

ESTs

1* 2 3 4 5
1,3,4,
1,3,4,5
1,2 가
kth2001@pusan.ac.kr

Biological Function Category of ESTs

Tae-Hyung Kim^{1*} Cheol-Goo Hur² Yeo-Jin Jeon³ Dae-Su Kim⁴ Heui-Soo Kim⁵

^{1,3,4} Interdisciplinary program of bioinformatics, Pusan National University

^{1,3,4,5} Dept of Biology, College of Natural Sciences, Pusan National University

^{1,2} Laboratory of Computational Biology, National Bioinformatics Center, Korea Research Institute of Bioscience and Biotechnology, Oun-dong 52, Yu-Sung, Daejeon, 305-333 Korea

ESTs
가
single linkage ESTs clustering BEC (Blast ESTs clustering)
ESTs virus small RNA
ESTs
Riken Fantom full-length mouse cDNA NCBI Refseq mouse
mRNA (Gene
Ontology) 가 hierarchical No-hit ESTs
3 Unknown
ESTs
clustering, quality ESTs 가

1.

ESTs [5,6], 가
가 [7] [8](virus like elements, rRNA like
elements) cDNA가
ESTs
cDNA 가 ESTs 가
ESTs
가 [1-4].
ESTs
(1/100bp) 가 ESTs

ESTs, formatdb, BLAST, UniGene, ESTs, dbEST Library, bone, 12,709 ESTs, BEC, bit score 200, ESTs, UniGene, bone, ESTs, CPU 700MHZ, 370MB, PC, 1, ESTs

1. ESTs

ESTs minimal linkage (transitive closure, single linkage)[10] BLAST (threshold) 가 ESTs cluster 가 cluster cluster

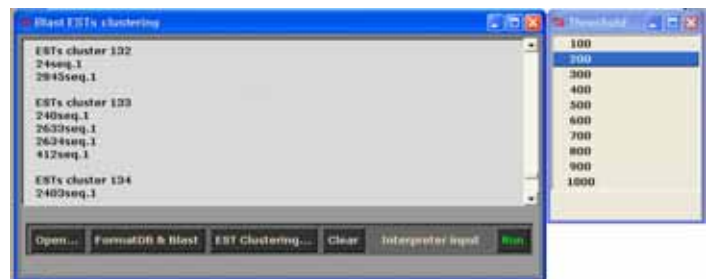


Figure 1. python BEC GUI

: cluster 가 (Si cluster i or Ci = i). : entry S0 query Si (1 ≤ i < N) MERGE (cluster C0) ← (cluster Ci) if d2 (S0, Si) < THRESHHOLD.

minimal linkage BEC mouse 3000 ESTs threshold (100~500) ESTs가 가 ESTs singleton ESTs

: S1 query가 cluster가 Si (2 ≤ i < N) MERGE (cluster C1) ← cluster Ci) if d2 (S1, Si) < THRESHHOLD.

(k) : (k - 1) . k- 1 Sk query cluster 가 (k + 1 ≤ i < N) MERGE (cluster Ck) ← cluster Ci) if d2 (Sk, Si) < THRESHHOLD.

Table 1. Score ESTs clustering

Threshold	Clusters	Total ESTs of clusters	Singleton ESTs
100	156	406	2,594
150	154	388	2,612
200	149	363	2,637
250	133	323	2,677
300	122	295	2,705
350	99	239	2,761
400	88	211	2,789
450	73	176	2,824

ESTs GUI ESTs가

500	65	157	2,843
-----	----	-----	-------

2.

rRNA 가 tRNA

(Cellular Gene)

ESTs cDNA

ESTs

[11].

ESTs

[12].

Repeatmasker[13]

Rebase[14]

Table 2

3,000

ESTs 628 elements가 SINEs (B1s, B2-B4, IDs, MIRs), LINEs (LINE1, LINE2), LTR elements (MaLRs, ERV_classI, ERV_classII)

DNA elements (MER1_type, MER2_type) DNA transposon, Small RNA, tandem repeat, Low complexity가

SINEs, LINEs, LTR elements repeat elements가 284 ESTs가

Table 2. 3,000 ESTs repeat

```

file name: mouse_3000
sequences: 3000
total length: 1337959 bp (1337446 bp excl N-runs)
GC level: 47.33 %
bases masked: 70809 bp ( 5.29 %)
=====
number of length percentage
elements* occupied of sequence
-----
SINEs:
  B1s 233 28074 bp 2.10
  B2-B4 116 13081 bp 0.98
  IDs 88 12481 bp 0.93
  MIRs 17 1294 bp 0.10
  12 328 bp 0.03
LINEs:
  LINE1 32 6573 bp 0.49
  LINE2 5 6126 bp 0.46
  L3/CR1 0 449 bp 0.03
LTR elements:
  MaLRs 64 15427 bp 1.15
  ERVL 24 4041 bp 0.30
  ERV1 0 0 bp 0.00
  ERV_classI 0 0 bp 0.00
  ERV_classII 25 2156 bp 0.16
  25 8386 bp 0.63
DNA elements:
  MER1_type 11 1471 bp 0.11
  MER2_type 8 344 bp 0.03
  2 362 bp 0.03
Unclassified: 2 634 bp 0.05 %
Total interspersed repeats: 52179 bp 3.90 %

Small RNA: 4 273 bp 0.02 %
Satellites: 1 73 bp 0.01 %
Simple repeats: 213 8908 bp 0.67 %
Low complexity: 220 9411 bp 0.70 %
=====

```

Retrotransposons ESTs 가

300bp ESTs

Table 3. 21 ESTs

L1, ERVK Retrotransposons ESTs

가

가

Table 3. ESTs 300bp

EST	Repeat	Family	Length	Direction
313_seq	MTA	LTR/MaLR	325	C
333_seq	IAPEz-int	LTR/ERVK	448	+
387_seq	Lx2	LINE/L1	398	+
404_seq	IAPEz-int	LTR/ERVK	480	+
1015_seq	L1	LINE/L1	638	+
1692_seq	Lx2	LINE/L1	355	C
1933_seq	ETnERV3	LTR/ERVK	830	+
1934_seq	L1	LINE/L1	396	C
1946_seq	Lx4	LINE/L1	438	C
2019_seq	ORR1A1	LTR/MaLR	346	+
2095_seq	IAPEz-int	LTR/ERVK	376	+
2100_seq	L1	LINE/L1	361	C
2134_seq	IAPEz-int	LTR/ERVK	878	+
2162_seq	L1	LINE/L1	420	+
2178_seq	IAP-d	LTR/ERVK	448	+
2414_seq	RMER17A2	LTR/ERVK	577	+
2429_seq	MTD	LTR/MaLR	409	C
2440_seq	RMER4B	LTR/ERVK	405	+
2509_seq	MMETn-int	LTR/ERVK	382	+
2670_seq	MER34B	LTR/ERV1	429	+
2964_seq	MMETn-int	LTR/ERVK	475	+

3. ESTs

ESTs RIKEN Fantom cDNA NCBI

Refseq mRNA (NM_XXXXXX)

ESTs

MGI annotation data

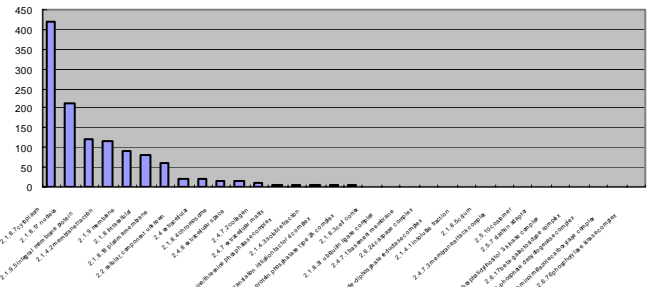
가

GO

Gene Ontology 3가 GO (biological process, cellular component, molecular function)

Figure 2.

3



3가

가

Biological

process 1.4.12.6 (Biosynthesis), 1.4.12.21 (nucleic acid metabolism), 1.4.12.9 (catabolism)

cellular component 2.1.8.7 (cytoplasm), 2.1.8.17 (nucleus), 2.1.9.5 (integral membrane protein),

Molecular function 3.12.25.41

(protein kinase), 3.6.18.3 (purine nucleotide binding),

3.6.17.1 (DNA binding) ESTs가

가

가

GO

MGI[15]

GO

GO node

biological

process 1, cellular component 2, molecular

function 3, 3가

1.xx.xx....

GO

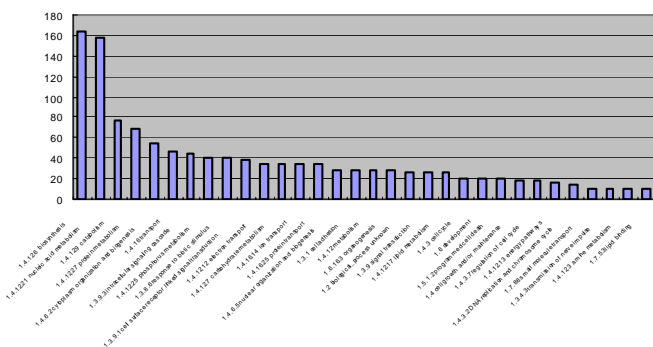
ESTs

가

3,000 ESTs

가

A. Biological process



B. Cellular component

C. Molecular function

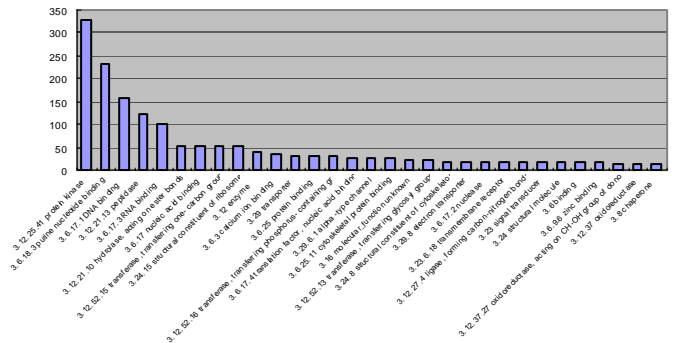


Figure 2. 3가 GO(Gene Ontology) category

ESTs

GO

No-hit

ESTs

BLASTX

ESTs

NR

(Non-redundant)

100bp

90% identity 가 ESTs

가

90%~70% identity 가 ESTs

70%~40% identity 가 ESTs

identity

3

ESTs

Interproscan

[16]

Profile scan

Fingerprint

scan

ESTs

가

가

가 가

Alternative splicing

3,000 ESTs

가

Figure 3.

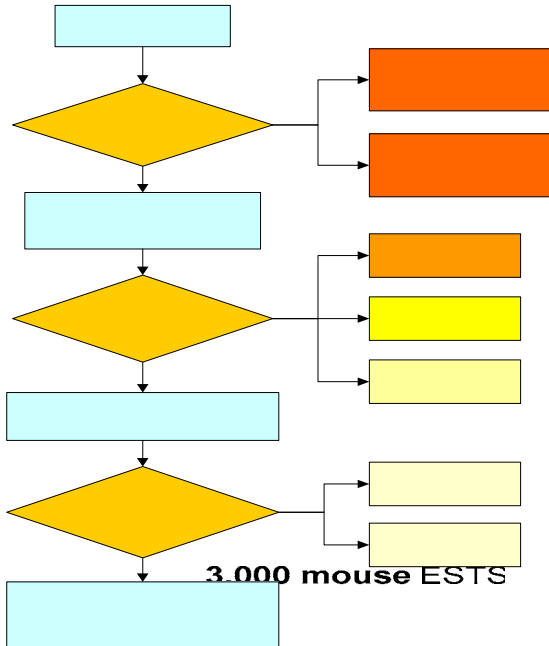


Figure 3. ESTs

BLASTN

Identifying known Gene
(Fantom, Refseq)

3,000

Did not confirmed to gene

No hit ESTs
1,219

minimum linkage ESTs

BLASTX

ESTs matching against
with NR proteins
(>100bp)

ESTs
Did not assigned to NR proteins
No hit ESTs

828

ESTs

Interproscan

Searching
motifs & domains
within ESTs

Did not categorized to several

ESTs

가 가

가

GO

ESTs

3가

가

ESTs

GO

가

가

가

RIKEN fantom cDNA

Identity >= 97%
e-value <= 1e-50
1,392 hits

가

NCBI Refseq mRNA

Identity >=95%
e-value <= 1e-50
1,318 hits

ESTs

1. Adams,M.D., Dubnick,M., Kerlavage,A.R., Moreno,R., Kelley,J.M., Utterback,T.R., Nagle,J.W., Fields,C. and Venter,J.C. Sequence identification of 2,375 human brain genes. *Nature*, **355**, 632-634. 1992.

2. Adams,M.D., Kerlavage,A.R., Fields,C. and Venter,J.C. 3,400 new expressed sequence tags identify diverse transcripts in human brain. *Nature Genet.*, **4**, 256-267. 1993.

3. Boguski,M.S. The turning point in genome research. *Trends Biochem. Sci.*, **20**, 295-296. 1995.

FingerprintsScan
452 hits

4. Marra,M.A., Hillier,L. and Waterston,R.H. Expressed sequence tags—ESTablishing bridges between genomes. *Trends Genet.*, **14**, 4–7. 1998.
5. Ewing,B. and Green,P. Base-calling of automated sequencer traces using Phred. II. Error probabilities. *Genome Res.*, **8**, 186–194. 1998.
6. Hillier,L.D., Lennon,G., Becker,M., Bonaldo,M.F., Chiapelli,B., Chissoe,S., Dietrich,N., DuBuque,T., Favello,A., Gish,W. *et al.* Generation and analysis of 280,000 human expressed sequence tags. *Genome Res.*, **6**, 807–828. 1996.
7. Croft,L., Schandorff,S., Clark,F., Burrage,K., Arctander,P. and Mattick,J.S. ISIS, the intron information system, reveals the high frequency of alternative splicing in the human genome. *Nature Genet.*, **24**, 340–341. 2000.
8. Weber,G., Shendure,J., Tanenbaum,D.M., Church,G.M. and Meyerson,M Identification of foreign gene sequences by transcript filtering against the human genome. *Nature Genet.*, **30**, 141–142. 2002.
9. Ashburner,M., Ball,C.A., Blake,J.A., Botstein,D., Butler,H., Cherry,J.M., Davis,A.P., Dolinski,K., Dwight,S.S., Eppig,J.T., Harris,M.A., Hill,D.P., Issel-Tarver,L., Kasarskis,A., Lewis,S., Matese,J.C., Richardson,J.E., Ringwald,M., Rubin,G.M. and Sherlock,G. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25** 25-29. 2000.
10. Miller,R.T., Christoffels,A.G., Gopalakrishnan,C., Burke,J., Ptitsyn,A.A., Broveak,T.R. and Hide,W.A. A comprehensive approach to clustering of expressed human gene sequence: the sequence tag alignment and consensus knowledge base. *Genome Res.*, **9**, 1143-1155. 1999.
11. Rotem Sorek, Gil Ast and Dan Graur. Alu-containing exons are alternatively spliced. *Genome Res.* **12**, 1060-1067. 2002.
12. Kim,H.S., Crow,T.J. and Hyun,B.H. Assignment of the endogenous retrovirus HERV-R (ERV3) to human chromosome 7q11.2 by radiation hybrid mapping. *Cytogenet. Cell Genet.* **89**, 10. 2000.
13. Smit, A. and P. Green. RepeatMasker at <http://ftp.genome.washington.edu/RM/RepeatMasker.html>. 1991.
14. Jurka,J. Repbase update: a database and an electronic journal of repetitive elements. *Trends Genet.* **16**, 418-420. 2000.
15. Jackson Laboratory. GO association data at ftp://ftp.informatics.jax.org/pub/reports/gene_association.mgi. 2002.
16. Zdobnov,E.M., Lopez,R., Apweiler,R. and Etzold,T. The EBI SRS server - recent developments. *Bioinformatics*, **18** 368-373. 2002.