

Annotating Non-coding Regions of the Human Genome

Robert Cedergren Bioinformatics Colloquium Montreal, Canada 2008.11.03, 16:30-17:30

> Mark B Gerstein Yale (Comp. Bio. & Bioinformatics)

Slides from Lectures.GersteinLab.org





(Please read permissions statement.) Paper references mostly from Papers.GersteinLab.org. See streams.gerstein.info on photos & images

> (Genome tech and Genome annotation talk, including: mismatch, uniarray, DART, pgenes-general, pseudofam, encode-pgenes, sdcnvcorr, gen-encode, sirna-pgene, what-is-gene [I:CEDERGREN]. Fits into ~50'.)



 % exon).

 regions,

 nes,

 How should

 [IHGSC, Nature 409, 2001]

 2001: Most of the genome is not coding (only ~1.2% exon). It consists of elements such as repeats, regulatory regions, non-coding RNAs, origins of replication, pseudogenes, segmental duplications....What do these elements do? How should they be annotated? [Venter et al. *Science* 29, 2001]



2007 : Pilot results from ENCODE Consortium on decoding what the bases do

- 1% of Genome (30 Mb in 44 regions)
- Tiling Arrays to assay Transcription & Binding
- Multi-organism sequencing and alignment
- Careful Annotation
- Variation Data

[IHGSC, *Nature* 409, 2001] [ENCODE Consortium, *Nature* 447, 2007]

- Binding Sites and Regulatory Nets
- Novel Transcribed Regions
- Pseudogenes
- SDs and SVs (Segmental Duplications & Structural Variations)

Different Views of the Function of Junk DNA

[NY Times, 26-Jun-07]

Human DNA, the Ultimate Spot for Secret Messages (Are Some There Now?)

By DENNIS OVERBYE

ESSAY

In Douglas Adams's science fuction classic, "The Bitchikar's Gaide to the Galaxy," there is a character by the name of Elarthartfast, who designed the type ds of Norway and left his signature in a glacier.

I was reminded of Startburthast recently as I was rying to group the implications of the test of a team of Japanese generics who annuanced that they had might relativity to a bacteriam, our of.

Using the same code that computer keyboards use, the Japanese group, led by Massers Younds (Neto Usi-Versit), write four copies of Albert Eaststeir's famous hermain, E-sec', along with "LRO," the date blait the young Einstein derived it, indo the bacterium's genome, the 40-million-long string of A's, G's, T's and C's that determine everything the limits long is and everything if is ever going in be

The point was not to celetrate Elements. The feat, they said in a paper published in the journal Bistechned on Progress, was a demonstration of DNA as the althmate softernation storage material, able to withstand floads, servorises, time and the changing fashnaos in nechenalogy, not to mention the ability to be improved with firstle ambitrarise trademark fabries — limite "Made by Monsaito" cage, say.

As an design they have accomplished at least a part of the dream that Jarme Lanter, a computer scientus and monocian, and David Salzer, a biologist at Columbia, enanciated in 1990. To create the ultimate time capsule as part of the millenniam feativities at this newspaper, they proposed to encode a year's worth of the New York Times magazine into the jusk DNA of a cockroach. "The archival acceleracies will be a robust repository," Mr. Latier wrone, "able to survive almost all conceivable scenarios."

If cockmaches can be archives, why not us? The haman persons, for example, consists of some 2.8 billion of those letters - the equivalent of about 750 megabytes of data - but only about 2 percent of it goes into companing the 22,000 or so genes that make so what we are.

The remaining 97 percent, so-called junk DNA, loads like globerish. It's the dark matter of anner apace, We durk know what it is sogring to or about us, but within that sea of megalaytes there is plenty of room for the imagination to roads, for trademark labels and much marr. The King James Bille, to pick one obvious example, only amounts to about five megalaytes.



If a bacterium can be encoded with E=mc², if cockroaches can be archives, why not us?

Insvitably, if you are me, you begin to wonder if there is already something written in the warm wet archive, whether or not some blartinatriat has already been here and we curselves are walking around with his le irademark tags or mover wrigging and sontaging and folded inside us. Gill Bejerana, a genetican at the University of California, Santa Ciru, who mentioned Slartibortflast to me, pointed out that the prohlem with raising this question is that people who look will see messages in the generic even if they aren't there — the way people have claimed in recent years to have found screet codes in the filleb.

Nevertheless, no less a personage than Prancis Crick, the co-discoverer of the double helix, writing with the chemist Leslie organ, new at the Sali Instatute in San Diego, suggested in 1073 that the primitive Earth was infected with DNA broadcast through space by an alten species.

As a result, it has been suggested that the search for extranerrestrial inselligence, or SETL should look inward as well as ourward. In an article in New Scientist, Paul Davies, a cosmologist at Artizona State University. Using the same code that computer keyboards use, the Japanese group... wrote four copies of Albert Einstein's famous formula, E=mc2... into the bacterium's genome... In so doing they have accomplished at least a part of the dream that Jaron Lanier, a computer scientist and musician, and David Sulzer, a biologist at Columbia, enunciated in 1999. To create the ultimate time capsule as part of the millennium festivities at this newspaper, they proposed to encode a year's worth of the New York Times magazine into the junk DNA of a cockroach. "The archival cockroach will be a robust repository," Mr. Lanier wrote, "able to survive almost all conceivable scenarios."

mice, chickens and dogs for at least 300 million years

of them had turned out to be playing important com-

tery." he said, toying that even regular genes that do

hits of DNA that neither help nor annuy an organism

mand and control functions.

mutane even more rapidly.

But Dr. Bejerano, one of the discoverers of these

"Why they need to be so conserved remains a mys-

thing undergo more change over time. Most junk

The Japanese team proposed to sidestep the mula-

on problem by inserting redundant copies of their mass

sage into the genome. By comparing the readouts, they

said, they would be able to recover Einstein's formula

even when up to 15 percent of the original letters in the

"altrasonserved" strings of the genome, said that many

sections of junk DNA seem to be markedly resistant to Start

How might we annotate a human text ? Ethnologue.com. **Color** is **Function** Lines are

[B Hayes, Am. Sci. (Jul.- Aug. '06)]

Similarity

The Semicolon Wars Brian Hayes

F YOU WANT TO BE a thorough-_____going world traveler, you need to learn 6,912 ways to say "Where is the toilet, please?" That's the number of languages known to be spoken by the peoples of planet Earth, according to

If you want to be the complete polyglot programmer, you also have quite a challenge ahead of you, learning all the ways to say:

printf("hello, world\n");

(This one is in C.) A catalog maintained by Bill Kinnersley of the University of Kansas lists about 2,500 programming languages. Another survey, compiled by Diarmuid Piggott, puts the total even higher, at more than 8,500. And keep in mind that whereas human languages have had millennia to evolve and diversify, all the computer languages have sprung up in just 50 years. Even by the more-conservative standards of the Kinnersley count, that means we've been inventing one language a week, on average, ever since Fortran.

For ethnologists, linguistic diversity is a cultural resource to be nurtured and preserved, much like biodiversity. All human languages are valuable; the

Every programmer knows there is one true programming language. A new one every week

a good-enough notation-for expressing an algorithm or defining a data structure.

There are programmers of my acquaintance who will dispute that last statement. I expect to hear from them. They will argue-zealously, ardently, vehemently-that we have indeed found the right programming language, and for me to claim otherwise is willful ignorance. The one true language may not yet be perfect, they'll concede, but it's built on a sound foundation and solves the main problems, and now we should all work together to refine and improve it. The catch, of course, is that each of these friends will favor a different language. It's Lisp,

cide which end of a boiled egg to crack. This famous tempest in an egg cup was replayed 250 years later by designers of computer hardware and communications protocols. When a block of data is stored or transmitted, either the leastsignificant bit or the most-significant bit can go first. Which way is better? It hardly matters, although life would be easier if everyone made the same choice. But that's not what has happened, and so quite a lot of hardware and software is needed just to swap ends at boundaries between systems.

This modern echo of Swift's Endian wars was first pointed out by Danny Cohen of the University of Southern California in a brilliant 1980 memo, "On holy wars and a plea for peace." The memo, subsequently published in Computer, was widely read and admired; the plea for peace was ignored.

Another feud-largely forgotten, think, but never settled by truce or treaty-focused on the semicolon. In Algol and Pascal program statements have to be separated by semicolons. For example, in x := 0; y := x+1; z := 2 the semicolons tell the compiler where one statement ends and the next begins. C programs are also peppered with semi-

Overview of the Process of

Annotation of non-coding Regions

• Basic Inputs

- 1. Doing large-scale similarity comparison, looking for repeated or deleted regions
- 2. Determining experimental signals for biochemical activity (e.g. transcription) across each base of genome

Results of Analyzing Similarity Comparison

- 1. Finding large repeated or deleted blocks (e.g. CNVs) as a function of degree of similarity
 - 1. within reference human genome
 - 2. within human population
 - 3. between related organisms (e.g. mouse)
- 2. Finding smaller "exon-level" similarities (e.g. pseudogenes)



Determining experimental signals for biochemical activity (e.g. transcription) across each base of genome

- Development of Sequence (and Array) Technology
 - Normalizing & Scoring Signal, Correcting Artifacts, Segmenting to create Small Annotation Blocks
- Output of Production Pipelines and Surveying a Single Type of Annotation on a Large-scale
 - Clustering Small Blocks into Larger
 Ones, Surveying
- Integrated Analysis Connecting Different Types of Annotation
 - Building networks and beyond

Outline

- Calling Blocks
 - Issues in Signal
 Processing Transcriptional
 Data from Tiling Arrays
- Clustering Blocks
 - Scoring transcribed regions and clustering these together into larger blocks
- Formal Annotation
 - ◊ Pseudogenes

- Integration with conservation
 SDs and CNVs
 - pseudogenes and seq. constraint
- Integration with activity
 - pseudogenes and transcriptions
- Future of Gene Annotation
 - ◊ What is a gene post encode?

Signal Processing: Normalizing Signal and Finding Initial Annotation Blocks ("Hits")



pers. **photo**, see streams.gerstein.info

A Starting Point: Noisy Raw Signal from Tiling Arrays (Transcription)



| | | <u>Specifi</u> | <u>c & Non-</u> |
|---|---------------|--|---|
| | | specific (| Cross-Hyb. |
| | | Perfect match (probe binding in Specific cross-h targets with a si Non-specific cro targets with mat general sticking | PM): Intended target hyb.: probes binding non-PM mall number of mismatches pss-hyb.: probes binding ny mismatches, due to ess of oligos |
| R | | | |
| | Perfect Match | Specific Cross-hyb. | Non-specific Cross-hyb. |

<u>Non-Specific Cross Hyb.</u> (Sequence Effects)



[Seringhaus et al., BMC Genomics (in press)]



Yeast ACT1 Gene

Human HBG2 Gene



Number of Mismatches







Types of Mispairs (probe on array is first)

18 Lectures.C

Observing Non-specific Cross-hyb. (Probe sequence effects)

Nimblegen 50th Quantile



Source: Royce, T.E., et al (2007), *Bioinformatics*, 23, 988-97

Quantile Normalization

Gene expression quantile normalization

Quantile normalization has proven to be the most effective way to normalize replicate gene expression arrays.



Source: Bolstad, B.M., et al (2003), Bioinformatics, 19, 185-93.



Iterated Quantile Normalization to Correct for Non-specific Cross-hyb.

- Adapt Bolstad et al (2003) approach to tiling arrays
- Force distributions with a given nt at each position to be same
- Distributions at other positions now different so iterate
- Also, robust adaptation of Naef & Magnasco (2003)

Measuring Specific Cross-Hyb

Source: Royce, T.E., et al (2007), *Nucleic Acids Res.*, 23, 98-97

Measuring Specific Cross-Hyb

- Start with Cheng et al. (2005) tiling of human genome at 5 nt resolution giving expression profiles across various cell lines
- Correlation betw. probe pairs computed across cell lines' expression profiles and tabulated vs. number of mismatches
- The mean correlation coefficient was computed for each mismatch bin (blue series).
- The number of pairs is plotted as orange bars.



Source: Royce, T.E., et al (2007), *Nucleic Acids Res.*, 23, 98-9%

Proof of principle test to "exploit" this



- Using Cheng et al. (2005), predict gene expression levels (and profiles across tissues) for genes on part of chr. #6
- ...Based on closest cross-hyb tiles on part of chr. #7
- Then compare to measured levels and profile on #6

Nearest Nbr Search on Virtual Tiling

b microarray hybridizations

a virtual tiling

GTECATGEGANCTTGETGEGACCCC TCGGACAGGTAGACAGGCCGCGAGA TAGACATAGACAATGCATACGCT AACUSTTOCKCOCCCGGGATGTAA 25nt TAGACGTAGACAATTCATACCCTAA c similarity search d profile assignment from nearest-neighbor Lectures.GersteinLab.org (c) 2007 **Query Tile** AGACUTAGACAAT INCLUCIONA CATACCCTAA **Probe Sequence** Database ₽ТАС GTCCATERCARCTTCCTOEGGACCCC INNERACACHINECCOCAGA nearest-neighbor search TACACMENCRICATECHERCOCTAC AACGCTTGCGCGCCCGGGBTGTAAC

Agreement between predicted tile expression profile and actual one

- Correlated predicted profiles with the actual profiles of gene expression across cell lines
- Much more correlation than expected by chance (dist. centered on 0)



Very Strong ROC Curve: Most genes are accurately detected using nearest-neighbor features' signals

Illustrates great magnitude of cross-hyb. on hi-density arrays

- High feature density arrays inadvertently resurrecting generic n-mer concept (van Dam & Quake, 2003)
- Suggests that tiling arrays could be exploited to create universal arrays
- Gold std. set of known expressed genes. How well do we find.
- A set of known positives was defined as the Refseq genes with at least 75% transfrag coverage. A set of known negatives was constructed by permuting the sequences in the set of known positives. For various thresholds, sensitivity and specificity were computed and then plotted.



Royce, T. E. et al. Nucl. Acids Res. 2007 35:e99



Annotating a single type of signal on a large-scale: Clustering and Classifying Unannotated Transcription (TARs)

pers. **photo**, see streams.gerstein.info









ENCODE Regions (30 Mb)

Locations of TARs

Of the approx 7,000 Novel TARs

- 955 are assigned to known genes
- 1,463 are clustered into ~200 Novel Loci

•DART Classification has been experimentally validated with some small scale experiments

- ◊ RT-PCR & Sequencing
- ◊ 18/46 (39%) confirmed by RT-PCR
- ◊ 4/5 Sequenced Products Map uniquely to correct genomic region

Rozowsky et al. Genome Research (2007)

DART Classification has been experimentally validated with some small scale experiment ♦ RT-PCR & Sequencing

Results:

18/46 (39%) confirmed by RT-PCR

4/5 Sequenced Products Map uniquely to correct genomic region

Rozowsky et al. Genome Research (2007)







Chr21

Example predicted structured RNAs (using RNAz)

Overlap of predicted structured RNAs with the union of TARs/Transfrags and the "moderate" set of sequence-constrained elements



[>700 candidate structured RNAs predicted in 1% of the reference genome] Stefan Washietl, Jakob Pedersen, Jan Korbel *et al.* (2007) *Genome Res* 17:852-864



Integrative Analyses:

Annotating Pseudogenes and relating them to functional signal and measures of conservation

<u>Pseudogenes are among the most</u> <u>interesting intergenic elements</u>

- Formal Properties of Pseudogenes (Ψ G)
 - ◊ Inheritable
 - Or Homologous to a functioning element
 - ◊ Non-functional*
 - No selection pressure so free to accumulate mutations
 - Frameshifts & stops
 - Small Indels
 - Inserted repeats (LINE/Alu)
 - What does this mean? no transcription, no translation?...

Identifiable Features of a Pseudogene (ψRPL21)



Distribution of Human Pseudogenes (for RPL21) across the chromosomes



Two Major Genomic Remodeling Processes Give Rise to Distinct Types of Pseudogenes



Gerstein & Zheng. Sci Am 295: 48 (2006).

Why Study Pseudogenes?

Cause errors in sequence databases

- > 8000 retropseudogenes in human
- Contamination in Ensembl
- 25% in C. elegans? [Mounsay, Genome Research, 2002]

Interfere with study on functional genes

- Cross-hybridation in micro-array and RT-PCR.
- Some pseudogenes have regulatory roles

\blacktriangleright Ψ G are "genomic fossils"

- Study the evolution of genes and genomes
- Measure mutation/insertion rates



[Ruud, Int. J. Cancer 1999]
Pseudogenes: Tools



- Integrating heterogeneous, **Dynamically Changing** Annotation
 - ♦ Changing sequences, gene predictions, repeats
- Track (slightly) changing objects across genome builds
 - Versioning and exact temporal reconstructability
- Fixed <u>Sets</u> of Pseudogenes
 - Orresponding to particular types of analyses or papers
- Generalizable Class Structure
 - ♦ fragments, alignments, collections, pseudogenes
- EAV
 - ♦ Elexible Annotation for extended characteristics
- Interface with Uniprot & UCSC



Karro et al., NAR (2007)



Pseudogene Accession Number (Isid format) Start Step Type 22 22 4464240 14464833 ene.org 9606.Pseudogene 5542 NSP00000347298 Human chr22 mb1

Pseudofam

- A Pseudogene Families Database
- Highlighted Features:
 - Browse families
 - Statistics
 - Enrichment
 - P-value, etc
 - Search families
 - Pfam ID/Acc
 - Ensembl ID
 - Pgene ID
 - Correlate families
 - Genes
 - Pgenes
 - Parents

| 2000 C 100 C | | |
|--|------------------|----------------------------------|
| Overview | | |
| | | View Genome: (All |
| | | |
| Summary Statistics | | |
| | F AR 4 | |
| Species Name: | 10 m | |
| Protein Families: | 3,820 | |
| Pseudogene Femilies: | 2,985 | |
| Total Genes: | 219,663 | |
| Total Parents: | 32,660 | |
| Total Pseudogenes: | 125,272 | |
| Pseudopene to gene Ratio | 0.57 | |
| Pseudogene-to-parent Ratio: | 3.84 | |
| Parent to-gene Ratio: | 0.15 | |
| | | Vew Chart: { Failure Composition |
| femily Composition | | |
| Non-pseudogene vs Pseudogene Family | and the state of | Non-parent Cene us Parent Gene |
| | | |

[Lam et al., NAR DB Issue (in press, '09)]

Pseudofam Data Sources

- Ensembl Database
 - Peptides
 - Exons
 - DNAs
- BioMart
 - ENSP-Pfam Mapping
- Pfam
 - Pfam Alignments
- PseudoPipe
 - Pseudogene Identification



Pseudofam Data Generation & Alignment



- Identify pseudogenes by proteins
- Map parent proteins to protein families
- Assign pseudogenes to their parent families
- Align the pseudogenes in pseudogene families.
- Calculate the key statistics and organize the data into database.



58 Lectures.GersteinLab.org (c) 2007

[Lam et al., NAR DB Issue (in press, '09)]

The Enrichment Statistics

• Hypergeometric Dist.

$$f(k; N, m, n) = \frac{\binom{m}{k}\binom{N-m}{n-k}}{\binom{N}{n}}.$$

- ◊ k: parent genes in a fam
- ◊ n: genes in a fam
- In: overall total parent genes
- N: overall total genes
- For example, in fam x:
 - ◊ 3 parent genes in x
 - \diamond 5 genes in x
 - ◊ 8 total parent genes in genome
 - ◊ 20 genes in genome



- Expected parents: 2
 Enrichment: 2/2, 1 5
 - ♦ Enrichment: 3/2=1.5
 - \diamond P-value: Pr(k>=3)=0.3

Upper Ontology



60 Lectures.GersteinLab.org (c) 2007



Pseudogene Assignments to the Human Genome

- Draft Assignment to Whole Genome
 - Automatic pipeline currently gives
 ~22644 [pseudogene.org]
- Eventual Annotation Process in framework of ENCODE project
 - ♦ Hybrid
 - Pipeline runs, manual curation flagging difficult cases, pipeline improvement, further curation...
- Pilot ENCODE Annotation
 - 4 automatic pipelines: retroFinder+pseudoFinder (UCSC), PseudoPipe (Yale), GIS
 - ◊ HAVANA manual
 - ◊ 201 pseudogenes vs ~400 genes

- Issues in Pilot
 - Comparing protein or transcript v genomic DNA, filtering, application of rules
 - ◊ What is a pseudogene?
 - ◊ Different criteria
 - ◊ Conservative approach here
 - ◊ Can't overlap gene annotation
 - Need to have a protein alignment

Zheng et al. (2007) Gen. Res.







<u>Overall</u> <u>Results:</u> <u>Regional</u> <u>Distribution</u>

201 pseudogenes 77 non-processed 124 processed

Zheng et al. (2007) Gen. Res.

browser + pseudogene.org/ENCODE



Vast Amounts of Different Data Types to Integrate in pilot ENCODE

- Determining experimental signals for biochemical activity across each base of genome
- Large-scale sequence comparison in relation to the human genome

| Feature Class | Expt. Tech. | Numb. Expt. Data Pts. | |
|-------------------------------------|--|--------------------------|--|
| Transcription | Tiling array, Integrated annotation | 63,348,656 | |
| 5' Ends of transcripts | Tag sequencing | 864,964 | |
| Histone modifications | Tiling array | 4,401,291 | |
| Chromatin structure | QT-PCR, Tiling array | 15,318,324 | |
| Sequence- specific factors | Tiling array, tag sequencing, Promoter assays | 324,846,018 | |
| Replication | Tiling array | 14,735,740 | |
| Computational analysis | Computational methods | NA | |
| Comparative sequence analysis | Genomic sequencing, multi- sequence alignments, computational analyses | NA | |
| Polymorphisms | Resequencing, copy number variation | NA | |

Integration of Different Types of Annotation: Pseudogenes with Sequence Constraint



Using phastOdd value to examine neutral evolution of pseudogenes



representative pseudogenes drawn from 201 total С Β D Ε F Α \boxtimes \boxtimes \boxtimes \boxtimes \boxtimes \boxtimes human - \boxtimes X \bowtie chimp - \boxtimes \boxtimes \boxtimes \boxtimes \boxtimes baboon - \boxtimes \boxtimes \boxtimes \boxtimes \boxtimes macaque - \boxtimes \odot \boxtimes marmoset - \boxtimes \square \boxtimes (\cdot) galago - \square \boxtimes \odot rat -(·) \boxtimes \boxtimes \boxtimes \odot mouse - \boxtimes (\cdot) (\cdot) rabbit -(·) \boxtimes (\cdot) (\cdot) (\cdot) (·) cow - \boxtimes \boxtimes \odot (\cdot) \boxtimes dog - \boxtimes \boxtimes (\cdot) rfbat - (\cdot) \boxtimes \boxtimes \odot shrew - (\cdot) \boxtimes \boxtimes \odot (\cdot) armadillo - (\cdot) \odot \boxtimes \boxtimes elephant - (\cdot) \odot \boxtimes \boxtimes tenrec - (\cdot) (\cdot) \odot \boxtimes (\cdot) monodelphis - (\cdot) \odot \boxtimes (\cdot) platypus - (\cdot) \odot \odot chicken -(·) (\cdot) \odot \odot (\cdot) (\cdot) \odot \boxtimes xenopus - \boxtimes (\cdot) \boxtimes \odot tetraodon -(·) \boxtimes (\cdot) (\cdot) zebrafish -

<u>History</u> <u>of</u> <u>Pseudogene</u> <u>Preservation</u>

Based on alignment from ENCODE MSA group

Zheng et al. (2007) Gen. Res.

Absent

Present with Disablement

Present without Disablement

<u>Most Processed Pseudogenes are Primate</u> <u>Specific Created by Recent (<45 MYA)</u> <u>Retrotranspositional Activity</u>



Sequence Decay of Pseudogenes, Approximately Neutral



Sequence Decay of Pseudogenes Relative to their Immediate Genomic Context



Analyzing Repeated Blocks in the Genome (SDs & CNVs)



pers. **photo**, see streams.gerstein.info

SEGMENTAL DUPLCATIONS AND COPY NUMBER VARIANTS ARE RELATED PHENOMENA AND HAVE BEEN CREATED BY SEVERAL DIFFERENT MECHANISMS





NAHR (Non-allelic homologous recombination)

Flanking repeat (e.g. Alu, LINE...)



NHEJ

(Non-homologous-endjoining)

No (flanking) repeats. In some cases <4bp microhomologies

Association of SDs and CNVs with pseudogenes

- CNVs are the raw form of variation producing duplicated elements
- Segmental Duplications (SDs) are fixed forms of CNVs/SVs. They give rise to duplicated genes and (eventually) protein protein families
- Thus, we expect, duplicated pseudogenes (failed duplications) to occur in SDs.
- CNVs and SDs tend to be enriched in environmental response genes, matching a patterns previously found for duplicated pseudogenes

[Korbel et al., COSB (in press, '08)]





В

Successfully duplicated genes (SDs spanning entire genes



C. Unsuccessful duplicates (duplicated genes inactivated by disruption of coding sequence



Pseudogene families and Segmental Duplications (SDs)



- SDs comprise ~5-6% of the human genome but contain ~17.8% genes, 45.7% duplicated pgenes and 21.6% processed pgenes
- Relative values of correlation coefficients in the plots above consistent with the observation that SDs contain more pgenes than parent genes

PERFORM LARGE SCALE CORRELATION ANALYSIS TO DETECT REPEAT SIGNATURES OF SDs AND CNVs



OLDER SDs ARE MUCH MORE LIKELY TO BE FORMED BY ALU ELEMENTS



ANOTHER FUNCTION FOR PSEUDOGENES: SERVING AS REPEATS FOR MEDIATING NAHR



FOCUSSING ON SDS: SDS CAN PROPAGATE THEMSELVES, WHICH LEADS TO A POWER-LAW DISTRIBUTION



Corollary

- SDs can mediate NAHR and lead to the formation of CNVs
- CNVs can become fixed and then be SDs
- Such mechanisms ("preferential attachment") are well studied in physics and should leads a very skewed ("power-law") distribution of SDs.



FOCUSSING ON SDS: SDs COLOCALIZE WITH EACH OTHER



Corollary

- SDs can mediate NAHR and lead to the formation of CNVs
- CNVs can become fixed and then be SDs
- SDs of similar age should co-localize better with each other:



ASSOCIATIONS ARE DIFFERENT FOR SDs AND CNVs

CNVs ARE LESS ASSOCIATED WITH SD association with repeats **SDs THAN THE GENERAL SD TREND** 0.27 CNV 0.21 Association 0.094 0.07 with SDs Microsatellite Pseudogenes LINE Alu 0.31 <0.001 (<0.001) 0.001 0.046 0.11 **CNV** association with repeats 0.0739 0.048 0.0466 0.0006 >99% SDs* CNVs Microsatellite **Pseudogenes** LINE Alu <0.001 0.92 0.046 0.001

[Kim et al. Gen. Res. (in press, '08), arxiv.org/abs/0709.4200v1]

82

82

ANALYZING SEQUENCED BREAKPOINTS CONFIRMS THE RESULTS FROM THE COARSE GRAINED ANALYSIS

| Repeat Type | Frequency | Global enrichment | p-value | Local enrichment | p-value |
|----------------|-------------|----------------------|-----------------|---------------------|----------|
| Alu | 0.09 | 0.94 | 3.24E-01 | 1.13 | 1.74E-01 |
| | | | | | |
| SD | 0.41 | 2.57 | 2.14E-07 | 1.17 | 2.64E-01 |
| L1 | 0.24 | 1.48 | 1.03E-07 | 1.12 | 7.16E-02 |
| L2 | 0.01 | 0.47 | 1.72E-02 | 0.52 | 2.31E-02 |
| Microsatellite | 0.03 | 3.91 | 6.74E-11 | 3.11 | 2.99E-07 |
| LTR | 0.09 | 1.14 | 1.71E-01 | 0.89 | 1.97E-01 |
| PPgene | 0.01 | 2.08 | 9.55E-02 | 1.66 | 1.98E-01 |
| GC | 0.39 | 0.96 | 7.24E-03 | 0.97 | 3.00E-02 |
| | | | | | |
| | | | .ATCAAGG CCGGAA | A | |
| Exact m | Exact match | | | | |
| Local e | nvironment | | | | |



ELEMENTS FOR GENOME About 40 million years ago

- there was a burst in retrotransposon activity
- The majority of Alu elements stem from that time
- This, in turn, led to rapid genome rearrangement via
- The resulting SDs, could create more SDs, but with Alu activity decaying, their creation slowed

Integration of Different Types of Annotation: Pseudogenes with Intergenic Biochemical Activity





Connecting TARs (TxFrags) in Integrative fashion to different types of Annotation

- Single Ex. of Pseudogene Intersecting with Transcriptional and Regulatory Evidence
- Are integrated experiments comparable -- i.e. done on consistent cell lines, on same coordinate sys., &c.

Intersection of Pseudogenes with Transcriptional Evidence

| 9 | | | | | |
|--------------------|--------------------|------|-------|-----------|---------------|
| | TAR / transfrag | CAGE | DiTag | RACEfrag | EST / mRNA |
| TAR / transfrag | 105 * | 8 | 2 | 5 | 14 |
| CAGE | | 8 | 1 | 0 | 1 |
| DiTag | | | 2 | 0 | 0 |
| RACEfrag | | | | <u>14</u> | 5 |
| EST / mRNA | | | | | 21 |
| | | | | | |

Excluding TARs (due to cross-hyb issues)

Targeted RACE expts to 160 pseudogenes, gives <u>14</u>

Total Evidence from Sequencing is 38 of 201 (with 5 having cryptic promotors)

Integrated

Integrating Transcriptional Evidence with Gene Annotation and Sequence Constraints



(c) 2007 88 Lectures.GersteinLab.org

Zheng et al. (2007) Gen. Res.
Extension to Whole Genome

- 233 Transcribed from ~8000 Processed Pseudogenes
- Evidence for Transcription
 - ◊ 8% Refseq mRNAs
 - ♦ 32% Unigene consensus sequences
 - ♦ 72% dbEST expressed sequence tags
 - 32% Oligonucleotide microarray data (extra support)
- Highly decayed
 - ◊ Fraction with Ka/Ks ≥ 0.5 is 54%



Biochemically Active Regions Don't all Appear to be Under Constraint

- Integrating & averaging results over larger and larger sets
- Comparison of integrated quantities

Grand Summary: Biochemical Activity vs. Sequence Constraints

- Constrained sequence
- Experimental annotation
- Not all constrained sequence annotated in some fashion
- Exactly how things are defined in terms of overlap?

"At the outset of the ENCODE Project, many believed that the broad collection of experimental data would nicely dovetail with the detailed evolutionary information derived from comparing multiple mammalian sequences to provide a neat 'dictionary' of conserved genomic elements, each with a growing annotation about their biochemical function(s). In one sense, this was achieved; the majority of constrained bases in the ENCODE regions are now associated with at least some experimentally-derived information about function. However, we have also encountered a remarkable excess of unconstrained experimentally-identified functional elements, and these cannot be dismissed for technical reasons. This is perhaps the biggest surprise of the pilot phase of the ENCODE Project, and suggests that we take a more 'neutral' view of many of the functions

and suggests that we take a more 'neutral' view of many of the functions conferred by the genome. "



[ENCODE Consortium, *Nature* 447, 2007]



pers. photo, see streams.gerstein.info



What are Active Pseudogenes Doing?

Potential for <u>Gene</u> <u>Regulation via</u> <u>endo-siRNA</u>

Recent Discoveries in Mouse & Fly

Czech, B. *et al. Nature* 453, 798–802 (2008). Ghildiyal, M. *et al. Science* 320, 1077–1081 (2008). Kawamura, Y. *et al. Nature* 453, 793–797 (2008). Okamura, K. *et al. Nature* 453, 803–806 (2008). Tam, O. H. *et al. Nature* 453, 534–538 (2008). Watanabe, T. *et al. Nature* 453, 539–543 (2008).

[Sasidharan & Gerstein, Nature ('08)]

Very Speculatively, Papers Blur Boundaries betw. siRNAs and miRNAs



Genes & Pseudogenes

(b) Dead Pseudogene



Genes or Pseudogenes?

(b) Dead Pseudogene



Zheng & Gerstein, TIG (2007)

Promoter Exon Pseudo-Exon RNA 🛪 Mutations disrupting protein coding

Genes or Pseudogenes?



*: unannotated spliced transcript products detected sequences(+) primer regions refSeq (+) DRG1 5' q12.2 30,130,000 30,120,000 refSeq (-) detected 1+ sequences (+) 30,917,600 30,918,000 30,918,200 30,918,400 30,918,600 30,918,800 30,919,000 30,919,2 detected sequences (+) * detected sequences(-) 5 refSeq (+) 5' TIMP3 FBX07 **q12.3** 31,300,000 31,600,000 31,700,000 31,200,000 31,500,000 primer regions refSeq SYN3 (-) detected sequences (+) 70,700,000 70,700,500 70,700,000 ш ш detected detected sequences (-) sequences (-)

Systematic analysis of transcribed loci in ENCODE regions using RACE sequencing reveals extensive transcription in the human genome

Source: Wu , Du, et al. (2007) Genome Biology



Biological complexity revealed by ENCODE: Long Interleaved Transcripts and Distributed Regulation

What is a Gene? and What is not a <u>Gene?</u>

> [Gerstein et al. Genome Res. 2007; 17: 669-681]

<u>Proposed Re-definition of a Gene: "Gene is a union of genomic sequences</u> <u>encoding a coherent set of potentially overlapping functional products."</u>



100 Lectures.GersteinLab.org (c) 2007



101 Lectures.GersteinLab.org (c) 2007

Overview of the Process of Intergenic Annotation

Basic Inputs

- 1. Doing large-scale similarity comparison, looking for repeated or deleted regions
- 2. Determining experimental signals for biochemical activity (e.g. transcription) across each base of genome
- Results of Analyzing Similarity
 Comparison
 - 1. Finding large repeated or deleted blocks (e.g. CNVs) as a function of degree of similarity
 - 1. within reference human genome
 - 2. within human population
 - 3. between related organisms (e.g. mouse)
 - 2. Finding smaller "exon-level" similarities (e.g. pseudogenes)

- Results of Processing Raw Expt. Signals
 - 1. Signal Processing: removing artifacts, normalizing, window averaging
 - Segmenting signal into larger "hits" ("Active Regions" or ARs)
 - 3. Clustering together active regions into even larger features at different length scales and classifying them
 - 4. Building networks and beyond....

Outline

- Calling Blocks
 - Issues in Signal
 Processing Transcriptional
 Data from Tiling Arrays
- Clustering Blocks
 - Scoring transcribed regions and clustering these together into larger blocks
- Formal Annotation
 - ◊ Pseudogenes

- Integration with conservation
 A SDa and CNIV/a
 - ♦ SDs and CNVs
 - pseudogenes and seq.
 constraint
- Integration with activity
 - pseudogenes and transcriptions
- Future of Gene Annotation
 - ◊ What is a gene post encode?

Processing the Raw Experimental Signal

- Normalizing and Correcting Artifacts
 - ◊ Characterizing, correcting and exploiting specific and non-specific cross-hybridization
 - Measuring the effect of types of mismatches
 - Iterated quantile norm. to correct non-specific cross
 - Towards a universal array based on specific cross-hyb.

Annotating the Human Genome: <u>First-Pass Annotation Clustering and</u> <u>Characterizing Novel Transcribed Regions</u> <u>and Groups of Binding Sites</u>

- DART classification of TARs
 - ◊ 1300 TARs in ~200 novel ENCODE loci
 - based on expression and phylogenetic clustering

Annotating the Human Genome: Integrative Annotation of Pseudogenes in Relation to Conservation, Transcription, and Duplication

- Pseudogene Assignment Technology
 - ◊ Pipeline + DB
 - Pseudofam analysis of Pseudogene Families, highlight outliers
 - ◊ Ontology
- Annotation of Human Genome
 - Pipeline draft (20K) + Hybrid
 Approach
 - Pilot Phase: Consensus annotation from automatic pipelines & manual curation gives 201

- Integration with Conservation and Seq. Constraint
 - ~2/3 processed are primate specific
 - Evidence for selection operating on a few but most neutral
- Pseudogene Activity
 - >20% appear to be transcribed (38/201)
 - No obvious selection on transcribed ones

Analysis of Duplication in the Genome: SVs and SDs

- Large-scale analysis of existing CNVs & SDs in human genome
 - ◊ Process giving rise to pseudogenes
- SDs assoc. with Alu, pseudogenes and older SDs
- CNVs assoc. other repeats (microsat.) and not as much with SDs
- Suggestion: Alu burst 40 MYA triggered much NAHR rearrangement, then dupl. feed on itself in hotspots but now dying down and NAHR assoc. with other repeats and CNVs also from NHEJ
 - ◊ Role for proc. pseudogenes in giving rise to duplicated ones

ENCODE Acknowledgements

<u>Adam Frankish, Robert Baertsch,</u> Philipp Kapranov, Alexandre Reymond, <u>Siew Woh Choo,</u> Y Fu, <u>Yontao Lu</u>, France Denoeud, Stylianos Antonarakis, <u>Yijun Ruan, Chia-Lin Wei</u>, Z Weng, Thomas Gingeras, Roderic Guigo, <u>Tim Hubbard, Jennifer Harrow</u>

Sanger, UCSC, GIS, AFFX, Geneva, IMIM, BU + SU

+

various consortia (ENCODE, modENCODE, 1000 Genomes)

Acknowledgements

pseudogene.org, tiling.gersteinlab.org, sv.gersteinlab.org



Acknowledgements

pseudogene.org, tiling.gersteinlab.org, sv.gersteinlab.org



Acknowledgements

pseudogene.org, tiling.gersteinlab.org, sv.gersteinlab.org





Permissions Statement

This Presentation is copyright Mark Gerstein, Yale University, 2007.

Feel free to use images in it with **PROPER acknowledgement**

(via citation to relevant papers or link to gersteinlab.org).