

TRANSCRIPTOME 2002

Seattle, WA

**The Integrated Molecular Analysis of
Genomes and their Expression
Consortium's Data Mining Tools:
Introducing the IQ**

Peg Folta

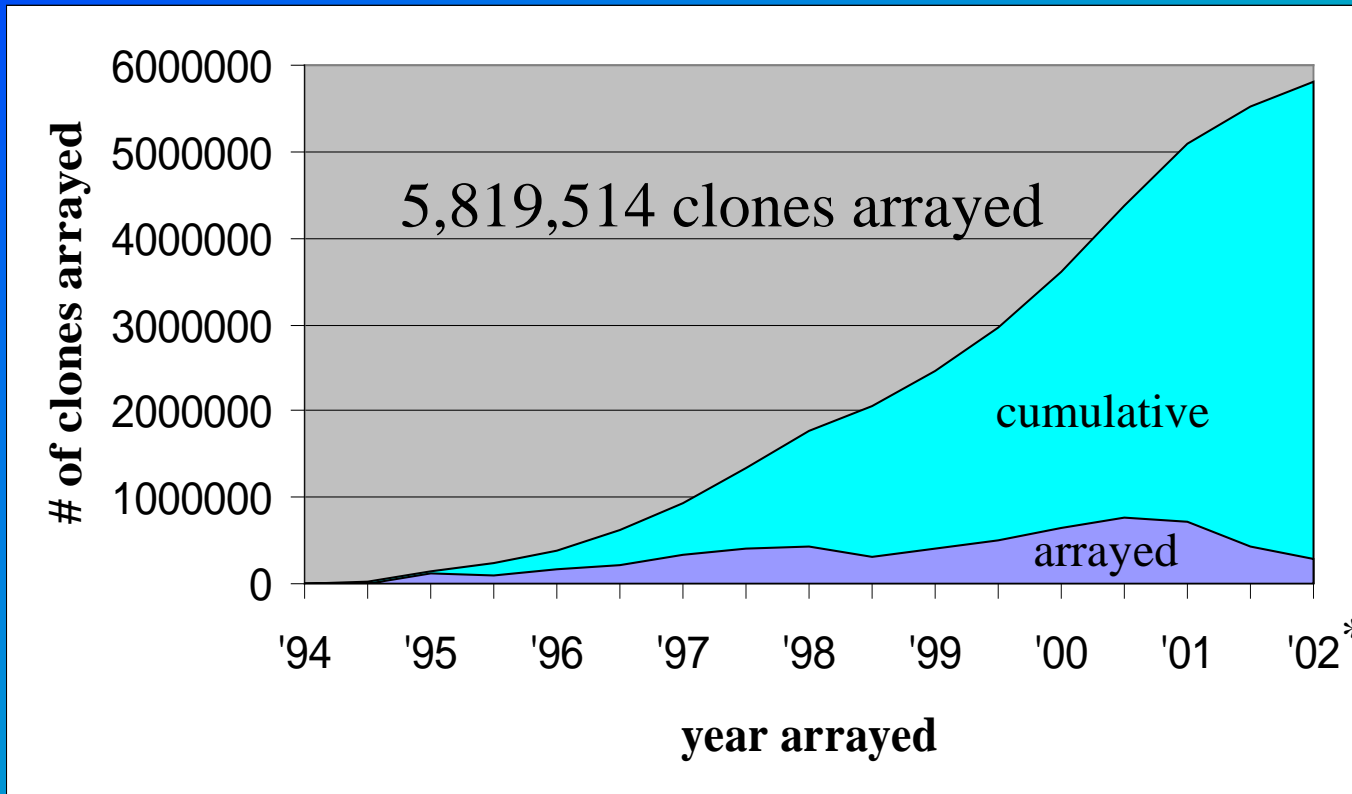
Lawrence Livermore National Laboratory

3/12/02





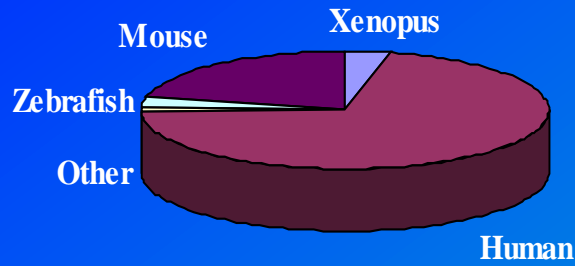
I.M.A.G.E. maintains world's largest publicly available cDNA collection



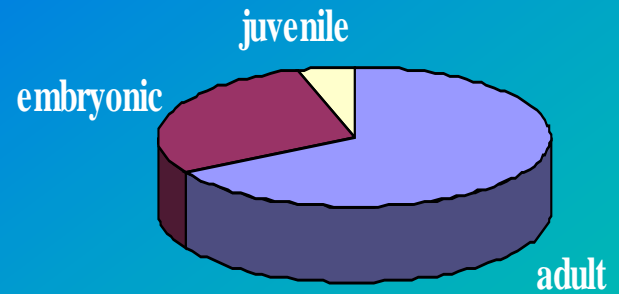
I.M.A.G.E. clones account for 64% of human ESTs in GenBank



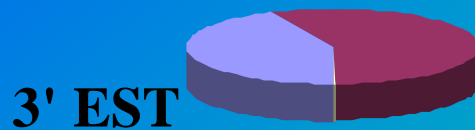
The I.M.A.G.E. collection has been shaped by projects (C-GAP, MGC...)



Species

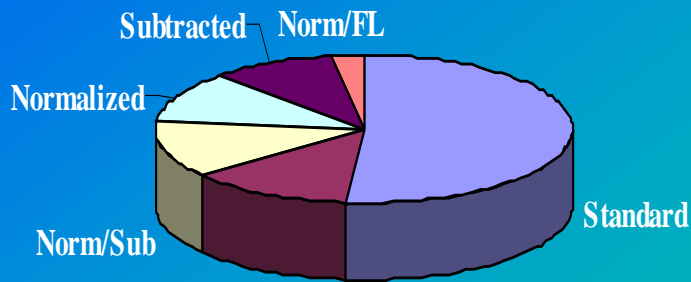


5' EST Developmental state

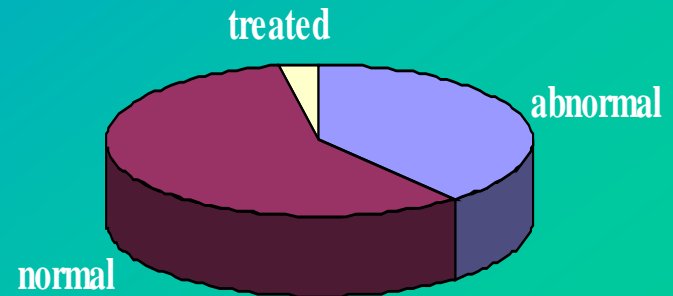


Full length

Clone sequence



Library Method



Tissue



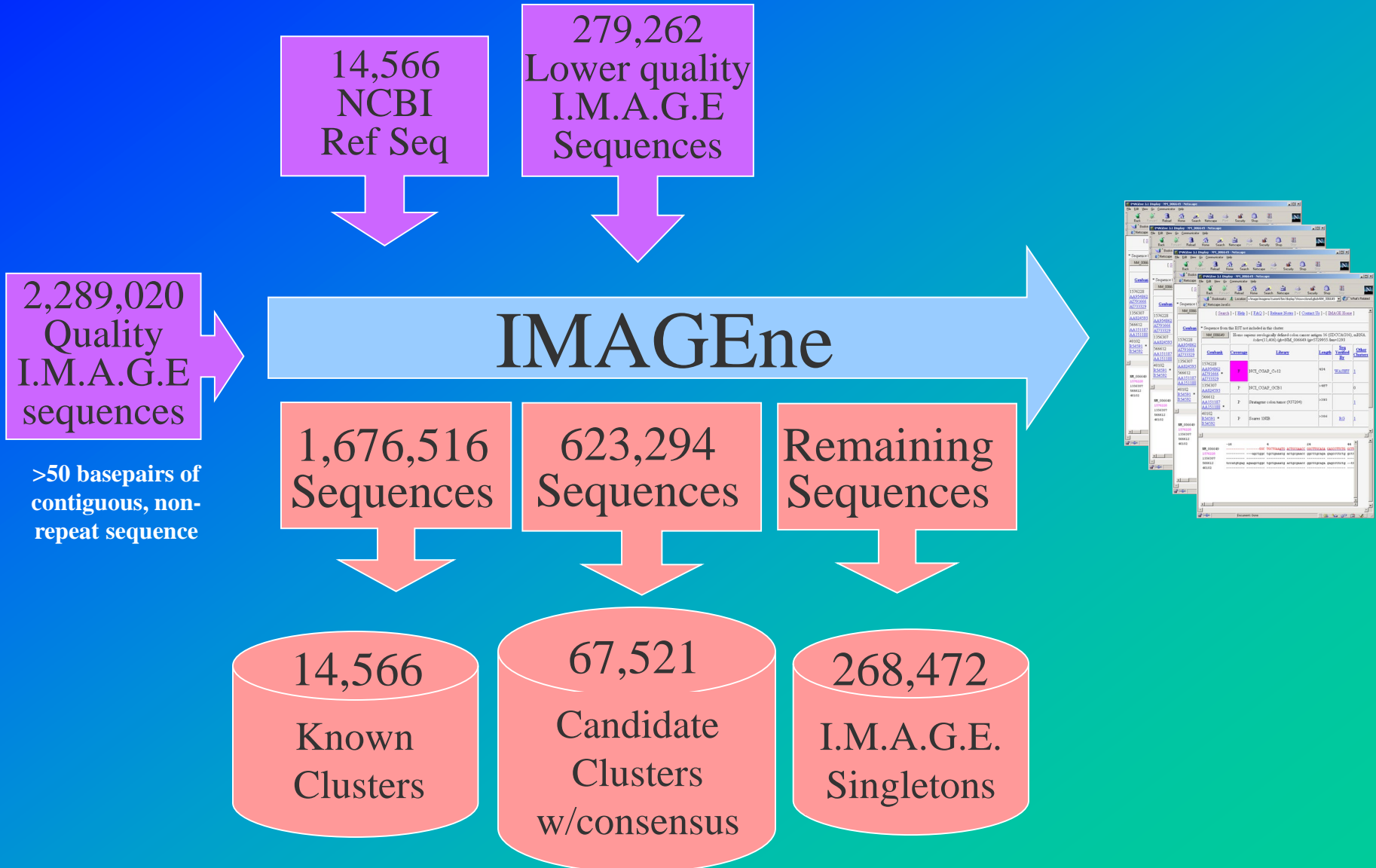
Informatics focus this year was on tools to characterize and query the collection.

- IMAGEne – mature clustering tool
- IMAGEne Tissue – allows searching of tissue type dominance in clusters
- IQ – Intelligent Query tool allows mining of I.M.A.G.E. data
- Library/plate query – allows selective searching of libraries and plates
- Problem report and query – allows users to report or query problems related to I.M.A.G.E. clones

Redesign of data management system



IMAGEnE-Human Process





Initial query page, construct the query.

[[Search](#)] - [[Help](#)] - [[FAQ](#)] - [[Release Notes](#)] - [[Contact Us](#)] - [[IMAGE Home](#)]

All sequence from I.M.A.G.E. human cDNA clones have been used to create two types of gene clusters: Known Gene Clusters based on NCBI's Reference Genes, and Candidate Gene Clusters having no known gene association. TIGR clones are used to supplement the Candidate Gene clustering. Those clones whose sequence does not match any other cDNA are grouped as Singletons. Submit a query against either cluster set to obtain a ranked display of available I.M.A.G.E. clones aligned with the corresponding known human gene or consensus sequence. *Alternative methods to query against the Singleton database will be provided at a later date.*

IMAGEne Release 3.3 (data retrieved from Genbank 1-8-02)

14566 Known gene clusters					75805 Multi-Member Candidate gene clusters	284749 ESTs from 268472 Singleton Clones (no query yet)
8834 Full	1493 Predicted Full	348 Unknown	3252 Partial	639 Empty		

Search by:

Query:

- Gene Name/Keyword
- IMAGE Clone ID
- GB Accession
- Sequence
- Cluster ID

Minimum Blast2 Score

Show Alignments by

Document: Done



Clusters matching query results, chose your cluster.

[[Search](#)] - [[Help](#)] - [[FAQ](#)] - [[Release Notes](#)] - [[Contact Us](#)] - [[IMAGE Home](#)]

Matches: 63

cluster id	description	fulls	predicted fulls	unknowns	partials
NM_000059	Homo sapiens breast cancer 2, early onset (BRCA2), mRNA. /cds=(229,10485) /gb=Nm_000059 /len=10987	0	0	5	19
NM_000249	Homo sapiens mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2) (MLH1) mRNA. /cds=(22,2292) /gb=Nm_000249 /len=2484	2	1	46	47
NM_000251	Homo sapiens mutS (E. coli) homolog 2 (colon cancer, nonpolyposis type 1) (MSH2) mRNA. /cds=(69,2873) /gb=Nm_000251 /len=3145	0	0	45	28
NM_000314	Homo sapiens phosphatase and tensin homolog (mutated in multiple advanced cancers 1) (PTEN) mRNA, and translated products. /cds=(1035,2246) /gb=Nm_000314 /len=3160	2	0	42	72
NM_001327	Homo sapiens cancer/testis antigen (CTAG1), mRNA. /cds=(89,631) /gb=Nm_001327 /len=806	6	0	4	6
NM_001635	Homo sapiens amphiphysin (Stiff-Mann syndrome with breast cancer 128kD autoantigen) (AMPH), mRNA. /cds=(75,2162) /gb=Nm_001635 /len=3260	0	1	41	19
NM_002387	Homo sapiens mutated in colorectal cancers (MCC) mRNA. /cds=(221,2710) /gb=Nm_002387 /len=4181	1	1	1	7
NM_003087	Homo sapiens synuclein, gamma (breast cancer-specific protein 1) (SNCG), mRNA. /cds=(49,432) /gb=Nm_003087 /len=720	13	0	21	8
NM_003225	Homo sapiens trefoil factor 1 (breast cancer, estrogen-inducible sequence expressed in) (TFF1) mRNA. /cds=(41,295) /gb=Nm_003225 /len=540	9	0	18	1
NM_003567	Homo sapiens breast cancer anti-estrogen resistance 3 (BCAR3) mRNA. /cds=(61,2538) /gb=Nm_003567 /len=3004	0	0	29	33
NM_003627	Homo sapiens prostate cancer overexpressed gene 1 (POV1) mRNA, and translated products. /cds=(77,1756) /gb=Nm_003627 /len=2326	1	1	34	20
	Homo sapiens Deleted in oral cancer-1 (DOC1) mRNA. /cds=(523,870)				

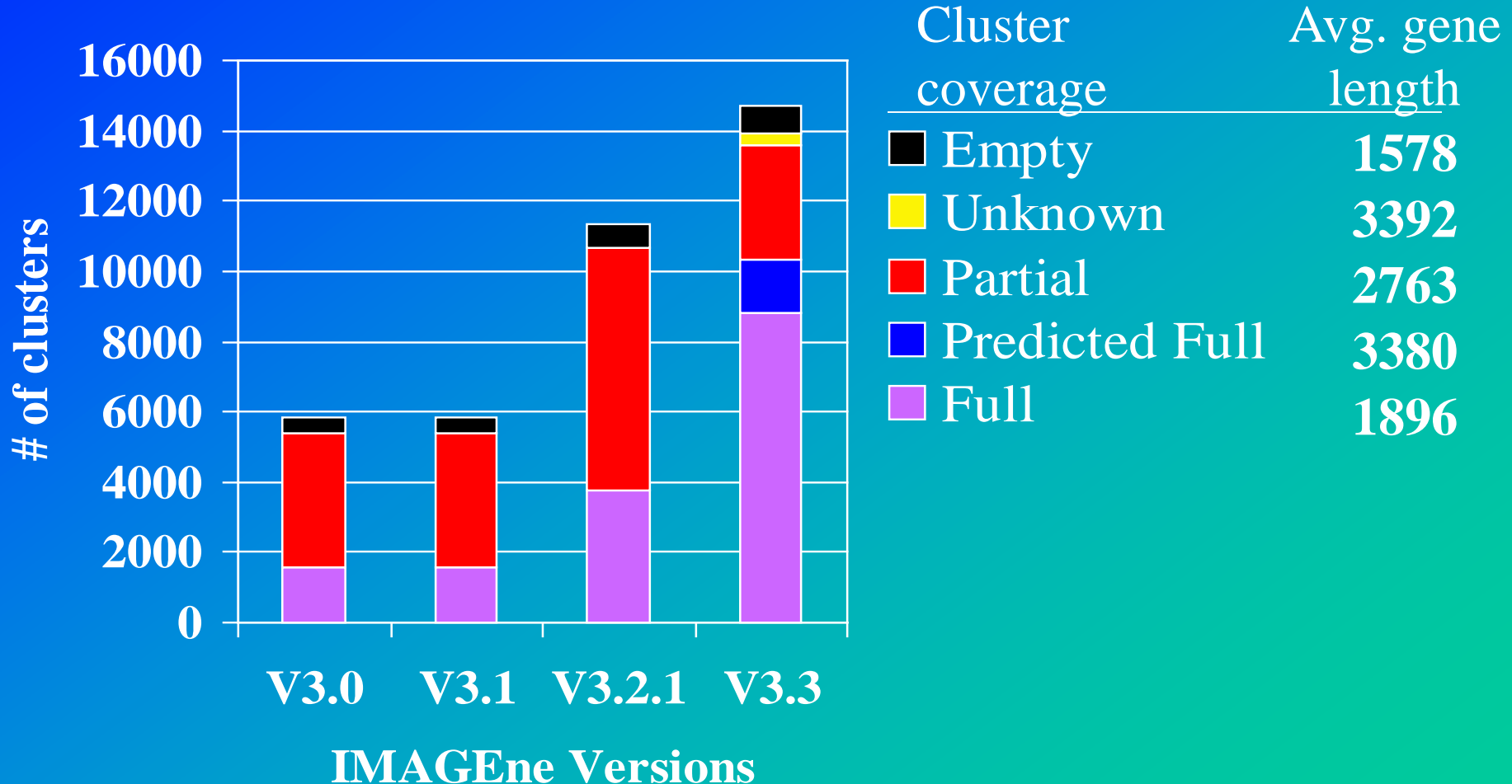


Document: Done



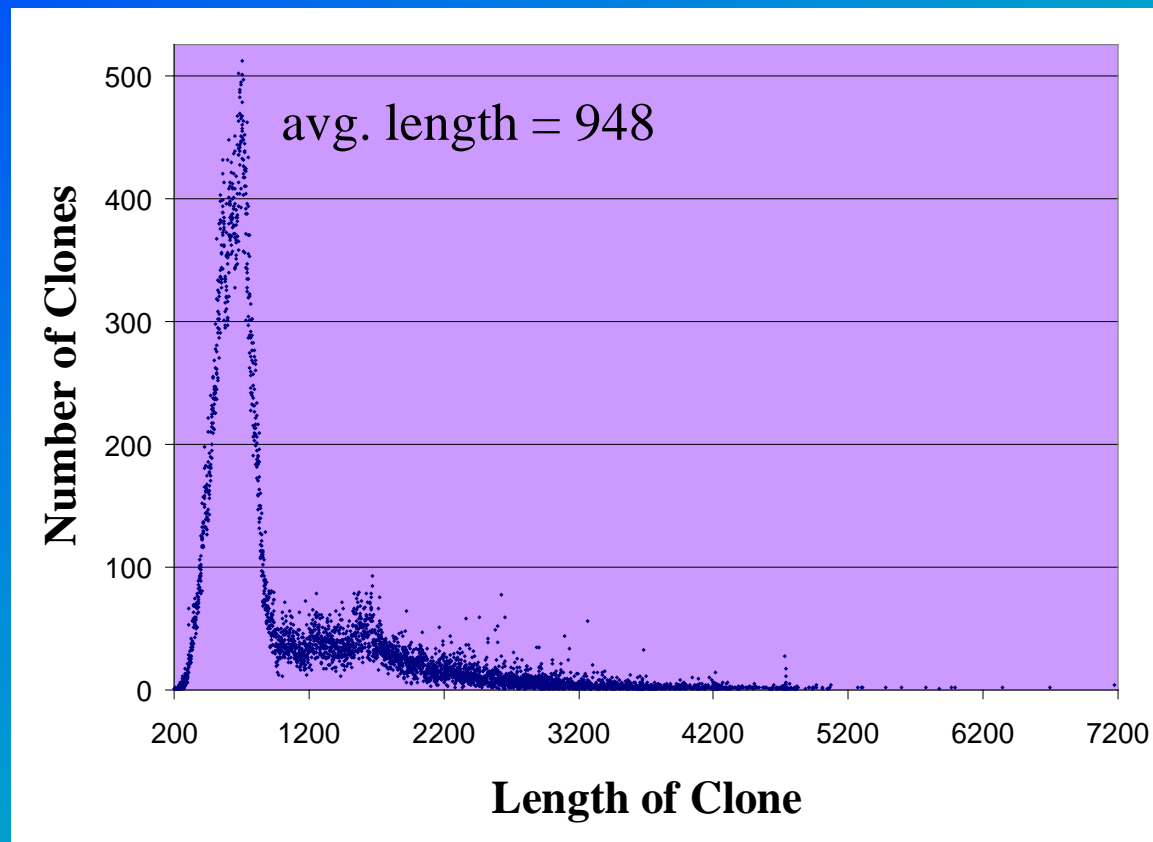


Known gene clusters with full length I.M.A.G.E. clones have doubled in number.



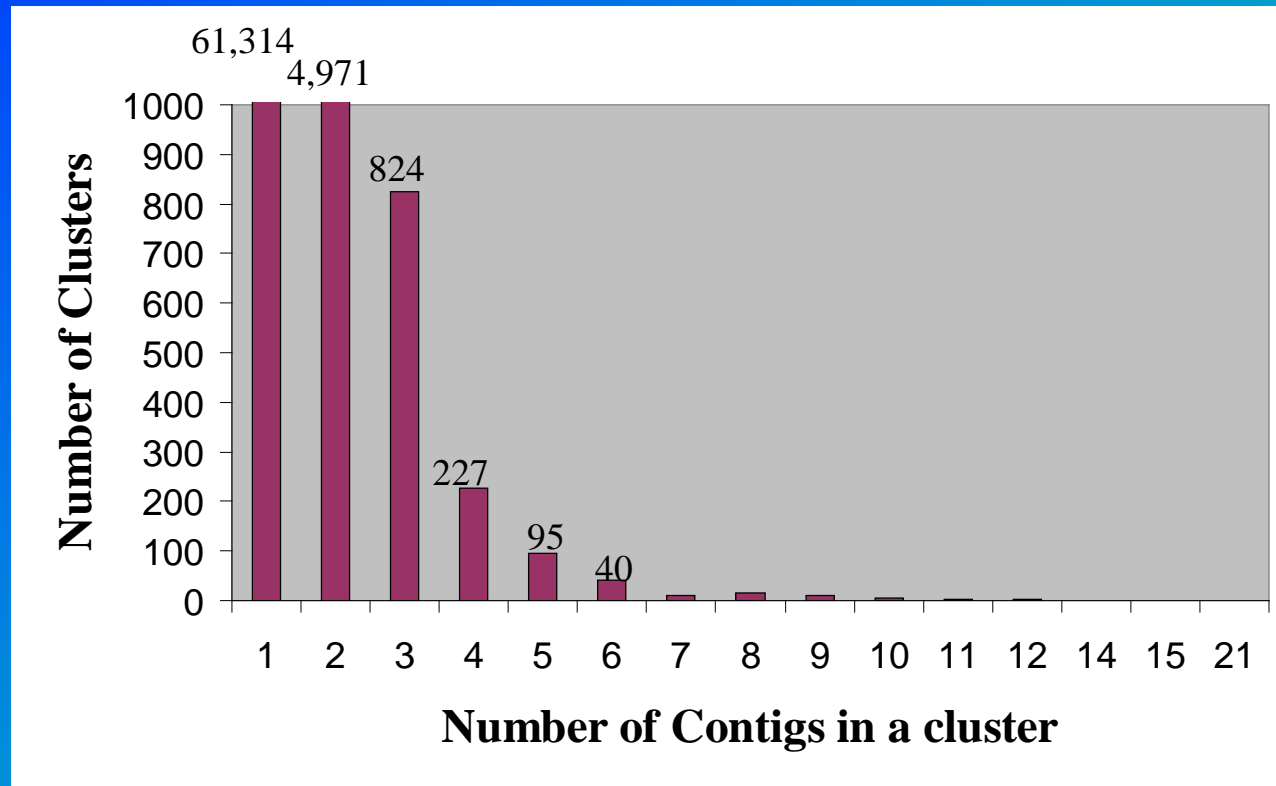


Known Gene Cluster distribution of full length clones



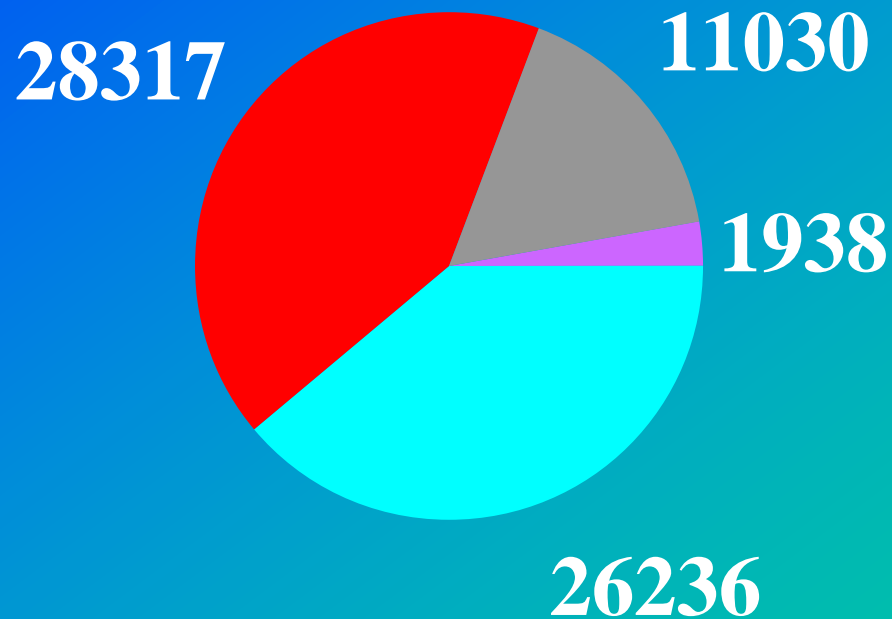


Candidate gene clusters consensus sequence and contigs are generated by CAP4





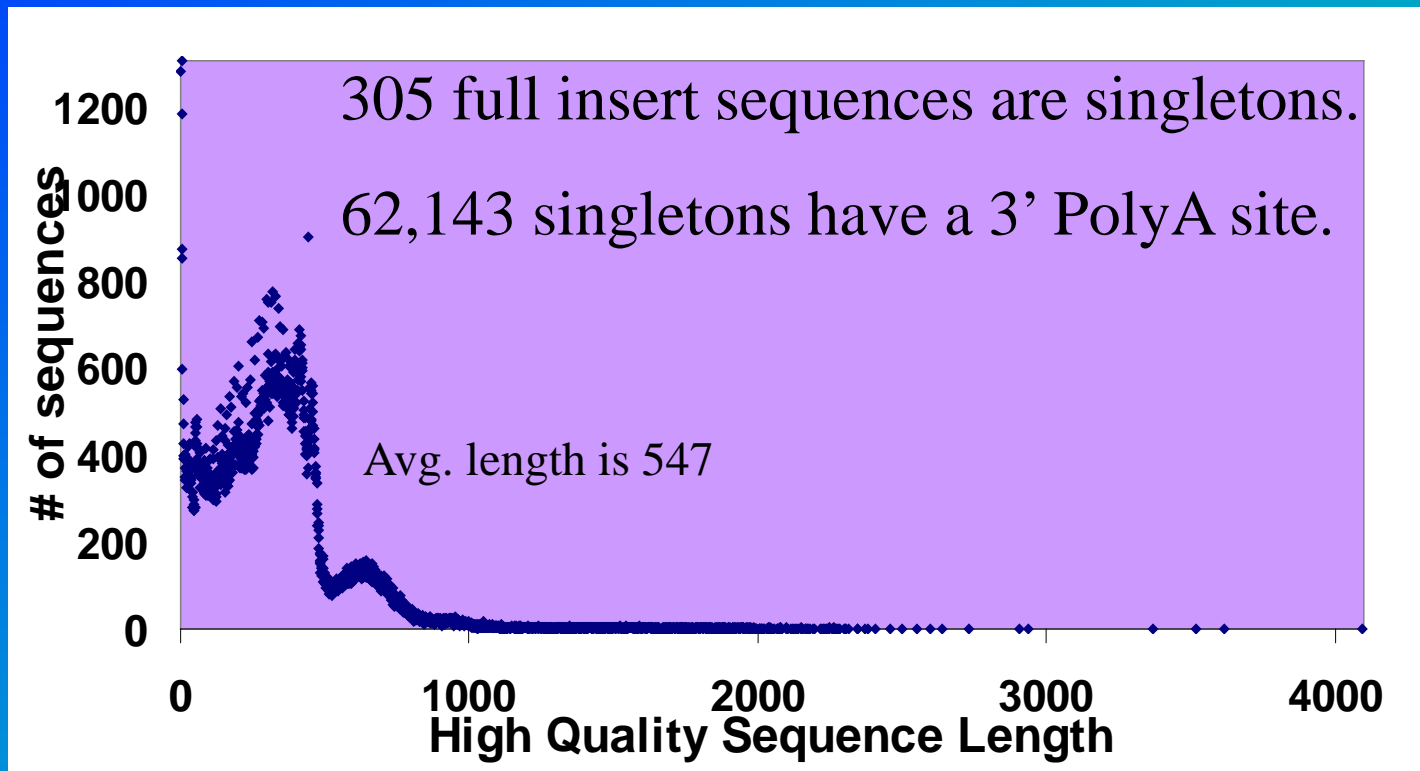
Candidate Gene cluster characteristics.



■ full insert ■ 3' & 5' ■ 3' only ■ 5' only



Singleton: Wheat within the chaff





IMAGeNe Tissue query allows searching for tissue proportions within clusters.

Please select a result format:

View results in an html page

Download results into a plain text file

Select desired file format:

[Imagene Version:](#)

Please choose a [cutoff](#) or [cutoff range](#):

Minimum: to Maximum

Tissues in Imagene

[Help](#)

<u>Contig ID</u>	<u>GB AccNum</u>	<u>Clone ID</u>	<u>Tissue</u>	<u>Ratio</u>
C000053-02	R42986	32492	brain/CNS	1.0000
C000053-02	R44600	33222	brain/CNS	1.0000
C000077-01	AA223686	651002	brain/CNS	1.0000
C000077-01	AA223694	651003	brain/CNS	1.0000
C000246-04	AI498073	2163041	brain/CNS	1.0000
C000246-04	BI831355	5166226	brain/CNS	1.0000
C000354-02	R15298	29768	brain/CNS	1.0000
C000354-02	T74017	22595	brain/CNS	1.0000
C000522-02	BI915961	5241183	brain/CNS	1.0000
C000522-02	BI457868	5277635	brain/CNS	1.0000
C000581-03	BG910399	4938321	brain/CNS	1.0000



Introducing the Intelligent Query - IQ

- For a given category (currently clone and library) a user can specify a query based on key database attributes.
- The user can specify the fields returned.
- Various result format options (HTML, text)
- Initial version was rolled out last summer
- New functionality to be added this year (additional categories, etc.)

Specify a clone-based query.

Field Name	Data Input
Clone ID	<input type="text"/> Convert
Plate Number	<input type="text"/> Convert
I.M.A.G.E. Library ID	<input type="text"/> Convert
I.M.A.G.E. Library Name	<div style="border: 1px solid gray; padding: 2px;"><p>Appel/Eisen 15-19 hr embryo Baker mouse embryo e6.5 Baker mouse embryo e7.5 Barstead HPL-RB1 Barstead HPL-RB2 Barstead HPL-RB3 Barstead HPL-RB4 Barstead HPL-RB4.1 Barstead HPL-RB5 Barstead HPL-RB6</p></div> <input type="button" value="Clear List"/>
Row Position	<input type="text" value="none selected"/>
Column Position	<input type="text" value="none selected"/>
Genbank Accession Number	<input type="text"/> Convert
Genbank Identifier	<input type="text"/> Convert
Tissue Type	<input type="text" value="none selected"/>
Is Reversed	<input type="radio"/> yes <input type="radio"/> no
Is Low Quality	<input type="radio"/> yes <input type="radio"/> no
Clone collection	<input type="text" value="none selected"/>
Species	<div style="border: 1px solid gray; padding: 2px;"><p>Xenopus laevis Xenopus tropicalis human mouse rat rhesus monkey zebrafish</p></div> <input type="button" value="Clear List"/>



Next specify what clone centric results will be provided and in what format.

Select the results you want to see back from the following list:

Check all

Clear all

Submit

- [Row Position](#)
- [Column Position](#)
- [Plate Number](#)
- [Clone collection](#)
- [Clone ID](#)
- [Sequence](#)
- [Repeat-masked Sequence](#)
- [Gb date created](#)
- [Date modified](#)
- [Genbank Accession Number](#)
- [Genbank Identifier](#)
- [Endedness](#)
- [Is Low Quality](#)
- [Is Reversed](#)
- [Seq length](#)
- [I.M.A.G.E. Library ID](#)
- [I.M.A.G.E. Library Name](#)
- [Tissue Count](#)
- [Species](#)
- [Tissue Type](#)
- [Vector name](#)



HTML version of clone-based query results.

I.M.A.G.E. Clone Query Results:

Your search returned 34 results! Here they are:

Result Number	PLATE	COLLECTION NAME	CLONE ID	SEQ	GB ACCNUM	ENDEDNESS	SPECIES	TISSUE TYPE
1	1651	LLAM	676104	<pre>gcggggttatctgtaaaggcctct aactttgtgaactgagtagcaag tagaagctaaggttaccacgaa ccccactttgcagttccccctcg tcttgctttgcagaatgaggatt ctcgcttgctgtccaccattac cttgggtggtcagttttgttga</pre>	AA208983	5	mouse	liver
2	1652	LLAM	676488	<pre>ttcggcacgtagggaaaaaactg cattagattttcccattaaacct tggatccaggtggacatgcagaa ggaagttagtcaccgggatac aaaccaaggtgctaaacactac ctaaagtctgctttaccacgga gtccaagtggcttacagctctg</pre>	AA208846	5	mouse	liver
3	1653	LLAM	676872	NULL	NULL	NULL	mouse	liver
4	1654	LLAM	677256	NULL	NULL	NULL	mouse	liver
5	1655	LLAM	677640	NULL	NULL	NULL	mouse	liver
6	1656	LLAM	678024	NULL	NULL	NULL	mouse	liver



Specify a library-based query.

Field Name		Data Input	
I.M.A.G.E. Library ID		<input type="text"/>	Convert
I.M.A.G.E. Library Name		<ul style="list-style-type: none">Appel/Eisen 15-19 hr embryoBaker mouse embryo e6.5Baker mouse embryo e7.5Barstead HPL-RB1Barstead HPL-RB2Barstead HPL-RB3Barstead HPL-RB4Barstead HPL-RB4.1Barstead HPL-RB5Barstead HPL-RB6	Clear List
Method		<ul style="list-style-type: none">full-length enrichedlarge-insertnormalizationnormalized/full-lengthnormalized/subtractedsagestandardsubtraction	Clear List
Species		<ul style="list-style-type: none">Xenopus laevisXenopus tropicalishumanmouseratrhesus monkeyzebrafish	Clear List
Tissue Type		<ul style="list-style-type: none">B-cellT-celladiposeadrenal glandbladderbloodbonebone marrowbowelbrain/CNS	Clear List

Submit

Clear Form



Similarly
specify what
library centric
results will be
provided.

Check all


Clear all

Submit

- [I.M.A.G.E. Library Name](#)
- [Tissue Count](#)
- [Species](#)
- [Strain](#)
- [Libr priming](#)
- [Vector digest](#)
- [Seq tag](#)
- [I.M.A.G.E. Library ID](#)
- [Method](#)
- [Host](#)
- [Libr provider](#)
- [Tissue Type](#)
- [Source age](#)
- [Source sex](#)
- [Source desc](#)
- [Source stage](#)
- [Tissue desc](#)
- [Vector name](#)



HTML version of library-based query results.



The I.M.A.G.E. Consortium

"Sharing resources to achieve a common goal - the discovery of all genes"


[Begin new search](#) | [Begin new Library search](#)


I.M.A.G.E. Library Query Results:


Your search returned **2** results! Here they are:

<u>Result Number</u>	<u>LIBR NAME</u>	<u>SPECIES</u>	<u>LIBR PRIMING</u>	<u>METHOD</u>	<u>SOURCE AGE</u>	<u>SOURCE DESC</u>	<u>TISSUE DESC</u>
1	NIH_MGC_87	human	oligo dT	full-length enriched	NULL	NULL	mammary adenocarcinoma, cell line
2	NIH_MGC_107	human	oligo dT	full-length enriched	NULL	NULL	adenocarcinoma, cell line

[Begin new search](#) | [Begin new Library search](#)

 [I.M.A.G.E. Consortium home page](#)

 [BBRP home page](#)

 [LLNL Programs, Projects, Centers and Consortia](#)



Other tools to mine I.M.A.G.E. information

Biotechnology Research PROGRAM **The I.M.A.G.E. Consortium**
 "Sharing resources to achieve a common goal - the discovery of all genes"

cDNA Plate and Collection Query

This page allows you to find plate and collection information about the libraries of your choice. Please note that you are limited to no more than 500 search results.

[Help using this page](#)

View results as HTML page
 Download results to a text file

Select Libraries to View:

Load Range:
Use First Letter of Library Name:

A B C D E F G
H I J K L M N
O P Q R S T U
V W X Y Z

All

Use Range of Letters:
 to Find

Search Library Name Using Keyword:
 Find

Current Range:

SAGE DL7
 SJD adult pectoral fin
 Schiller AWT-20
 Schiller AWT-20, RESP18 induced
 Schiller MAC13
 Schiller MAC16
 Schiller astrocytoma
 Schiller glioblastoma multiforme
 Schiller meningioma
 Schiller oligodendroglioma
 Schneider fetal brain 00004
 Soares 1NFLS
 Soares 1NFLS-S1
 Soares 1NLE
 Soares 1NHP6

Add-> Select All

Running List:

Soares 1NFLS
 Soares 3NBMS
 Stratagene mouse kidney (937315)

Remove Select All

View Reset

OR You may search by I.M.A.G.E. library ID:
 View

I.M.A.G.E. Consortium - Problem Clone Interface

Search Add Update Delete Help Quit

I.M.A.G.E. Consortium Problem Clones

Choose a Search method and enter one or more comma separated IDs if required.

Clone ID Submit

527061,1153242,3398810,938642

[LLNL Disclaimer](#)
 Web page maintained by webmaster@image.llnl.gov
 UCRL-MI-119648

Click on a problem_id to receive more detailed information.

problem_id	clone_id	problem_type	date_modified
53946	527061	Unknown Species	Sep 05 2001 12:10PM
53947	527061	Unknown Species	Sep 05 2001 12:10PM
52189	527061	Other (specify in comment field)	May 22 2001 09:55AM
68446	938642	Other (specify in comment field)	Feb 25 2002 01:53PM
51576	1153242	Incorrect Sequence (specify 3 or 5 in comment field)	Mar 26 2001 02:18PM
51577	3398810	Incorrect Sequence (specify 3 or 5 in comment field)	Mar 30 2001 11:21AM

Messages:

Query plates from libraries.

Query for reported problems.

Quality control for historical collection



Plates	Source	Well Error Rate
1-3705	Incyte	13
	LLNL Master	10
	Research Genetics	12
	Resource Center of HumanGenome Project	10
	ATTC	11
3,796-6000	Incyte	7
	LLNL Master	7
	Research Genetics	10
	Resource Center of Human Genome Project	12

QC on-going

LLNL Replication

Master vs. GenBank



Months	Well error rate	Plate Error Rate	Well error rate	Plate Error Rate
6/2000	1 (1,3)	0	7 (4,11)	2
10/2000	1 (0,3)	0	1 (0,3)	2
12/00	0 (0,2)	2	1 (0,3)	2
1/01	2 (1,4)	0	6 (4,11)	3
2/01	1 (0,3)	0	2 (1,5)	2
3/01	2 (1,5)	2	2 (1,5)	0
4/01	1 (0,3)	2	2 (1,4)	0
5/01	0 (0,1)	0	2 (1,5)	0
6/01	1 (0,3)	0	1 (0,4)	0
7/01	1 (0,4)	0	2 (1,6)	0
8/01	2 (1,3)	0	3 (2,6)	0



Ongoing QC results

RO (LLNL Master Plate) compared to R3 (from RO, by LLNL, for LLNL)

Month	Number of Clones		Well Error (%)		Number of Plates		Plate Error (%)	Calc Date	Link to Details
	Compared	In Error	Rate	95% C.I.	Total	In Error			
6 / 2000	391	5	1	(0, 3)	57	0	0	Aug 15 2001	6 / 2000
10 / 2000	409	5	1	(0, 3)	56	0	0	Oct 30 2001	10 / 2000
11 / 2000	*	*	1	(*,*)	60	0	0	Feb 20 2001	11 / 2000
12 / 2000	407	2	0	(0, 2)	58	1	2	Aug 15 2001	12 / 2000
1 / 2001	344	7	2	(1, 4)	57	0	0	Aug 15 2001	1 / 2001
2 / 2001	385	5	1	(0, 3)	59	0	0	Aug 15 2001	2 / 2001
3 / 2001	323	8	2	(1, 5)	50	1	2	Aug 15 2001	3 / 2001
4 / 2001	372	5	1	(0, 3)	57	1	2	Aug 15 2001	4 / 2001
5 / 2001	345	0	0	(0, 1)	48	0	0	Oct 30 2001	5 / 2001
6 / 2001	359	3	1	(0, 3)	52	0	0	Nov 14 2001	6 / 2001
7 / 2001	304	4	1	(0, 4)	49	0	0	Nov 14 2001	7 / 2001
8 / 2001	389	6	2	(1, 3)	58	0	0	Jan 22 2002	8 / 2001
9 / 2001	154	3	2	(1, 6)	25	1	4	Jan 22 2002	9 / 2001

Error in replication @ LLNL

On-going
Comparing master to GenBa



Month	Number of Clones		Well Error (%)		Number of Plates		Plate Error (%)	Calc Date	Link to Details
	Compared	In Error	Rate	95% C.I.	Total	In Error			
6 / 2000	230	15	7	(4, 11)	40	1	2	Aug 15 2001	6 / 2000
10 / 2000	310	3	1	(0, 3)	50	1	2	Oct 30 2001	10 / 2000
11 / 2000	198	4	2	(1, 5)	32	0	0	Aug 15 2001	11 / 2000
12 / 2000	276	3	1	(0, 3)	46	1	2	Aug 15 2001	12 / 2000
1 / 2001	187	12	6	(4, 11)	36	1	3	Aug 15 2001	1 / 2001
2 / 2001	253	5	2	(1, 5)	45	1	2	Aug 15 2001	2 / 2001
3 / 2001	266	6	2	(1, 5)	46	0	0	Aug 15 2001	3 / 2001
4 / 2001	245	4	2	(1, 4)	43	0	0	Aug 15 2001	4 / 2001
5 / 2001	201	4	2	(1, 5)	34	0	0	Oct 30 2001	5 / 2001
6 / 2001	239	3	1	(0, 4)	41	0	0	Oct 30 2001	6 / 2001
7 / 2001	152	3	2	(1, 6)	27	0	0	Oct 31 2001	7 / 2001
8 / 2001	287	9	3	(2, 6)	46	0	0	Jan 30 2002	8 / 2001
9 / 2001	16	12	75	(47, 92)	38	34	89	Jan 22 2002	9 / 2001



Next for I.M.A.G.E. Informatics

- Extensive expansion of query tools and data access
- IMAGEne non-species specific
- Analysis of human cluster candidate genes and singletons
- Redo of web site, easier to navigate



MUCH influenced by public needs.....you have a say!



Acknowledgements

- LLNL
 - Christa Prange, I.M.A.G.E. PI
 - Tim Harsch, Amber Johnston, Julie Amundson
- Sponsors
 - DOE, Marv Stodolsky
 - NIH, Bob Strausberg

image.llnl.gov

This work was partially funded by the NIH and was performed under the auspices of the U.S. Department of Energy by the University of California, Lawrence Livermore National Laboratory under contract no. W-7405-Eng-48.