that, due to the existence of the polymorphism, some members may have the original sequence (i.e., the original "allele") whereas other members may have the variant sequence (i.e., the variant "allele"). In the simplest case, only one variant sequence may exist, and the polymorphism is thus said to be di-allelic. In other cases, the population may contain multiple alleles, and the polymorphism is termed tri-allelic, etc. A single gene may have multiple different unrelated polymorphisms. For example, it may have a di-allelic polymorphism at one site, and a multi-allelic polymorphism at another site.

The variation that defines the polymorphism may range from a single nucleotide variation to the insertion or deletion

of extended regions within a gene. In some cases, the DNA sequence variations are in regions of the genome that are characterized by short tandem repeats (STRs) that include tandem di- or tri-nucleotide repeated motifs of nucleotides. Polymorphisms characterized by such tandem repeats are referred to as "variable number tandem repeat" ("VNTR") polymorphisms. VNTRs have been used in identity analysis (Weber, U.S. Pat. No. 5,075,217; Armour et al., *FEBS Lett.* 307:113-115 (1992); Jones et al., *Eur. J. Haematol.* 39:144-147 (1987); Horn et al., PCT Application WO91/14003; Jeffreys, European Patent Application 370,719; Jeffreys, U.S. Pat. No. 5,175,082; Jeffreys et al., *Amer. J. Hum. Genet.* 39:11-24 (1986); Jeffreys et al., *Nature* 316:76-79 (1985); Gray et al., *Proc. R. Acad. Soc. Lond.* 243:241-253 (1991); Moore et al., *Genomics* 10:654-660 (1991); Jeffreys et al., *Anim. Genet.* 18:1-15 (1987); Hillel et al., *Anim. Genet.* 20:145-155 (1989); Hillel et al., *Genet.* 124:783-789 (1990)).

The detection of polymorphic sites in a sample of DNA may be facilitated through the use of nucleic acid amplification methods. Such methods specifically increase the concentration of polynucleotides that span the polymorphic site, or include that site and sequences located either distal or proximal to it. Such amplified molecules can be readily detected by gel electrophoresis or other means.

The most preferred method of achieving such amplification employs the polymerase chain reaction ("PCR") (Mullis et al., *Cold Spring Harbor Symp. Quant. Biol.* 51:263-273 (1986); Erlich et al., *European Patent Appln.* 50,424; *European Patent Appln.* 84,796, *European Patent Application* 258,017, *European Patent Appln.* 237,362; Mullis, *European Patent Appln.* 201,184; Mullis, et al., U.S. Pat. No. 4,683,202; Erlich, U.S. Pat. No. 4,582,788; and Saiki et al., U.S. Pat. No. 4,683,194), using primer pairs that are capable of hybridizing to the proximal sequences that define a polymorphism in its double-stranded form.

In lieu of PCR, alternative methods, such as the "Ligase Chain Reaction" ("LCR") may be used (Barany, *Proc. Natl. Acad. Sci. USA* 88:189-193 (1991)). LCR uses two pairs of oligonucleotide probes to exponentially amplify a specific target. The sequences of each pair of oligonucleotides is selected to permit the pair to hybridize to abutting sequences of the same strand of the target. Such hybridization forms a substrate for a template-dependent ligase. As with PCR, the resulting products thus serve as a template in subsequent cycles and an exponential amplification of the desired sequence is obtained.

LCR can be performed with oligonucleotides having the proximal and distal sequences of the same strand of a polymorphic site. In one embodiment, either oligonucleotide will be designed to include the actual polymorphic site of the polymorphism. In such an embodiment, the reaction conditions are selected such that the oligonucleotides can be ligated together only if the target molecule either contains or lacks the specific nucleotide that is complementary to the polymorphic site present on the oligonucleotide. Alternatively, the oligonucleotides may be selected such that they do not include the polymorphic site (see, Segev, PCT Application WO 90/01069).

The "Oligonucleotide Ligation Assay" ("OLA") may alternatively be employed (Landegren et al., *Science* 241:1077-1080 (1988)). The OLA protocol uses two oligonucleotides which are designed to be capable of hybridizing to abutting sequences of a single strand of a target. OLA, like LCR, is particularly suited for the detection of point mutations. Unlike LCR, however, OLA results in "linear" rather than exponential amplification of the target sequence.

Nickerson et al. have described a nucleic acid detection assay that combines attributes of PCR and OLA (Nickerson et al., *Proc. Natl. Acad. Sci. USA* 87:8923-8927 (1990)). In this method, PCR is used to achieve the exponential amplification of target DNA, which is then detected using OLA. In addition to requiring multiple, and separate, processing steps, one problem associated with such combinations is that they inherit all of the problems associated with PCR and OLA.

Schemes based on ligation of two (or more) oligonucleotides in the presence of nucleic acid having the sequence of the resulting "di-oligonucleotide", thereby amplifying the di-oligonucleotide, are also known (Wu et al., *Genomics* 4:560 (1989)), and may be readily adapted to the purposes of the present invention.

Other known nucleic acid amplification procedures, such as allele-specific oligomers, branched DNA technology, transcription-based amplification systems, or isothermal amplification methods may also be used to amplify and analyze such polymorphisms (Malek et al., U.S. Pat. No. 5,130,238; Davey et al., European Patent Application 329,822; Schuster et al., U.S. Pat. No. 5,169,766; Miller et al., PCT Application WO 89/06700; Kwok et al., *Proc. Natl. Acad. Sci. USA* 86:1173-1177 (1989); Gingeras et al., PCT Application WO 88/10315; Walker et al., *Proc. Natl. Acad. Sci. USA* 89:392-396 (1992)).

The identification of a polymorphism can be determined in a variety of ways. By correlating the presence or absence of it in a plant with the presence or absence of a phenotype, it is possible to predict the phenotype of that plant. If a polymorphism creates or destroys a restriction endonuclease cleavage site, or if it results in the loss or insertion of DNA (e.g., a VNTR polymorphism), it will alter the size or profile of the DNA fragments that are generated by digestion with that restriction endonuclease. As such, individuals that possess a variant sequence can be distinguished from those having the original sequence by restriction fragment analysis. Polymorphisms that can be identified in this manner are termed "restriction fragment length polymorphisms" ("RFLPs"). RFLPs have been widely used in human and plant genetic analyses (Glassberg, UK Patent Application 2135774; Skolnick et al., *Cytogen. Cell Genet* 32:58-67 (1982); Botstein et al., *Ann. J. Hum. Genet.* 32:314-331 (1980); Fischer et al. PCT Application WO90/13668; Ublen, PCT Application WO90/11369).

Polymorphisms can also be identified by Single Strand Conformation Polymorphism (SSCP) analysis. The SSCP technique is a method capable of identifying most sequence variations in a single strand of DNA, typically between 150 and 250 nucleotides in length (Elles, *Methods in Molecular Medicine: Molecular Diagnosis of Genetic Diseases*, Humana Press (1996)); Orita et al., *Genomics* 5:874-879 (1989)). Under denaturing conditions a single strand of DNA will adopt a conformation that is uniquely dependent on its sequence conformation. This conformation usually will be different, even if only a single base is changed. Most conformations have been reported to alter the physical configuration or size sufficiently to be detectable by electrophoresis. A number of protocols have been described for SSCP including, but not limited to Lee et al., *Anal. Biochem.* 205:289-293 (1992); Suzuki et al., *Anal. Biochem.* 192:82-84 (1991); Lo et al., *Nucleic Acids Research* 20:1005-1009 (1992); Sarkar et al., *Genomics* 13:441-443 (1992)). It is understood that one or more of the nucleic acids of the present invention, may be utilized as markers or probes to detect polymorphisms by SSCP analysis.

Polymorphisms may also be found using a DNA fingerprinting technique called amplified fragment length polymor-

phism (AFLP), which is based on the selective PCR amplification of restriction fragments from a total digest of genomic DNA to profile that DNA. Vos et al., *Nucleic Acids Res.* 23:44074414 (1995). This method allows for the specific co-amplification of high numbers of restriction fragments, which can be visualized by PCR without knowledge of the nucleic acid sequence.

AFLP employs basically three steps. Initially, a sample of genomic DNA is cut with restriction enzymes and oligonucleotide adapters are ligated to the restriction fragments of the DNA. The restriction fragments are then amplified using PCR by using the adapter and restriction sequence as target sites for primer annealing. The selective amplification is achieved by the use of primers that extend into the restriction fragments, amplifying only those fragments in which the primer extensions match the nucleotide flanking the restriction sites. These amplified fragments are then visualized on a denaturing polyacrylamide gel.

AFLP analysis has been performed on *Salix* (Beismann et al., *Mol. Ecol.* 6:989-993 (1997)); *Acinetobacter* (Janssen et al., *Int. J. Syst. Bacteriol* 47:1179-1187 (1997)), *Aeromonas popoffi* (Huys et al., *Int. J. Syst. Bacteriol.* 47:1165-1171 (1997)), rice (McCouch et al., *Plant Mol. Biol.* 35:89-99 (1997)); Nandi et al., *Mol. Gen. Genet.* 255:1-8 (1997); Cho et al., *Genome* 39:373-378 (1996)), barley (*Hordeum vulgare*) (Simons et al., *Genomics* 44:61-70 (1997); Waugh et al., *Mol. Gen. Genet.* 255:311-321 (1997); Qi et al., *Mol. Gen. Genet.* 254:330-336 (1997); Becker et al., *Mol. Gen. Genet.* 249:65-73 (1995)), potato (Van der Voort et al., *Mol. Gen. Genet.* 255:438447 (1997); Meksem et al., *Mol. Gen. Genet.* 249:74-81 (1995)), *Phytophthora infestans* (Van der Lee et al., *Fungal Genet. Biol.* 21:278-291 (1997)), *Bacillus anthracis* (Keim et al., *J. Bacteriol.* 179:818-824 (1997)), *Astragalus cremnophylax* (Travis et al., *Mol. Ecol.* 5:735-745 (1996)), *Arabidopsis* (Cnops et al., *Mol. Gen. Genet.* 253:3241 (1996)), *Escherichia coli* (Lin et al., *Nucleic Acids Res.* 24:3649-3650 (1996)), *Aeromonas* (Huys et al., *Int. J. Syst. Bacteriol.* 46:572-580 (1996)), nematode (Folkertsma et al., *Mol. Plant Microbe Interact.* 9:47-54 (1996)), tomato (Thomas et al., *Plant J.* 8:785-794 (1995)), and human (Latorra et al., *PCR Methods Appl.* 3:351-358 (1994)). AFLP analysis has also been used for fingerprinting mRNA (Money et al., *Nucleic Acids Res.* 24:2616-2617 (1996); Bachem, et al., *Plant J.* 9:745-753 (1996)). It is understood that one or more of the promoter sequences of the present invention, may be utilized as markers or probes to detect polymorphisms by AFLP analysis for fingerprinting mRNA.

Polymorphisms may also be found using random amplified polymorphic DNA (RAPD) (Williams et al., *Nucl. Acids Res.* 18:6531-6535 (1990)) and cleavable amplified polymorphic sequences (CAPS) (Lyamichev et al., *Science* 260:778-783 (1993)). It is understood that one or more of the promoter sequences of the present invention, may be utilized as markers or probes to detect polymorphisms by RAPD or CAPS analysis.

Promoter sequences of the present invention can be used to in a microarray-based method for high-throughput screening of plant genomic DNA. This 'chip'-based approach involves using microarrays of nucleic acid molecules as gene-specific hybridization targets to identify and quantitatively measure the corresponding plant genes (Schena et al., *Science* 270:467-470 (1995); Shalon, Ph.D. Thesis. Stanford University (1996)). Every nucleotide in a large sequence can be queried at the same time. Hybridization can be used to efficiently analyze nucleotide sequences.

Several microarray methods have been described. For example, microarrays of BACs may be prepared to suffi-

ciently cover 3× of an entire genome. Such microarrays can be used in a variety of genomics experiments including gene mapping, DNA fingerprinting and promoter identification. Microarrays of genomic DNA can also be used for parallel analysis of genomes at single gene resolution (Lemieux et al., *Molecular Breeding* 277-289 (1988)). It is understood that one or more of the molecules of the present invention, preferably one or more of the promoter sequences of the present invention may be utilized in a genomic microarray based method. In a preferred embodiment of the present invention, one or more of the rice genomic promoter sequences may be utilized in a genomic microarray based method. For example, Genomic Mismatch Scanning (GMS), a hybridization-based method of linkage analysis that allows rapid identification of regions of identity-by-descent between two related individuals, can be carried out with microarrays. GMS is reported to have been used to identify genetically common chromosomal segments based on the ability of these DNA sequences to form extensive regions of mismatch-free heteroduplexes. A series of enzymatic steps, coupled with filter binding, is used to selectively remove heteroduplexes that contain mismatches (i.e., chromosomal regions that do not share identity-by descent.). Fragments of chromosomal DNA representing inherited regions are hybridized to a microarray of ordered genomic clones and positive hybridization signals pinpoint regions of identity-by-descent at high resolution (Lemieux et al., *Molecular Breeding* 277-289 (1988)).

It is understood that one or more of the nucleic acid molecules of the present invention may be utilized in a GMS microarray based method to locate regions of identity-by-descent between related individuals. In a preferred embodiment of the present invention, one or more of the nucleic acid molecules of the present invention may be utilized in a GMS microarray based method to locate regions of identity-by-descent between related individuals.

A particularly preferred microarray embodiment of the present invention is a microarray comprising nucleic acid molecules that are homologues of known sequences but elicit only limited or no matches to known nucleic acid molecules. A further preferred microarray embodiment of the present invention is a microarray comprising genomic nucleic acid molecules of the present invention that elicit only limited or no matches to known genes.

It is understood that one or more of the molecules of the present invention, preferably one or more of the promoter sequences of the present invention may be utilized in a microarray based method.

In a preferred embodiment of the present invention, one or more of the nucleic acid molecules of the present invention may be utilized in a microarray based method.

#### Computer Related Uses of the Invention

A nucleic acid molecule comprising SEQ ID NO:1 through SEQ ID NO: 57,467, complements thereof and fragments of either, or a nucleic acid molecule that hybridizes under stringent conditions with SEQ ID NO: 1 through SEQ ID NO:57, 467, or any complement thereof; or exhibits 85% or greater identity, and more preferably at least 86 or greater, 87 or greater, 88 or greater, 89 or greater, 90 or greater, 91 or greater, 92 or greater, 93 or greater, 94 or greater, 95 or greater, 96 or greater, 97 or greater, 98 or greater, or 99% or greater identity to SEQ ID NO:57,467; can be "provided" in a variety of mediums to facilitate its. Such a medium can also provide a subset thereof in a form that allows a skilled artisan to examine the sequences.

In a preferred embodiment, at least 20, 50, 100, 500, 1,000, 2,000, 3,000, or 4,000 of the nucleic acid sequences of the present invention are provided in a variety of mediums. In one

application of this embodiment, a nucleotide sequence of the present invention can be recorded on computer readable media. As used herein, "computer readable media" refers to any medium that can be read and accessed directly by a computer. Such media include, but are not limited to: magnetic storage media, such as floppy discs, hard disc, storage medium, and magnetic tape; optical storage media such as CD-ROM; electrical storage media such as RAM and ROM; and hybrids of these categories such as magnetic/optical storage media. A skilled artisan can readily appreciate how any of the presently known computer readable mediums can be used to create a manufacture comprising computer readable medium having recorded thereon a nucleotide sequence of the present invention.

As used herein, "recorded" refers to a process for storing information on computer readable medium. A skilled artisan can readily adopt any of the presently known methods for recording information on computer readable medium to generate media comprising the nucleotide sequence information of the present invention. A variety of data storage structures are available to a skilled artisan for creating a computer readable medium having recorded thereon a nucleotide sequence of the present invention. The choice of the data storage structure will generally be based on the means chosen to access the stored information. In addition, a variety of data processor programs and formats can be used to store the nucleotide sequence information of the present invention on computer readable medium. The sequence information can be represented in a word processing text file, formatted in commercially-available software such as WordPerfect and Microsoft Word, or represented in the form of an ASCII file, stored in a database application, such as DB2, Sybase, Oracle, or the like. A skilled artisan can readily adapt any number of data processor structuring formats (e.g., text file or database) in order to obtain computer readable medium having recorded thereon the nucleotide sequence information of the present invention.

By providing one or more of nucleotide sequences of the present invention, a skilled artisan can routinely access the sequence information for a variety of purposes. Computer software is publicly available which allows a skilled artisan to access sequence information provided in a computer readable medium. The examples which follow demonstrate how software which implements the BLAST (Altschul et al., *J. Mol. Biol.* 215:403410(1990)) and BLAZE (Brutlag et al., *Comp. Chem.* 17:203-207 (1993)) search algorithms on a Sybase system can be used to identify open reading frames (ORFs) within the genome that contain homology to ORFs or proteins from other organisms. Such ORFs are protein-encoding fragments within the sequences of the present invention and are useful in producing commercially important proteins such as enzymes used in amino acid biosynthesis, metabolism, transcription, translation, RNA processing, nucleic acid and a protein degradation, protein modification, and DNA replication, restriction, modification, recombination, and repair.

The present invention further provides systems, particularly computer-based systems, which contain the sequence information described herein. Such systems are designed to identify commercially important fragments of the nucleic acid molecule of the present invention. As used herein, "a computer-based system" refers to the hardware means, software means, and data storage means used to analyze the nucleotide sequence information of the present invention. The minimum hardware means of the computer-based systems of the present invention comprises a central processing unit (CPU), input means, output means, and data storage means. A skilled artisan can readily appreciate that any one of

the currently available computer-based system are suitable for use in the present invention.

As indicated above, the computer-based systems of the present invention comprise a data storage means having stored therein a nucleotide sequence of the present invention and the necessary hardware means and software means for supporting and implementing a search means. As used herein, "data storage means" refers to memory that can store nucleotide sequence information of the present invention, or a memory access means which can access manufactures having recorded thereon the nucleotide sequence information of the present invention. As used herein, "search means" refers to one or more programs which are implemented on the computer-based system to compare a target sequence or target structural motif with the sequence information stored within the data storage means. Search means are used to identify fragments or regions of the sequence of the present invention that match a particular target sequence or target motif. A variety of known algorithms are disclosed publicly and a variety of commercially available software for conducting search means are available and can be used in the computer-based systems of the present invention. Examples of such software include, but are not limited to, MacPattern (EMBL), BLASTIN and BLASTIX (NCBIA). One of the available algorithms or implementing software packages for conducting homology searches can be adapted for use in the present computer-based systems.

As used herein, "a target structural motif," or "target motif," refers to any rationally selected sequence or combination of sequences in which the sequence(s) are chosen based on primary sequence composition or a three dimensional configuration which is formed upon folding of the target motif. There are a variety of target motifs known in the art. Target motifs include, but are not limited to, transcription factor binding sites, repressor binding sites, inducible expression elements, transcriptional activation sites, transcription initiation sites, untranslated leaders, intron splicing sites, methylation sites, histone binding sites, RNA processing sites, non-histone structural protein binding sites, replication sites, sites which influence the stability of transcribed mRNA message and hairpin sites.

Thus, the present invention further provides an input means for receiving a target sequence, a data storage means for storing the target sequences of the present invention sequence identified using a search means as described above, and an output means for outputting the identified homologous sequences. A variety of structural formats for the input and output means can be used to input and output information in the computer-based systems of the present invention. A preferred format for an output means ranks fragments of the sequence of the present invention by varying degrees of homology to the target sequence or target motif. Such presentation provides a skilled artisan with a ranking of sequences which contain various amounts of the target sequence or target motif and identifies the degree of homology contained in the identified fragment.

A variety of comparing means can be used to compare a target sequence or target motif with the data storage means to identify sequence fragments sequence of the present invention. For example, implementing software which implement the BLAST and BLAZE algorithms (Altschul et al., *J. Mol. Biol.* 215:403-410 (1990)) can be used to identify open frames within the nucleic acid molecules of the present invention. A skilled artisan can readily recognize that any one of the publicly available homology search programs can be used as the search means for the computer-based systems of the present invention.

Having now generally described the invention, the same will be more readily understood through reference to the following examples which are provided by way of illustration, and are not intended to be limiting of the present invention, unless specified.

Each periodical, patent, and other document or reference cited herein is herein incorporated by reference in its entirety.

## EXAMPLES

The following examples are provided to better illustrate the practice of the present invention and should not be interpreted in any way to limit the scope of the present invention. Those skilled in the art will recognize that various modifications, truncations, etcetera can be made to the methods and genes described herein while not departing from the spirit and scope of the present invention.

### Example 1

#### Generating a Genomic Bacterial Artificial Chromosome (BAC) Library

BACs are stable, non-chimeric cloning systems having genomic fragment inserts (100-300 kb) and their DNA can be prepared for most types of experiments including DNA sequencing. BAC vector, pBeloBAC11, is derived from the endogenous *E. coli* F-factor plasmid, which contains genes for strict copy number control and unidirectional origin of DNA replication. Additionally, pBeloBAC11 has three unique restriction enzyme sites (Hind III, Bam HI and Sph I) located within the LacZ gene which can be used as cloning sites for megabase-size plant DNA. Indigo, another BAC vector contains Hind III and Eco RI cloning sites. This vector also contains a random mutation in the LacZ gene that allows for darker blue colonies.

As an alternative, the P1-derived artificial chromosome (PAC) can be used as a large DNA fragment cloning vector (Ioannou, et al., *Nature Genet.* 6:84-89 (1994); Suzuki, et al., *Gene* 199:133-137 (1997)). The PAC vector has most of the features of the BAC system, but also contains some of the elements of the bacteriophage P1 cloning system.

BAC libraries are generated by ligating size-selected restriction digested DNA with pBeloBAC 11 followed by electroporation into *E. coli*. BAC library construction and characterization is extremely efficient when compared to YAC (yeast artificial chromosome) library construction and analysis, particularly because of the chimerism associated with YACs and difficulties associated with extracting YAC DNA.

There are general methods for preparing megabase-size DNA from plants. For example, the protoplast method yields megabase-size DNA of high quality with minimal breakage. A process involves preparing young leaves which are manually feathered with a razor-blade before being incubated for four to five hours with cell-wall-degrading enzymes. A second method developed by Zhange et al., *Plant J.* 7:175-184 (1995), is a universal nuclei method that works well for several divergent plant taxa. Fresh or frozen tissue is homogenized with a blender or mortar and pestle. Nuclei are then isolated and embedded. DNA prepared by the nucleic method is often more concentrated and is reported to contain lower amounts of chloroplast DNA than the protoplast method.

Once protoplasts or nuclei are produced, they are embedded in an agarose matrix as plugs or microbeads. The agarose provides a support matrix to prevent shearing of the DNA while allowing enzymes and buffers to diffuse into the DNA.

The DNA is purified and manipulated in the agarose and is stable for more than one year at 4° C.

Once high molecular weight DNA is prepared, it is fragmented to the desired size range. In general, DNA fragmentation utilizes two general approaches, 1) physical shearing and 2) partial digestion with a restriction enzyme that cuts relatively frequently within the genome. Since physical shearing is not dependent upon the frequency and distribution of particular restriction enzymes sites, this method should yield the most random distribution of DNA fragments. However, the ends of the sheared DNA fragments must be repaired and cloned directly or restriction enzyme sites added by the addition of synthetic linkers. Because of the subsequent steps required to clone DNA fragmented by shearing, most protocols fragment DNA by partial restriction enzyme digestion. The advantage of partial restriction enzyme digestion is that no further enzymatic modification of the ends of the restriction fragments are necessary. Four common techniques that can be used to achieve reproducible partial digestion of megabase-size DNA are 1) varying the concentration of the restriction enzyme, 2) varying the time of incubation with the restriction enzyme 3) varying the concentration of an enzyme cofactor (e.g., Mg<sup>2+</sup>) and 4) varying the ratio of endonuclease to methylase.

There are three cloning sites in pBeloBAC11, but only Hind III and Barn HI produce 5' overhangs for easy vector dephosphorylation. These two restriction enzymes are primarily used to construct BAC libraries. The optimal partial digestion conditions for megabase-size DNA are determined by wide and narrow window digestions. To optimize the optimum amount of Hind III, 1, 2, 3, 10, and 5- units of enzyme are each added to 50 ml aliquots of microbeads and incubated at 37° C. for 20 minutes.

After partial digestion of megabase-size DNA, the DNA is run on a pulsed-field gel, and DNA in a size range of 100-500 kb is excised from the gel. This DNA is ligated to the BAC vector or subjected to a second size selection on a pulsed field gel under different running conditions. Studies have previously reported that two rounds of size selection can eliminate small DNA fragments co-migrating with the selected range in the first pulse-field fractionation. Such a strategy results in an increase in insert sizes and a more uniform insert size distribution. A practical approach to performing size selections is to first test for the number of clones/microliter of ligation and insert size from the first size selected material. If the numbers are good (500 to 2000 white colony/microliter of ligation) and the size range is also good (50 to 300 kb) then a second size selection is practical. When performing a second size selection one expects a 80 to 95% decrease in the number of recombinant clones per transformation.

Twenty to two hundred nanograms of the size-selected DNA is ligated to dephosphorylated BAC vector (molar ratio of 10 to 1 in BAC vector excess). Most BAC libraries use a molar ratio of 5 to 15:1 (size selected DNA:BAC vector).

Transformation is carried out by electroporation and the transformation efficiency for BACs is about 40 to 1,500 transformants from one microliter of ligation product or 20 to 1000 transformants/ng DNA.

Several tests can be carried out to determine the quality of a BAC library. Three basic tests to evaluate the quality include: the genome coverage of a BAC library-average insert size, average number of clones hybridizing with single copy probes and chloroplast DNA content.

The determination of the average insert size of the library is assessed in two ways. First, during library construction every ligation is tested to determine the average insert size by assaying 20-50 BAC clones per ligation. DNA is isolated from

recombinant clones using a standard mini preparation protocol, digested with Not I to free the insert from the BAC vector and then sized using pulsed field gel electrophoresis (Maule, *Molecular Biotechnology* 9:107-126 (1998)).

To determine the genome coverage of the library, it is screened with single copy RFLP markers distributed randomly across the genome by hybridization. Microtiter plates containing BAC clones are spotted onto Hybond membranes. Bacteria from 48 or 72 plates are spotted twice onto one membrane resulting in 18,000 to 27,648 unique clones on each membrane in either a 4x4 or 5x5 orientation. Since each clone is present twice, false positives are easily eliminated and true positives are easily recognized and identified.

Finally, the chloroplast DNA content in the BAC library is estimated by hybridizing three chloroplast genes spaced evenly across the chloroplast genome to the library on high density hybridization filters.

There are strategies for isolating rare sequences within the genome. For example, higher plant genomes can range in size from 100 Mb/1C (*Arabidopsis*) to 15,966 Mb/C (*Triticum aestivum*), (Arumuganathan and Earle, *Plant Mol Bio Rep.* 9: 208219 (1991)). The number of clones required to achieve a given probability that any DNA sequence will be represented in a genomic library is  $N = (1/P) / (1 - L/G)$  where N is the number of clones required, P is the probability desired to get the target sequence, L is the length of the average clone insert in base pairs and G is the haploid genome length in base pairs (Clarke et al., *Cell* 9:91-100 (1976)).

The rice BAC library of the present invention is constructed in the pBeloBAC11 or similar vector. Inserts are generated by partial Eco RI or other enzymatic digestion of DNA. The 25x library provides 4-5x coverage sequence from BAC clones across genome.

#### Example 2

##### Sequencing Genomic DNA Inserts from a Genomic BAC Library

Two basic methods can be used for DNA sequencing, the chain termination method of Sanger et al., *Proc. Natl. Acad. Sci. USA* 74:5463-5467 (1977), and the chemical degradation method of Maxam and Gilbert, *Proc. Natl. Acad. Sci. USA* 74:560-564 (1977). Automation and advances in technology such as the replacement of radioisotopes with fluorescence-based sequencing have reduced the effort required to sequence DNA (Craxton, *Methods*, 2:20-26 (1991); Ju et al., *Proc. Natl. Acad. Sci. USA* 92:4347-4351 (1995); Tabor and Richardson, *Proc. Natl. Acad. Sci. USA* 92:6339-6343 (1995)). Automated sequencers are available from, for example, Pharmacia Biotech, Inc., Piscataway, N.J. (Pharmacia ALF), LI-COR, Inc., Lincoln, Nebr. (LI-COR 4,000) and Millipore, Bedford, Mass. (Millipore BaseStation).

In addition, advances in capillary gel electrophoresis have also reduced the effort required to sequence DNA and such advances provide a rapid high resolution approach for sequencing DNA samples (Swerdlow and Gesteland, *Nucleic Acids Res.* 18:1415-1419 (1990); Smith, *Nature* 349:812-813 (1991); Luckey et al., *Methods Enzymol.* 218:154-172 (1993); Lu et al., *J. Chromatog. A.* 680:497-501 (1994); Carson et al., *Anal. Chem.* 65:3219-3226 (1993); Huang et al., *Anal. Chem.* 64:2149-2154 (1992); Kheterpal et al., *Electrophoresis* 17:1852-1859 (1996); Quesada and Zhang, *Electrophoresis* 17:1841-1851 (1996); Baba, *Yakugaku Zasshi* 117: 265-281 (1997)).

A number of sequencing techniques are known in the art, including fluorescence-based sequencing methodologies. These methods have the detection, automation and instru-

mentation capability necessary for the analysis of large volumes of sequence data. Currently, the 377.DNA Sequencer (Perkin-Elmer Corp., Applied Biosystems Div., Foster City, Calif.) allows the most rapid electrophoresis and data collection. With these types of automated systems, fluorescent dye-labeled sequence reaction products are detected and data entered directly into the computer, producing a chromatogram that is subsequently viewed, stored, and analyzed using the corresponding software programs. These methods are known to those of skill in the art and have been described and reviewed (Birren et al., *Genome Analysis: Analyzing DNA*, 1, Cold Spring Harbor, N.Y. 1999).

PHRED is used to call the bases from the sequence trace files ([www-mbt.washington.edu](http://www-mbt.washington.edu)). Phred uses Fourier methods to examine the four base traces in the region surrounding each point in the data set in order to predict a series of evenly spaced predicted locations. That is, it determines where the peaks would be centered if there were no compressions, drop-outs, or other factors shifting the peaks from their "true" locations. Next, PHRED examines each trace to find the centers of the actual, or observed peaks and the areas of these peaks relative to their neighbors. The peaks are detected independently along each of the four traces so many peaks overlap. A dynamic programming algorithm is used to match the observed peaks detected in the second step with the predicted peak locations found in the first step.

After the base calling is completed, contaminating sequences (E. coli, BAC vector sequences >50 bases and sub-cloning vector are removed and constraints are made for the assembler. Contigs are assembled using CAP3 (Huang, et al., *Genomics* 46: 37-45 (1997)).

A two-step re-assembly process is employed to reduce sequence redundancies caused by overlaps between BAC clones. In the first step, BAC clones are grouped into clusters based on overlaps between contig sequences from different BACs. These overlaps are identified by comparing each sequence in the dataset against every other sequences, by BLASTN. BACs containing overlaps greater than 5,000 base pairs in length and greater than 94% in sequence identity are put into the same cluster. Repetitive sequences are masked prior to this procedure to avoid false joining by repetitive elements present in the genome. In the second step, sequences from each BAC cluster are assembled by PHRAP.longread, which is able to handle very long sequences. A minimum match is set at 100 bp and a minimum score is set at 600 as a threshold to join input contigs into longer contigs.

### Example 3

#### Identifying Genes within a Genomic BAC Library

This example illustrates the identification of combigens within the rice genomic contig library as assembled in Example 2. The genes and partial genes that are embedded in such contigs are identified through a series of informatic analyses. The tools to define genes fall into two categories: homology-based and predictive-based methods. Homology-based searches (e.g., GAP2, BLASTX supplemented by NAP and TBLASTX) detect conserved sequences during comparisons of DNA sequences or hypothetically translated protein sequences to public and/or proprietary DNA and protein databases. Existence of an *Oryza sativa* gene is inferred if significant sequence similarity extends over the majority of the target gene. Since homology-based methods may overlook genes unique to *Oryza sativa*, for which homologous nucleic acid molecules have not yet been identified in databases, gene prediction programs are also used. Predictive methods employed in the definition of the *Oryza sativa* genes included

the use of the GenScan gene predictive software program which is available from Stanford University (e.g., at the website: [www-gnomc.stanford.edu/GENSCANW.html](http://www-gnomc.stanford.edu/GENSCANW.html). and the Genemark.hmm for Eukaryotes program from Gene Probe, Inc (Atlanta, Ga.) [www-geneprobe.net/index.htm](http://www-geneprobe.net/index.htm)). GenScan, in general terms, infers the presence and extent of a gene through a search for "gene-like" grammar. GeneMark.hmm searches a file containing DNA sequence data for genes. It employs a Hidden Markov Model algorithm with a species-specific inhomogeneous Markov model of gene-encoding regions of DNA.

The homology-based methods that are used to define the *Oryza sativa* gene set included GAP2, BLASTX supplemented by NAP and TBLASTX. For a description of BLASTX and TBLASTX see Coulson, *Trends in Biotechnology* 12:76-80 (1994) and Birren et al., *Genome Analysis*, 1:543-559 (1997). GAP2 and NAP are part of the Analysis and Annotation Tool (AAT) for Finding Genes in Genomic Sequences which was developed by Xiaohu Huang at Michigan Tech University and is available at the web site [www-genome.cs.mtu.edu/](http://www-genome.cs.mtu.edu/). The AAT package includes two sets of programs, one set DPS/NAP (referred to as "NAP") for comparing the query sequence with a protein database, and the other set DDS/GAP2 (referred to as "GAP2") for comparing the query sequence with a cDNA database. Each set contains a fast database search program and a rigorous alignment program. The database search program identifies regions of the query sequence that are similar to a database sequence. Then the alignment program constructs an optimal alignment for each region and the database sequence. The alignment program also reports the coordinates of exons in the query sequence. See Huang, et al., *Genomics* 46: 37-45 (1997). The GAP2 program computes an optimal global alignment of a genomic sequence and a cDNA sequence without penalizing terminal gaps. A long gap in the cDNA sequence is given a constant penalty. The DNA-DNA alignment by GAP2 adjusts penalties to accommodate introns. The GAP2 program makes use of splice site consensus in alignment computation. GAP2 delivers the alignment in linear space, so long sequences can be aligned. See Huang, *Computer Applications in the Biosciences* 10 227-235 (1994). The GAP2 program aligns the *Oryza sativa* contigs with a library of 42,260 *Oryza sativa* cDNAs.

The NAP program computes a global alignment of a DNA sequence and a protein sequence without penalizing terminal gaps. NAP handles frameshifts and long introns in the DNA sequence. The program delivers the alignment in linear space, so long sequences can be aligned. It makes use of splice site consensus in alignment computation. Both strands of the DNA sequence are compared with the protein sequence and one of the two alignments with the larger score is reported. See Huang, and Zhang, *Computer Applications in the Biosciences* 12(6), 497-506 (1996).

NAP takes a nucleotide sequence, translates it in three forward reading frames and three reverse complement reading frames, and then compares the six translations against a protein sequence database (e.g. the non-redundant protein (ie., nr-aa database maintained by the National Center for Biotechnology Information as part of GenBank and available at the web site: [www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)).

The first homology-based search for genes in the *Oryza saliva* contigs is effected using the GAP2 program and the *Oryza saliva* library of clustered *Oryza saliva* cDNA. The *Oryza saliva* clusters are mapped onto an assembly of *Oryza saliva* contigs using the GAP2 program. GAP2 standards for

selecting a DNA-DNA match are  $\geq 92\%$  sequence identity with the following parameters:

gap extension penalty=1  
 match score=2  
 gap open penalty=6  
 gap length for constant penalty=20  
 mismatch penalty=2  
 minimum exon length=21  
 minimum total length of all exons in a gene (in nucleotide)=200

When a particular *Oryza sativa* cDNA aligns to more than one *Oryza sativa* contig, the alignment with the highest identity is selected and alignments with lower levels of identity are filtered out as surreptitious alignments. *Oryza sativa* cDNA sequences aligning to *Oryza sativa* contigs with exceptionally low complexity are filtered out when the basis for alignment included a high number of cDNAs with poly A tails aligning to genomic regions with extended repeats of A or T.

The second homology-based method used for gene discovery is BLASTX hits extended with the NAP software package. BLASTX is run with the *Oryza sativa* genomic contigs as queries against the GenBank non-redundant protein data library identified as "nr-aa". NAP is used to better align the amino acid sequences as compared to the genomic sequence. NAP extends the match in regions where BLASTX has identified high-scoring-pairs (HSPs), predicts introns, and then links the exons into a single ORF prediction. Experience suggests that NAP tends to mis-predict the first exon. The NAP parameters are:

gap extension penalty=1  
 gap open penalty=15  
 gap length for constant penalty=25  
 min exon length (in aa)=7  
 minimum total length of all exons in a gene (in nucleotide)=200  
 homology>40%

The NAP alignment score and GenBank reference number for best match are reported for each contig for which there is a NAP hit.

In the final homology-based method, TBLASTX, is used with cDNA information from four plant sequencing projects: 27,037 sequences from *Triticum aestivum*, 136,074 sequences from *Glycine max*, 71,822 sequences from *Zea mays* and 68,517 sequences from *Arabidopsis thaliana*. Conservative standards for inclusion of TBLASTX hits into the gene set are utilized. These standards are a minimal E value of  $1E-16$ , and a minimal match of 150 bp in *Oryza saliva* contig.

The GenScan program is "trained" with *Arabidopsis thaliana* characteristics. Though better than the "off-the-shelf" version, the GenScan trained to identify *Oryza saliva* genes proved more proficient at predicting exons than predicting full-length genes. Predicting full-length genes is compromised by point mutations in the unfinished contigs, as well as by the short length of the contigs relative to the typical length of a gene. Due to the errors found in the full-length gene predictions by GenScan, inclusion of GenScan-predicted genes is limited to those genes and exons whose probabilities are above a conservative probability threshold. The GenScan parameters are:

weighted mean GenScan P value>0.4  
 mean GenScan T value>0  
 mean GenScan Coding score>50  
 length>200 bp  
 minimum total length of all exons in a gene=500

The weighted mean GenScan P value is a probability for correctly predicting ORFs or partial ORFs and is defined as

the  $(l/SS li)(SS li Pi)$ , where "l" is the length of an exon and "P" is the probability or correctness for the exon.

The GeneMark.hmm for Eukaryotes program uses the Hidden Markov model for species *Oryza Sativa*. Minimum total length of all exons in a gene is 500 bp. Except for the model selection, there is no specific run-time parameter for GeneMark.hmm.

The gene predictions from these programs are stored in a database and then combigenes are derived from these predictions. A combigene is a cluster of putative genes which satisfy the following criteria:

All genes making up a single combigene are located on the same strand of a contig;

Maximum intron size of a valid gene is 4000 bp;

Maximum distance between any two genes in the same combigene is 200 bp, as measured by the bases between adjacent ending exons;

If an individual gene is predicted by NAP it has at least 40% sequence identity to its hit;

If an individual gene is predicted by GAP2 it has at least 92% sequence identity to its hit;

If an individual gene is predicted by Genscan the weighted average of the probabilities calculated for all of its exons is not less than 0.4. The gene boundaries of a Genscan-predicted gene are determined while taking into account only exons.

Since TBLASTX-predicted genes are standless the combigene which is made up of such genes can be assigned a strand only if there is a gene in the cluster that was predicted by a strand-defining gene-predicting program.

#### Example 4

Identifying Promoters in the Genomic BAC Library Using Bioinformatic Techniques

Candidate promoter sequences are selected by identifying the regions of DNA located immediately upstream of "combigenes" as described and defined in Example 3. The length of the region to be extracted from the corresponding contig's sequence is set to be 1500 nucleotides plus the very first nucleotide of a combigene. Thus, if a combigene is sufficiently far from the edge of a contig a 1501 nucleotide sequence is obtained, otherwise the sequence will be shorter. Only coding region predictions are considered when building combigenes. Therefore, the 5' UTR of the putative cDNA is included as part of the combigene upstream region.

If there is an AAT/NAP-predicted component in a combigene, then the putative promoter sequence is extracted upstream of the beginning of that component otherwise—the sequence is extracted upstream of the beginning of the combigene (which may correspond to Genscan, AAT/GAP or a TBLASTX prediction).

Promoter candidates are further selected using bioinformatic analysis of the candidate promoter sequence.

The candidate promoter regions listed in SEQ ID NO:1 through SEQ ID NO:57467 are analyzed for known promoter motifs listed in Table 2.

The identification of such motifs provides important information about the candidate promoter. For example, some motifs are associated with informative annotations such as "light inducible binding site" or "stress inducible binding motif" and can be used to select with confidence a promoter that is able to confer light inducibility or stress inducibility to an operably-linked transgene, respectively.

Putative promoter sequences are also searched with matrices for the TATA box, GC box (factor name: V\_GC\_0) and CCAAT box (factor name: F\_HAP234\_01). The matrix for the TATA box is from the Eukaryotic Promoter Database

(www.epd.isb-sib.ch/) and the matrices for the GC box and the CCAAT box are from Transfac (www-transfac.gbf.de/TRANSFAC/).

The algorithm that is used to annotate promoters searches for matches to both sequence motifs and matrix motifs. First, individual matches are found. For sequence motifs, a maximum number of mismatches is allowed (see Table 2). If the code M,R,W,S,Y, or K are listed in the sequence motif (each of which is a degenerate code for 2 nucleotides) 1/2 mismatch is allowed. If the code B, D, H, or V are listed in the sequence motif (each of which is a degenerate code for 3 nucleotides) 1/3 mismatch is allowed. p values are determined by simulation with a 5 Mb of random DNA with the same dinucleotide frequency as the test set is generated and the probability of a given matrix score is determined (number of hits/5e7). Once the individual hits have been found, the putative promoter sequence is searched for clusters of hits in a 250 bp window. The score for a cluster is found by summing the negative natural log of the p value for each individual hit. Using 100 Mb simulations as described above, the probability of a window having a cluster score greater than or equal to the given value is determined. Clusters with a p value more significant than  $p < 1e-6$  are reported. Only the top 287 hits are taken and are ranked by p value. Effects of repetitive elements are screened. If the 287th ranked hit has the same p value as the first ranked hit, no results are reported for that factor.

For matrix motifs, a p value cutoff is used on a matrix score. The matrix score is determined by adding the path of a given DNA sequence through a matrix. P values are determined by simulation; 5 Mb of random DNA with the same dinucleotide frequency as a test set is generated to test individual matrix hits and 100 Mb is used to test clusters; the probability of a given matrix score and the probability scores for clusters are determined as are the sequence motifs. The usual cutoff for matrices is  $2.5e-4$ . No clustering is done for the TATA box, GC box or CCAAT box.

Candidate promoters are also selected based on the expression characteristics of the gene that is cis-associated with the candidate promoter, (i.e. the native gene). For example, a promoter region located 5' to a gene, which is expressed during a specific stage of development, likely plays a key role in the temporal regulation of that gene. Thus the promoter, when operably linked to a heterologous coding sequence, may similarly regulate the heterologous coding sequence.

Combining the motif analysis with the expression analysis, the list of candidate promoters having desired properties can be narrowed. This decreases the overall number of candidate promoters that must be screened to confirm the promoter's function. For example, one can start with seed-expressed transcription factors, identify candidate promoters that match the consensus regulation sites for seed-expressed transcription factors, and then test the identified candidate promoters to confirm the promoter sub-set which are capable of conferring seed-specific expression to a gene.

#### Example 5

Identifying Promoters in the Genomic BAC Library Using an Expression Assay

Promoters may also be identified based on quantitative analysis of genes that are cis-associated with candidate promoters, (i.e. the native genes). In this method, the native genes associated with SEQ ID NO:1 through SEQ ID NO:57,467 are analyzed on a digital northern blot. Digital northern data can be generated from EST sequencing, SAGE and other methods, which in effect count RNA molecules expressed in cell. This data can be generated as needed, or is generally

available to the public on a number of web sites (e.g., www.tigr.org). Data can be obtained from any plant species, although data on rice gene expression is particularly preferred. Promoters are selected based on the expression information of the digital northern. For example, identifying genes expressing genes under stress-related conditions would provide a group of promoters able to confer such stress-inducible expression to other genes.

#### Example 6

Identifying Promoters in the Genomic BAC Library Using Microarray Analysis

Promoters may also be selected based by transcriptional profiling or microarray analysis. Transcriptional profiling can be completed on large scale for each cis-linked gene associated with SEQ ID NO:1 through SEQ ID NO:57,467. Transcription profiling data can be obtained on RNA prepared from any plant species using a chip comprised of sequences from any plant species, although data generated from rice using a rice chip is preferred.

A comprehensive database of transcription profiling data narrows down the list of promoter candidates that confer a desired expression pattern. For example, a promoter that confers drought-specific expression can be selected by identifying a cis-linked gene that is induced under drought conditions (on the microarray), but is not expressed at other stages of plant growth and development. Such a promoter is likely to confer drought inducibility to an operably linked transgene. Public databases of transcript profiling data are becoming more comprehensive and thereby enabling this type of analysis.

#### Example 7

Functional Screening of Promoters in an Expression Assay

Promoters are screened in an expression assay. The promoters in SEQ ID NO:1 through SEQ ID NO:57,467 are amplified by PCR from rice genomic DNA and cloned into an expression vector containing a reporter transgene (e.g., GUS or GFP). The individual promoter or a collection of promoters ("promoter library") are then screened in an expression assay for the ability to express the reporter transgene. In a common expression assay for leaf promoters, the promoters are transfected into rice or maize leaf protoplasts. Reporter gene expression in the protoplasts indicates a promoter capable of conferring gene-expression in the leaf. The promoters are also transfected into protoplasts from other tissues or plant species to identify other regulatory features of the promoter.

Alternatively, promoters may be screened using a particle gun technique to bombard the cells, tissues or plants. The bombarded samples are visually inspected for reporter gene expression. Reporter gene expression observed in any bombarded samples indicates the presence of a promoter able to confer expression of a transgene in that cell, tissue or plant.

The promoters may also be screened in plants where transformation protocols have been greatly enhanced to facilitate the screening of large numbers of promoters. In this approach, the individual rice promoters or "promoter library" is transformed into Arabidopsis plants. The resulting transformed tissues or progeny are scored for reporter expression. Again, reporter gene expression in a given tissue indicates that a promoter is able to confer transgene expression in that tissue.

For some promoters, such as those providing constitutive expression, a reporter transgene can be replaced with a selectable marker transgene, such as a gene conferring glyphosate tolerance. Transformed cells, tissues or plants expressing the

selectable marker are selected, rather than visually scored. For example, the promoter is linked to a selectable marker, such as glyphosate resistance, and then screening for male sterile plants. The selected plants, in this case male sterile plants, may contain a promoter for male reproductive tissues.

The promoters described herein can also be used to ablate or kill cells expressing a gene from the promoter. In such cases, the promoter is operably linked to a negative selectable marker gene, including but not limited to the diphtheria toxin gene, or to a conditional lethal gene, including but not limited to the phosphonate ester hydrolase gene (pehA). The negative selectable marker gene is transformed into cells, tissues or plants. The cells, tissues or plants which express the negative selectable gene from the promoter are selectively killed. In the case of the conditional lethal gene, the transformed cells, tissues or plants which express the conditional lethal gene are only killed in the presence of the negative selective agent or negative selective condition. In the example of the phosphonate ester hydrolase gene, the transformed cells, tissues or plants which express the conditional lethal gene are only killed in the presence of glyceryl glyphoste.

Table 1

The data in Table 1 provides features relating to the putative promoter sequences.

\*column headings

5 Seq Num: Provides the SEQ ID NO. for the rice contigs on which the putative promoter sequences are found.

Contig ID: unique identifier of the rice contig

10 CmbGID: name of the putative promoter sequence. Putative promoters are named as cg\_“no”. The “no” refers to the combigene from which the putative promoter is selected.

CNTG LEN: The length of the contig

BEGN POS: starting position of the putative promoter sequence;

15 STRND: DNA strand on which the combigene is located

LEN: length of the putative promoter sequence

Table 2

Table 2 lists the sequence motifs that are searched in the putative promoter sequences

TABLE 2

SEQ ID NO.	Transcription Factor Name	Sequence Motif name	Sequence Motif	Sequence Motif Length	Maximum mis-matches allowed	Reference for transcription factors and sequence motifs
109670	Fac006	ABADES11	RTACGTGGCR	10	1	PLACE
109671	Fac037	ABADES12	GGACGCGTGGC	11	2	PLACE
109672	Fac010	ABFOS	GCATCTTTACTTTAGCATC	19	6	PLACE
109673	Fac016	ABRE3OSRAB16	GTACGTGGCGC	11	2	PLACE
109674	Fac016	ABREATRD22	RYACGTGGYR	10	0	PLACE
109675	Fac020	ABREOSRAB21	ACGTSSSC	8	0	PLACE
109676	Fac021	ABREOSRGA1	CCACGTGG	8	0	PLACE
109677	Fac021	ABRETAEM	GGACACGTGGC	11	2	PLACE
109678	Fac022	ACGTABOX	TACGTA	6	0	PLACE
109679	Fac031	AMYBOX1	TAACARA	7	0	PLACE
109680	Fac032	AMYBOX2	TATCCAT	7	0	PLACE
109681	Fac060	DREDR1ATRD29AB	TACCGACAT	9	1	PLACE
109682	Fac064	EREGCCNTCHN	TAAGAGCCGCC	11	2	PLACE
109683	Fac066	GARE2R	TAACARANTCYGG	14	2	PLACE
109684	Fac068	GBOXRELOSAMY3	CTACGTGGCCA	11	2	PLACE
109685	Fac070	GLUTAACAOS	AACAACACTCTAT	12	2	PLACE
109686	Fac071	1OSGT2GLUTEBOX	ATATCATGAGTCACTTCA	18	4	PLACE
109687	Fac071	1OSGT2GLUTEBOX	ATATCATGAGTCACTTCA	18	4	PLACE
109688	Fac071	1OSGT3GLUTEBOX	TATCTAGTGAGTCACTTCA	19	5	PLACE
109689	Fac071	1OSGT3GLUTEBOX	TATCTAGTGAGTCACTTCA	19	5	PLACE
109690	Fac072	2OSGT2GLUTEBOX	TCCGTGTACCA	11	2	PLACE
109691	Fac072	2OSGT2GLUTEBOX	TCCGTGTACCA	11	2	PLACE
109692	Fac072	2OSGT3GLUTEBOX	CTTTTGTGTACCTTA	15	3	PLACE
109693	Fac072	2OSGT3GLUTEBOX	CTTTTGTGTACCTTA	15	3	PLACE
109694	Fac073	GLUTEBP1OS	AAGCAACACACAAC	14	3	PLACE
109695	Fac074	GLUTEBP2OS	ATGCTCAATAGATATAAGT	19	5	PLACE
109696	Fac075	GLUTECOREOS	CTTTCGTGTAC	11	2	PLACE
109697	Fac079	GT2OSPBY	AGCGGTAATT	9	1	PLACE
109698	Fac105	MYBGAVH	TAACAAA	7	0	PLACE
109699	Fac129	PROLAMINBOX	CACATGTGTAAAGGT	15	4	PLACE
109700	Fac135	RGATAOS	CAGAAGATA	9	1	PLACE
109701	Fac136	RNF1OS	GATCATCGATC	11	2	PLACE
109702	Fac137	RNF2OS	CCAGTGTGCCCTGG	15	4	PLACE
109703	Fac139	RYREPEAT4	TCCATGCATGCAC	13	3	PLACE
109704	Fac139	RYREPEAT4	TCCATGCATGCAC	13	3	PLACE
109705	Fac139	RYREPEATGMGY2	CATGCAT	7	0	PLACE
109706	Fac139	RYREPEATGMGY2	CATGCAT	7	0	PLACE
109707	Fac139	RYREPEATLEGUMINBOX	CATGCAY	7	0	PLACE
109708	Fac139	RYREPEATLEGUMINBOX	CATGCAY	7	0	PLACE
109709	Fac139	RYREPEATVFLB4	CATGCATG	8	0	PLACE
109710	Fac139	RYREPEATVFLB4	CATGCATG	8	0	PLACE
109711	Fac149	SITEIIAOSPCNA	TGGGCCCGT	9	1	PLACE
109712	Fac150	SITEIIIBOSPCNA	TGGTCCAC	9	1	PLACE
109713	Fac151	SITEIOSPCNA	CCAGGTGG	8	1	PLACE
109714	Fac163	AACAOSGLUB1	CAACAACTATATC	14	3.5	PLACE
109715	Fac165	ACGTOSGLUB1	GTACGTG	7	0	PLACE
109716	Fac180	GT1CONSENSUS	GRWAAW	6	0	PLACE
109717	Fac201	PYRIMIDINEBOXOSRAMY1A	CCTTTT	6	0	PLACE
109718	Fac218	ABREMOTIFAOSOSEM	TACGTGTC	8	0.5	PLACE

TABLE 2-continued

SEQ ID NO.	Transcription Factor Name	Sequence Motif name	Sequence Motif	Sequence Motif Length	Maximum mis-matches allowed	Reference for transcription factors and sequence motifs
109719	Fac219	ABREMOTIFIIOS RAB16B	GCCGCGTGCC	10	1.5	PLACE
109720	Fac220	ABREMOTIFIOSRAB16B	AGTACGTGGC	10	1.5	PLACE
109721	Fac223	CE3OSESEM	AACGCGTGTC	10	1.5	PLACE
109722	Fac267	POLASIG2	AATTAATA	7	0	PLACE
109723	OS_A-box	OS_A-box	TATCCATCCATCC	13	3	PlantCARE
109724	OS_A-box2	OS_A-box2	AATAACAAACTCC	13	3	PlantCARE
109725	OS_AACA	OS_AACA	TAACAAACTCCA	12	2.5	PlantCARE
109726	OS_ABRE	OS_ABRE	GACACGTACGT	11	2	PlantCARE
109727	OS_ABRE2	OS_ABRE2	ACGTACGTGTCGCGC	15	4	PlantCARE
109728	OS_AP-2-like	OS_AP-2-like	CGCGCCGG	8	0.5	PlantCARE
109729	OS_AP-2-like2	OS_AP-2-like2	CGACCAGG	8	0.5	PlantCARE
109730	OS_ATGCAAAT	OS_ATGCAAAT	ATACAAAT	8	0.5	PlantCARE
109731	OS_CE3	OS_CE3	GACGCGTGTC	10	1.5	PlantCARE
109732	OS_GATT	OS_GATT	CT CCTGATTGGA	12	2.5	PlantCARE
109733	OS_GCN4	OS_GCN4	TGWGTCA	7	0	PlantCARE
109734	OS_GCN4_2	OS_GCN4_2	CAAGCCA	7	0	PlantCARE
109735	OS_P-box	OS_P-box	CGCTTTTGAGT	11	2	PlantCARE
109736	OS_P-box2	OS_P-box2	CCTTTTG	7	0	PlantCARE
109737	OS_Prolamin_box	OS_Prolamin_box	TGCAAAGT	8	0.5	PlantCARE
109738	OS_Skn-1	OS_Skn-1	GTCAAT	5	0	PlantCARE
109739	OS_TATC-box	OS_TATC-box	TATCCA	7	0	PlantCARE
109740	OS_TGGCA	OS_TGGCA	GACACCAAGTGCCA	14	3.5	PlantCARE
109741	OS_light	OS_light	AACCAATCTCATCCATCC	18	5.5	PlantCARE
109742	AS_RF2A_01	AS_RF2A_01	CCAGTGTGGCGCTGG	15	4	TRANS
109743	AT_RS1A_01	AT_RS1A_01	CTTCCACGTGGCA	13	3	TRANS
109744	PV_GRP18_01	PV_GRP18_01	TGGATGTGGAAAGACAGCA	18	5.5	TRANS
109745	RICE_ACT_01	RICE_ACT_01	GCCCAACCCAACCCAAC	17	5	TRANS
109746	RICE_AGB_01	RICE_AGB_01	GCCACGTAAG	10	1.5	TRANS
109747	RICE_AGB_03	RICE_AGB_03	GCCACGTCAG	10	1.5	TRANS
109748	RICE_EM_01	RICE_EM_01	TACGTGT	7	0	TRANS
109749	RICE_EM_02	RICE_EM_02	GACGTGT	7	0	TRANS
109750	RICE_GL51_01	RICE_GL51_01	AAGTCATAACTG	12	2.5	TRANS
109751	RICE_GL51_02	RICE_GL51_02	CCATGTCATATT	12	2.5	TRANS
109752	RICE_GL51_03	RICE_GL51_03	AATGATGTGTCAAT	14	3.5	TRANS
109753	RICE_GL51_04	RICE_GL51_04	TCCCGTGTACCAC	13	3	TRANS
109754	RICE_GL51_05	RICE_GL51_05	TGAGTCA	7	0	TRANS
109755	RICE_GLU2_01	RICE_GLU2_01	CCTTTCGTGTACC	13	3	TRANS
109756	RICE_GLUB1_01	RICE_GLUB1_01	CTGAGTCAT	9	1	TRANS
109757	RICE_NITR_01	RICE_NITR_01	CACGTAC	8	0.5	TRANS
109758	RICE_RAB16A_01	RICE_RAB16A_01	TACGTGGCNNNNCCGC	23	6	TRANS
109759	RICE_RAB16A_03	RICE_RAB16A_03	CGCGCCT	8	0.5	TRANS
109760	TAF-1AS_	TAF1_01	GCAACGTGGC	10	1.5	TRANS
109761	TAF-1RICE_	RAB16B_01	GGTACGTGGCG	11	2	TRANS
109762	WHEAT_H3_01	WHEAT_H3_01	CCACGTCA	8	0.5	TRANS
109763	Seed_sp	Seed_AACA_motif	AACAACTCTATC	13	3	lit1
109764	Seed_sp	Seed_GCN4	GTGAGTCAC	9	1	lit1
109765	SugRep	ACGTABOX	TACGTA	6	0	lit2
109766	SugRep	TCmotif	TATCCAY	7	0	lit2
109767	Amy3	DAMYBOX2	TATCCAT	7	0	lit3
109768	Amy3	DGBOXRELOSAMY3	CTACGTGGCCA	11	2	lit3

\*Column Headings for Table 2

Transcription Factor Name: Name of the transcription factor which binds to a sequence motif within a promoter region

Sequence Motif name: Name of the sequence motif to which the transcription factor binds

Sequence Motif: sequences searched to further annotate putative promoters

Maximum mismatches allowed: Number of mismatches that are permitted when searching for motifs to annotate putative promoters

Reference for transcription factors and sequence motifs:

Motifs and transcription factors are found in one of three databases: PLACE, PlantCARE or TRANS (respectively, [www-dna.affrc.go.jp/htdocs/PLACE/](http://www-dna.affrc.go.jp/htdocs/PLACE/), [www-sphinx.rug.ac.be:8080/PlantCARE/index.htm](http://www-sphinx.rug.ac.be:8080/PlantCARE/index.htm), [www-transfac.gbf.de/TRANSFAC/](http://www-transfac.gbf.de/TRANSFAC/), or Yoshihara et al., FEBS Letters 383, 1996, pp 213-218; or Toyofuku K et al.

FEBS Lett 428:275-280 (1998) or lit3 (Huang et al Plant Mol Biol 14:655-668 (1990)).

Table 3

Table 3 describes those putative promoter sequences containing TATA boxes, GC boxes or CCAT boxes as determined by matrix motifs.

\*Column headings for Table 3

Seq num

Provides the SEQ ID NO. for the listed sequences.

Seq ID

Arbitrarily assigned identifier for each putative promoter sequence

Start

Indicates the start position of the TATA box, GC box or CCAAT box.

End

Indicates the end position of the of the TATA box, GC box or CCAAT box.

p-Value  
Probability value is determined by simulation as described above.

-ln (p-Value)  
Indicates the negative natural log of the p-Value.

# of hits in cluster  
No clustering is done in Table 3. Therefore, all entries in this column are listed as "1".

Factor Name  
Transcription factors associated with the TATA box (TATA-plant), GC box (V\_GC\_01) and CCAAT box (F\_HAP234\_01) are listed under this column heading when the matrix motifs for the TATA box, GC box and CCAAT box are identified.

Site name  
List whether the search is done for the TATA box, GC box or CCAAT box.

Table 4 describes those putative promoter sequences containing specific sequence motifs as listed in Table 2.

Column headings for Table 4

Seq num  
Provides the SEQ ID NO. for the listed sequences.

Seq ID  
Arbitrarily assigned identifier for each putative promoter sequence

Start  
Indicates the start position of a sequence motif. Note that if multiple motifs are present in a window of sequence, the start position, the first time it is indicated, lists the start position of the first motif in a window. The second or subsequent times a start position is listed, this heading refers to the start position of the subsequent individual sequence motif within the window of sequence.

End  
Indicates the end position of a sequence motif. Note that if multiple motifs are present in a window of sequence, the end position, the first time it is indicated, lists the end position of the first motif in a window. The second or subsequent times an end position is listed, this heading refers to the end position of the subsequent individual sequence motif within the window.

Strand  
The strand of genomic DNA on which the sequence motif is located (+/-)

p-Value  
Probability value as determined by simulation as described above. Note that if multiple motifs are present in a window of sequence, the p-Value, the first time it is indicated, lists the p-Value position for all of the motifs in a window of sequence. The second or subsequent times a p-Value is listed, this heading refers to the p-Value of the subsequent individual sequence motif within the window.

-ln (p-Value)  
Indicates the negative natural log of the p-Value as described in the column heading above.

# of hits in cluster P1  
Clustering is described above. The number of hits in a cluster refer to the number of times sequence motifs appear in a window of sequence.

Factor Name  
Transcription factors associated with the sequence motifs listed in Table 2 are included under this column heading.

Site name  
Lists the sequence motif for which a search is done on a putative promoter sequence.

Table 5  
Table 5 lists the putative promoter sequences for which sequence or matrix motifs are not identified.

Column headings for Table 5

Seq num  
Provides the SEQ ID NO. for the listed sequences.

Seq ID  
Arbitrarily assigned identifier for each putative promoter sequence

Table 6  
Table 6 lists the contigs, combigenes and description information associated with each gene prediction.

\*Column Headings:

Seq num  
Provides the SEQ ID NO. for the listed sequences.

Contig id  
Arbitrarily assigned name for each contig.

CDS.  
The location of the exons found within the gene as determined by the gene-predicting program (Method).

CG ID  
Arbitrarily assigned name for each combigene.

CG Start  
Indicates the start position of the combigene gene.

CG End  
Indicates the end position of the combigene gene.

Strand  
Indicates the strand location of the gene (+/-)

Gene  
Indicates an arbitrarily assigned gene name based on the method used to predict the gene.

Method  
Indicates the gene-predicting program used. These programs are GenScan, AAT/NAP, AAT/GAP, TBLASTX or Genemark.hmm.

Gene Start  
The start position of the putative gene making up a combigene as predicted by the particular gene predicting program used.

Gene End  
The end position of the putative gene making up a combigene as predicted by the particular gene-predicting program used.

Hit Score  
The aat\_nap score (under Hit score in the rows where the method is AAT/NAP) is reported by the nap program in the aat package. It is an alignment score in which each match and mismatch is scored based on the BLOSUM62 scoring matrix. The aat\_gap score (under Hit score in the rows where the method is AAT/GAP) is the alignment score for each hit sequence, as reported by AAT/GAP.

For TBLASTX the Bit score for BLAST match score that is generated by the sequence comparison of the genomic contig with the Monsanto cDNA sequence named under the GI column is listed. The E-value corresponding to a given bit score is  $E = mn^2 \cdot S'$ . "m" and "n" are two proteins of length "m" and "n", "E" is the E value and S' is the bit score.

GI  
Each sequence in the GenBank public database is arbitrarily assigned a unique NCBI gi (National Center for Biotechnology Information GenBank Identifier) number. In this table, the NCBI gi number which is associated (in the same row) with a given contig or singleton refers to the particular GenBank sequence which is the best match for that sequence. If the hit is based on cDNAs from Monsanto's SeqDB, the name, of the cDNA sequence it hit to is named.

## Description

The Description column provides a description of the NCBI gi referenced in the "GI" column.

6. The nucleic acid molecule of claim 3, wherein the nucleic acid sequence exhibits a 99% or greater identity to the nucleic acid sequence of SEQ ID NO:1.

## SEQUENCE LISTING

The patent contains a lengthy "Sequence Listing" section. A copy of the "Sequence Listing" is available in electronic form from the USPTO web site ([http://seqdata.uspto.gov/?pageRequest=docDetail&DocID=\[US07365185B2\]USRE046033E1](http://seqdata.uspto.gov/?pageRequest=docDetail&DocID=[US07365185B2]USRE046033E1)). An electronic copy of the "Sequence Listing" will also be available from the USPTO upon request and payment of the fee set forth in 37 CFR 1.19(b)(3).

## We claim:

1. A substantially purified nucleic acid molecule that comprises the nucleotide sequence of SEQ ID NO: 1, or a complement thereof, *operably linked to a heterologous structural nucleic acid sequence.* 20
2. A substantially purified nucleic acid molecule that consists of the nucleotide sequence of SEQ ID NO: 1, or a complement thereof, *operably linked to a heterologous structural nucleic acid sequence.* 25
3. A substantially purified nucleic acid molecule comprising a nucleic acid sequence wherein the nucleic acid sequence: 30
  - i) hybridizes under high stringency conditions with the sequence of SEQ ID NO:1 or a complement thereof; or
  - ii) exhibits an 85% or greater identity to the sequence of SEQ ID NO:1;*wherein the nucleic acid sequence is operably linked to a heterologous structural nucleic acid sequence.* 35
4. The nucleic acid molecule of claim 3, wherein the nucleic acid sequence exhibits a 90% or greater identity to the nucleic acid sequence of SEQ ID NO:1.
5. The nucleic acid molecule of claim 3, wherein the nucleic acid sequence exhibits a 95% or greater identity to the nucleic acid sequence of SEQ ID NO:1. 40
7. The nucleic acid molecule of claim 3, wherein said nucleic acid sequence comprises the sequence of SEQ ID NO: 1.
8. The nucleic acid sequence of claim 3, wherein the nucleic acid molecule further comprises one or more cis-acting nucleic acid elements.
9. The nucleic acid molecule of claim 3, wherein the nucleic acid molecule further comprises a 5' leader sequence selected from the group consisting of dSSU 5', PetHSP70 5', and GmHSP17.9 5'.
10. The nucleic acid molecule of claim 3, wherein the nucleic acid molecule further comprises a 3' untranslated region.
11. The nucleic acid molecule of claim 10, wherein the 3' untranslated region is selected from the group consisting of NOS 3', E9 3', ADR12 3', 7S $\alpha$ 3', 11S 3', and albumin 3'.
12. A transgenic plant comprising a recombinant nucleic acid molecule having the nucleic acid sequence of claim 3.
13. A host cell comprising a recombinant nucleic acid molecule having the nucleic acid molecule of claim 3.
14. The host cell of claim 13, wherein said host cell is a plant cell.
15. A transgenic plant comprising the host cell of claim 13.

\* \* \* \* \*