

Supplementary Information 1: a list of regular expression used in our “uninformative rule” filter.

HYPOTHETICAL PROTEIN
HYPOTHETICAL PROTEIN %S+
HYPOTHETICAL PROTEIN KIAA%d+
HYPOTHETICAL PROTEIN FLJ%d+
HYPOTHETICAL PROTEIN CGI%-d+
HYPOTHETICAL PROTEIN DKFZP%S+
HYPOTHETICAL PROTEIN, MGC:-d+
HYPOTHETICAL GENE %S+
HA%d{4}
HYPOTHETICAL [%.%d]+%s*KDA PROTEIN
HYPOTHETICAL [%.%d]+%s*KDA PROTEIN PRECURSOR
DKFZP%S+ PROTEIN
UNKNOWN
PROTEIN FOR MGC:-d+
PROTEIN FOR IMAGE:-d+
R%d{5}%_d
R%d{5}%_d, PARTIAL PROTEIN
R%d{5}%_d, PARTIAL CDS
PRO%d{4}
PRO%d{4} PROTEIN
ORF %d
PUTATIVE KIAA%d+ HOMOLOGUE
KIAA%d+ GENE PRODUCT
KIAA%d+
KIAA%d+ PROTEIN
HSPC%d+
HSPC%d+ PROTEIN
CHROMOSOME %d+ OPEN READING FRAME %d+
C%d+ORF%d+
FLJ%d+ PROTEIN
DJ%d+[A-Z]%d+(%.%d+)*
DJ%d+[A-Z]%d+(%.%d+)* PROTEIN
NOVEL PROTEIN

PUTATIVE NOVEL PROTEIN
CGYd+ PROTEIN
CGYd+ GENE PRODUCT
CGYd+
CGIY-Yd+ PROTEIN
CGIY-Yd+
CDNA:? FLJYd+ FIS, CLONE Yw+
BAYd+[A-Z]Yd+[A-Z]?Y.Yd(Y.Yd)?
RIKEN CDNA .{10} GENE
MRNA, COMPLETE CDS, CLONE:Yd+(Y+Yd[A-Z])?Y-Yd+
MRNA, COMPLETE CDS, CLONE:SMAPYd+Y-Yw+
BRAIN CDNA, CLONE MNCB-Yd+
. {10}RIK PROTEIN
UNNAMED PROTEIN PRODUCT
MYYd{3}
MYYd{3} PROTEIN
BRAIN MYYd{3}
NPDYd{3} PROTEIN
[A-Z][0-9]{2}[A-Z0-9]+Y.[0-9]+ PROTEIN
WUGSC:H_Yw+Y.Yw+ PROTEIN
expressed sequence YS+
putative
unknown
ORF
open reading frame Yd+
DNA SEGMENT, CHR [0-9XY]+, WAYNE STATE UNIVERSITY Yd+,
EXPRESSED
DNA SEGMENT, CHR [0-9XY]+, KL MOHLKE Yd+
DNA SEGMENT, CHR [0-9XY]+, BAYLOR Yd+
PROTEIN HSPCYd+
ORF PROTEIN
unknown protein, Yd+Y-Yd+
HYPOTHETICAL [Y.Yd]+Ys*KDA PROTEIN YS+ IN CHROMOSOME YS+
EG:[0-9A-ZY.] + PROTEIN
GENOMIC DNA, CHROMOSOME Yd+, P1 CLONE:YS+
[^,]+, RIKEN FULL-LENGTH ENRICHED LIBRARY, CLONE:.{10}, FULL

INSERT SEQUENCE

ZK¥d+¥.¥d+ PROTEIN

EST ¥w+

B2 ELEMENT

hypothetical protein ¥S+