

Annotation Non-coding Regions of the Human Genome

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

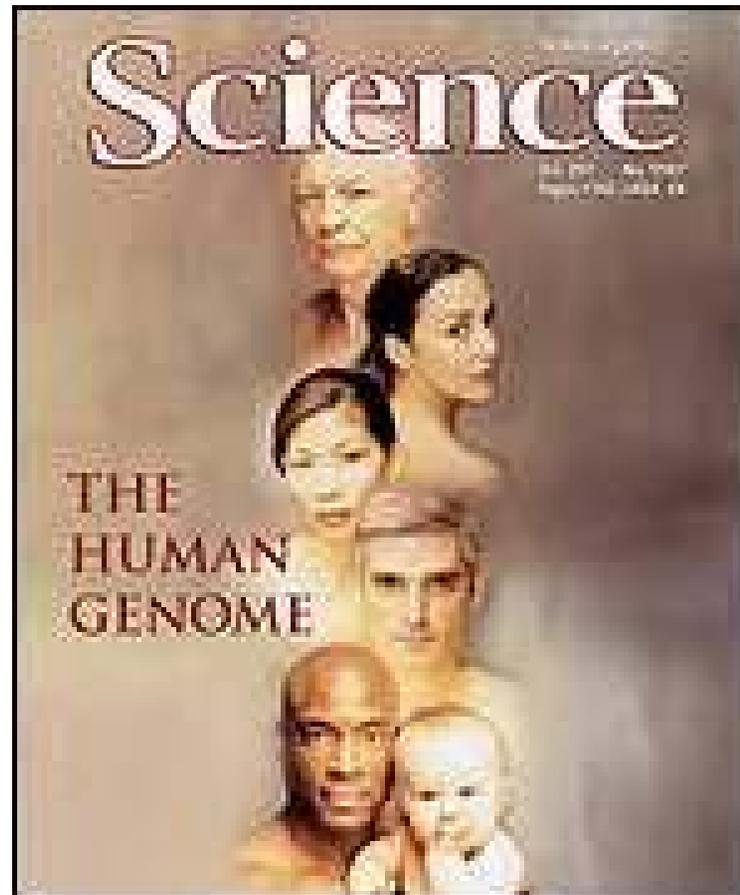
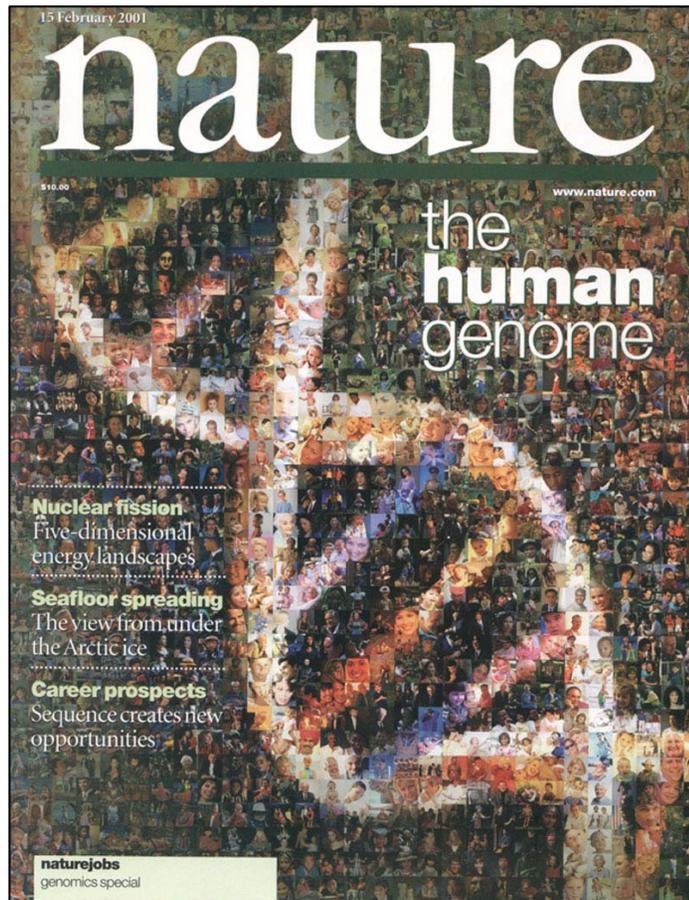
Bristol Myers

2008.08.15, 14:20-14:45

Slides downloadable from Lectures.GersteinLab.org.

(Please read permissions statement.)

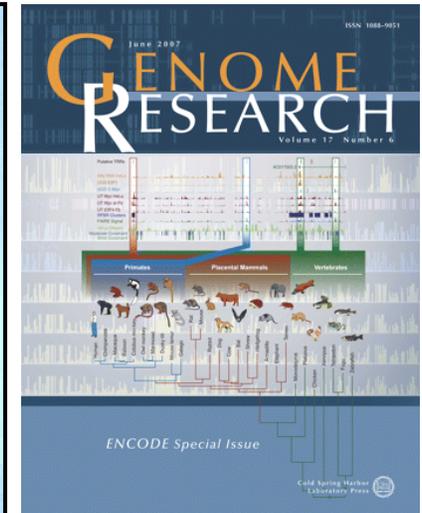
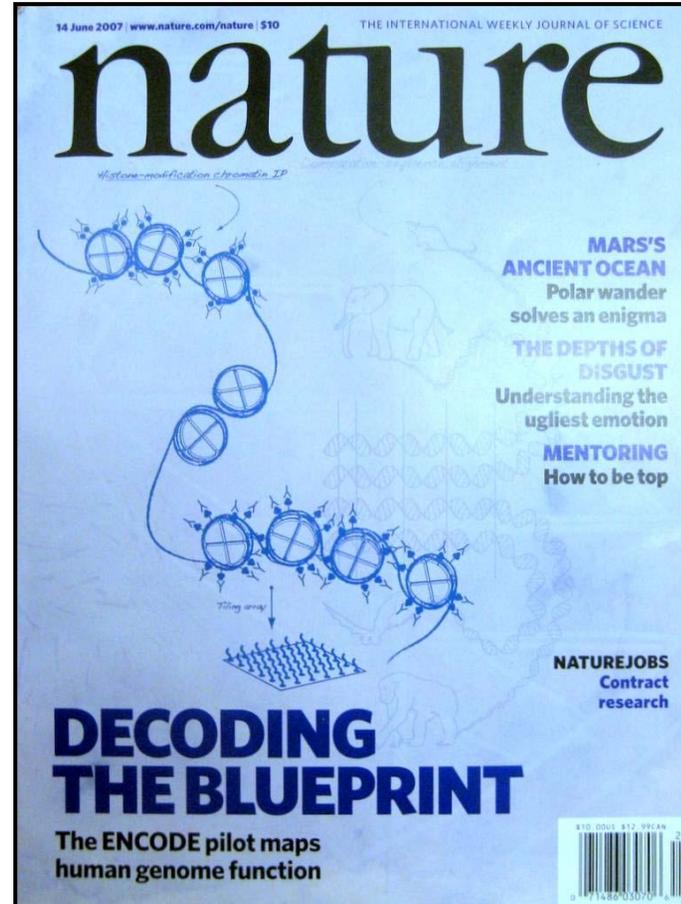
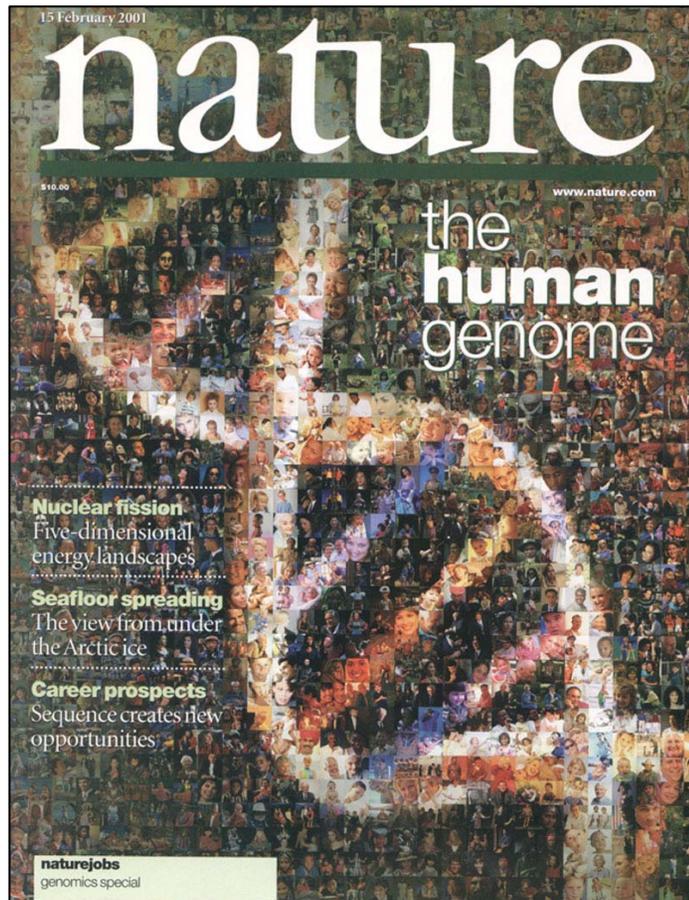
Paper references mostly from Papers.GersteinLab.org.



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?

[IHGSC, *Nature* 409, 2001]

[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]

How might we annotate a human text

?

Color is
Function

Lines are
Similarity

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

The Semicolon Wars

Brian Hayes

IF 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 Ethnologue.com.

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,

decide 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 least-significant 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, I 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 Annotation Process

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

- ◇ Development of Sequence (and Array) Technology

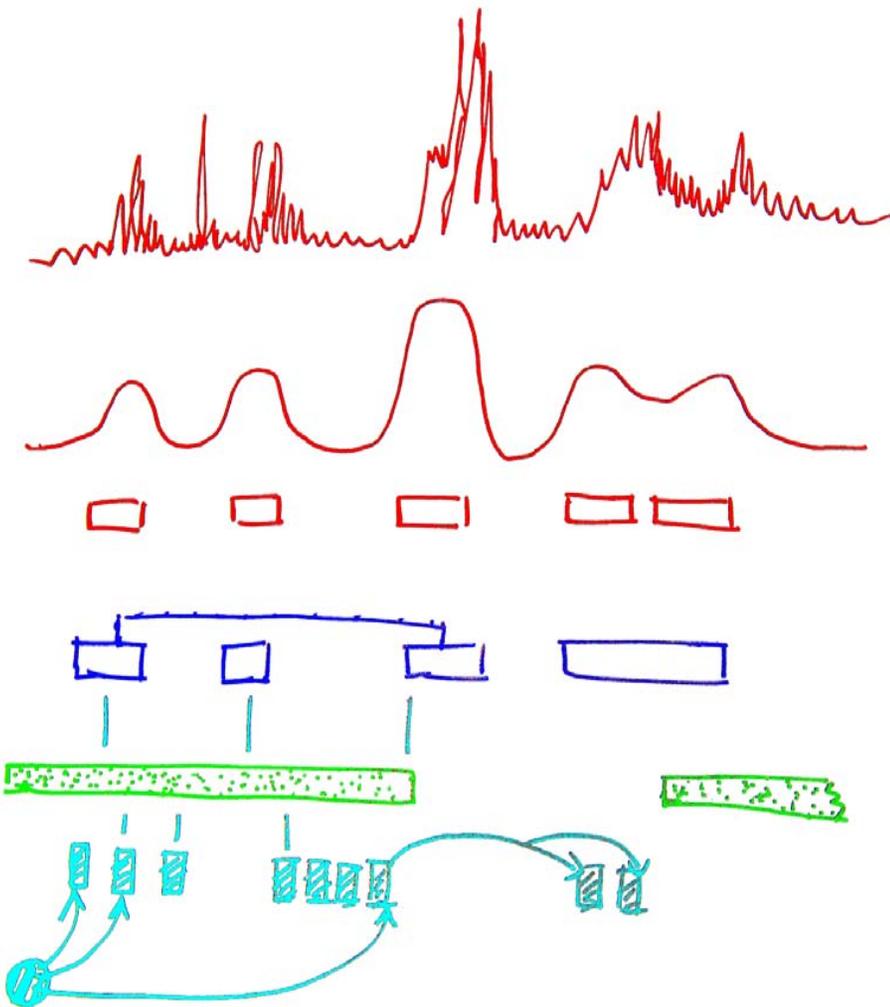
- Normalizing & Scoring Signal, Correcting Artifacts, Segmenting to create Small Annotation Blocks

- ◇ 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



ENCODE

&

mod

ENCODE

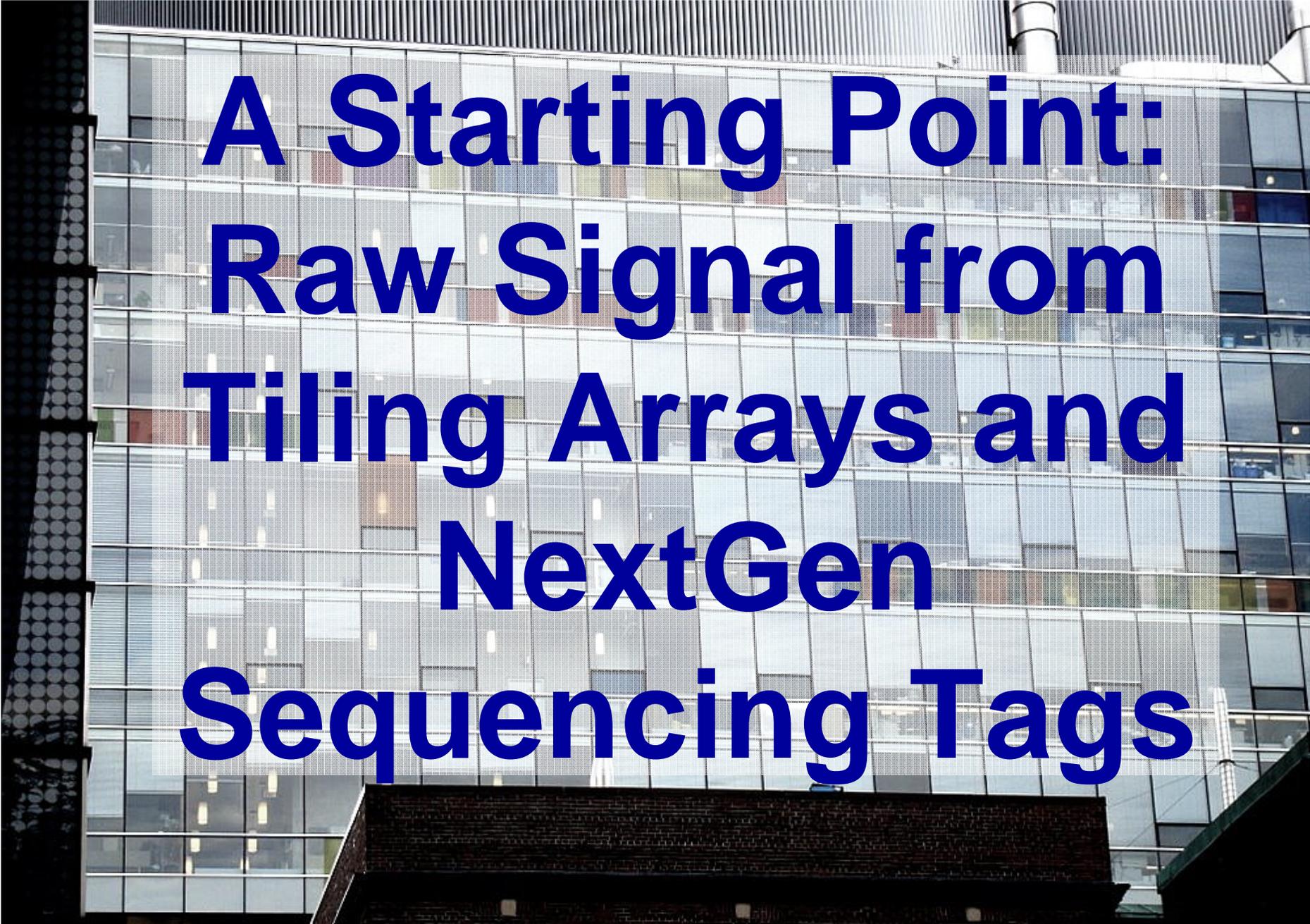
@ Yale

+

CEGS

Yale Center for Excellence
in Genome Sciences

- Array and NextGen Seq. Experiments
 - Mike Snyder & Sherman Weissman
 - Interrogation of small fragments of chromosomes to determine their function
 - Large-scale hybridization to find transcribed regions in unbiased fashion
 - TF binding sites (via ChIP-chip)
 - CNVs and SDs (from hires-aCGH)
- 1st Pass Computational Annotation
 - Classification of Novel Transcribed Regions
 - Grouping and Classification of Binding Sites
 - Characterization of SDs and CNVs
- Integrative Annotation
 - Pseudogenes (Zheng et al., GR,GB)
 - Inter-relation with Transcription & CNVs



A Starting Point: Raw Signal from Tiling Arrays and NextGen Sequencing Tags

High-Resolution CGH with Oligonucleotide Tiling Microarrays

Maskless Array Synthesis

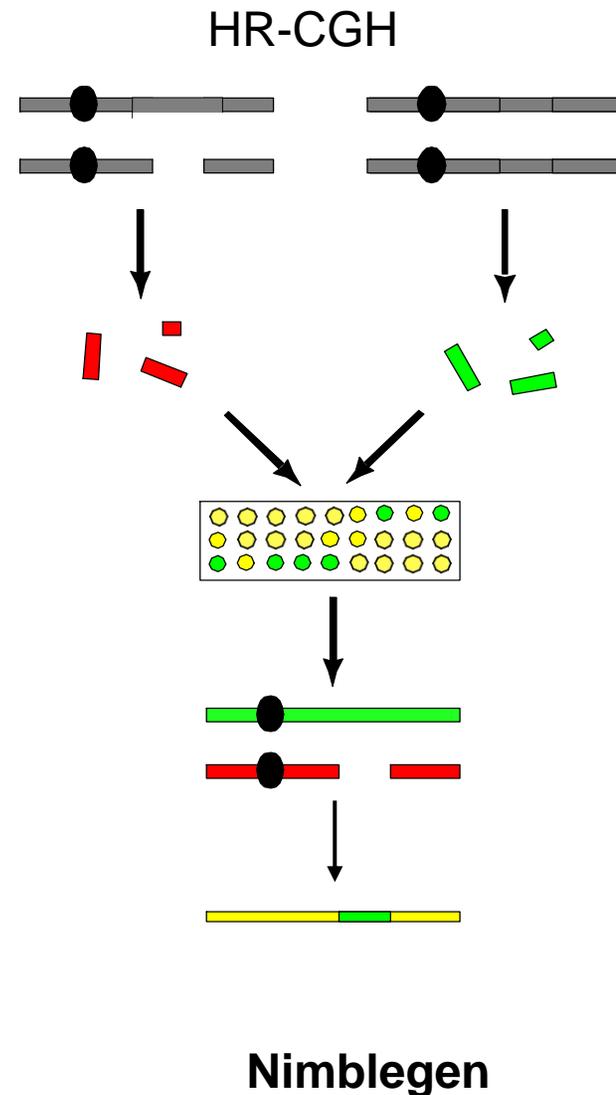
385,000 oligomers/chip

Isothermal oligomers,
45-85 bp

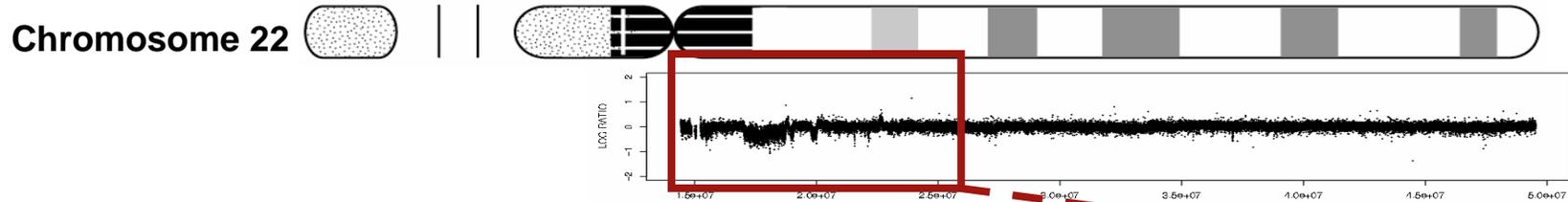
Tiling at ~1/100bp non-
repetitive genomic sequence

Detects CNVs at 1 kb
resolution

Urban et al., 2006



Representative Signal from aCGH with CNVs & Breakpoints



High Resolution Array
Comparative Genomic
Hybridization (aCGH)

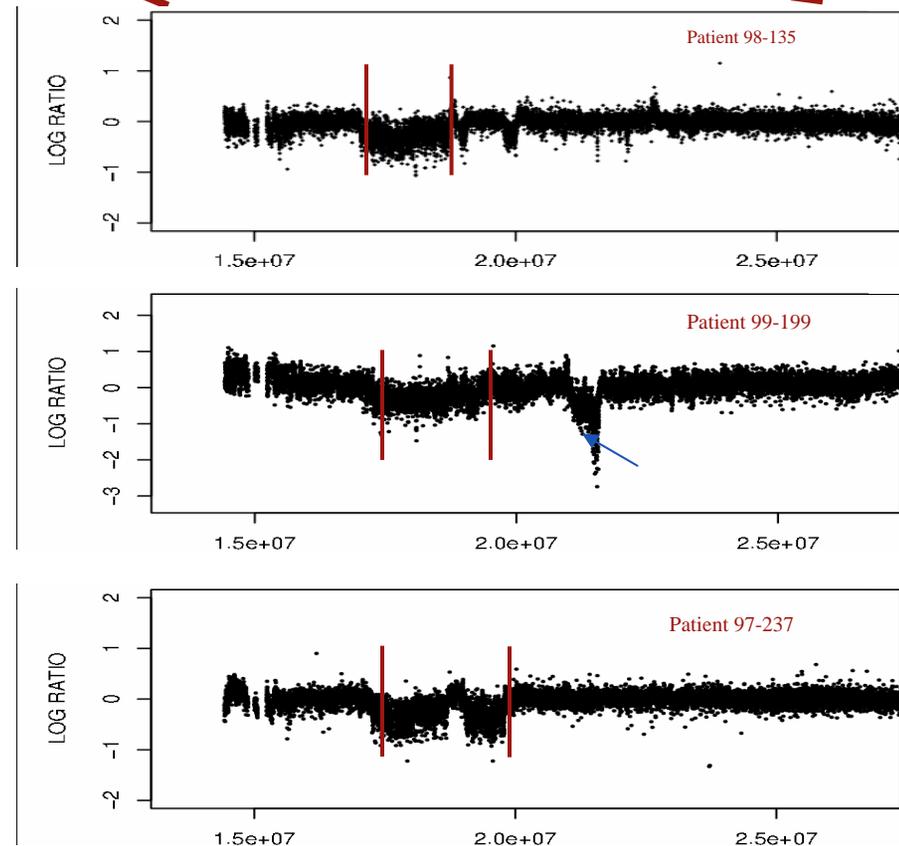
Calling Copy Number
Variants (CNVs) between
Breakpoints

Nimblegen/MAS Technology

Isothermal Arrays Covering
Chromosome 22

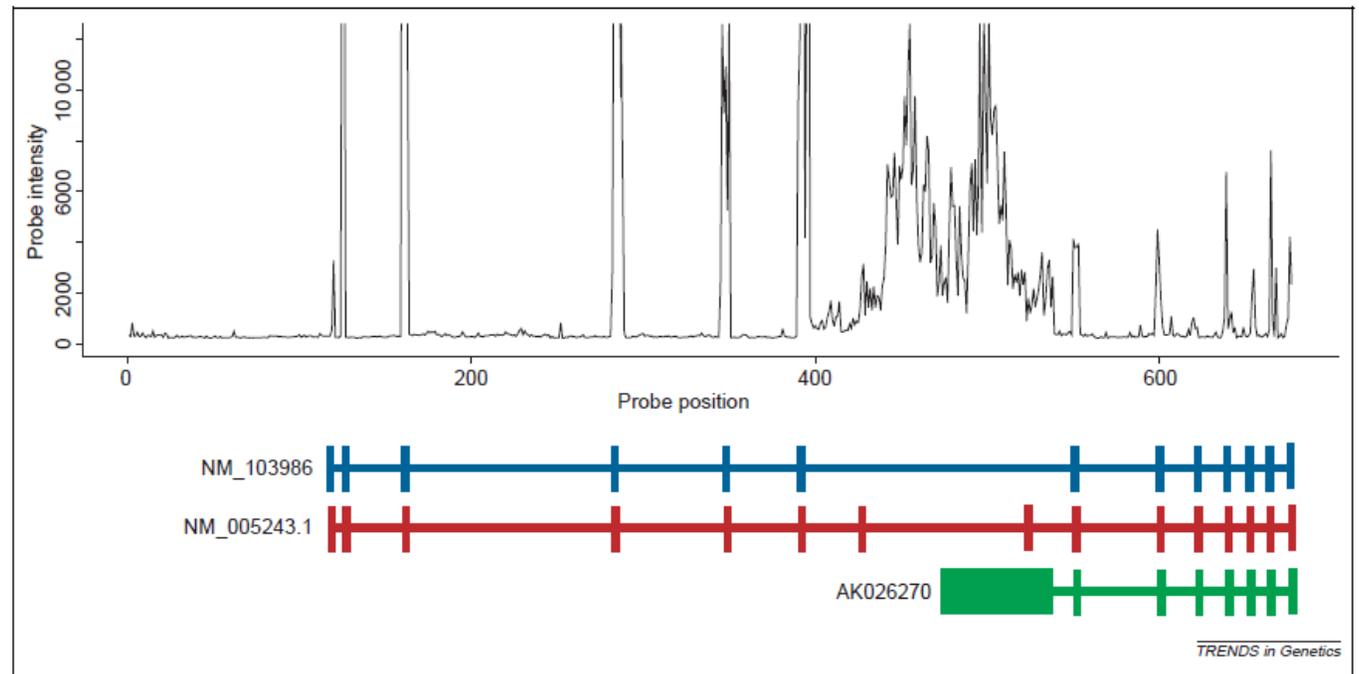
Resolution ~1 kb

Urban et al. (2006) PNAS

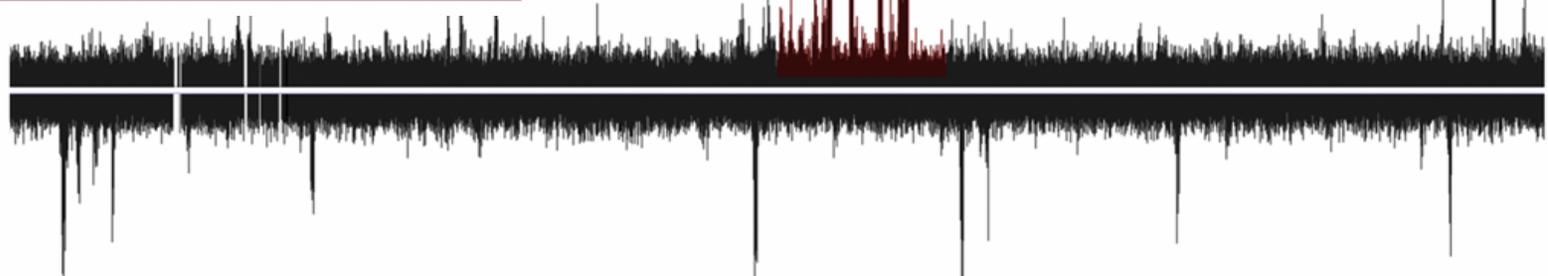
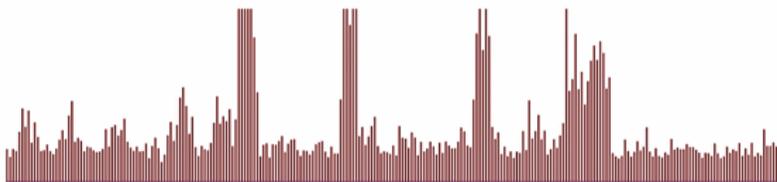


A Starting Point: Noisy Raw Signal from Tiling Arrays (Transcription)

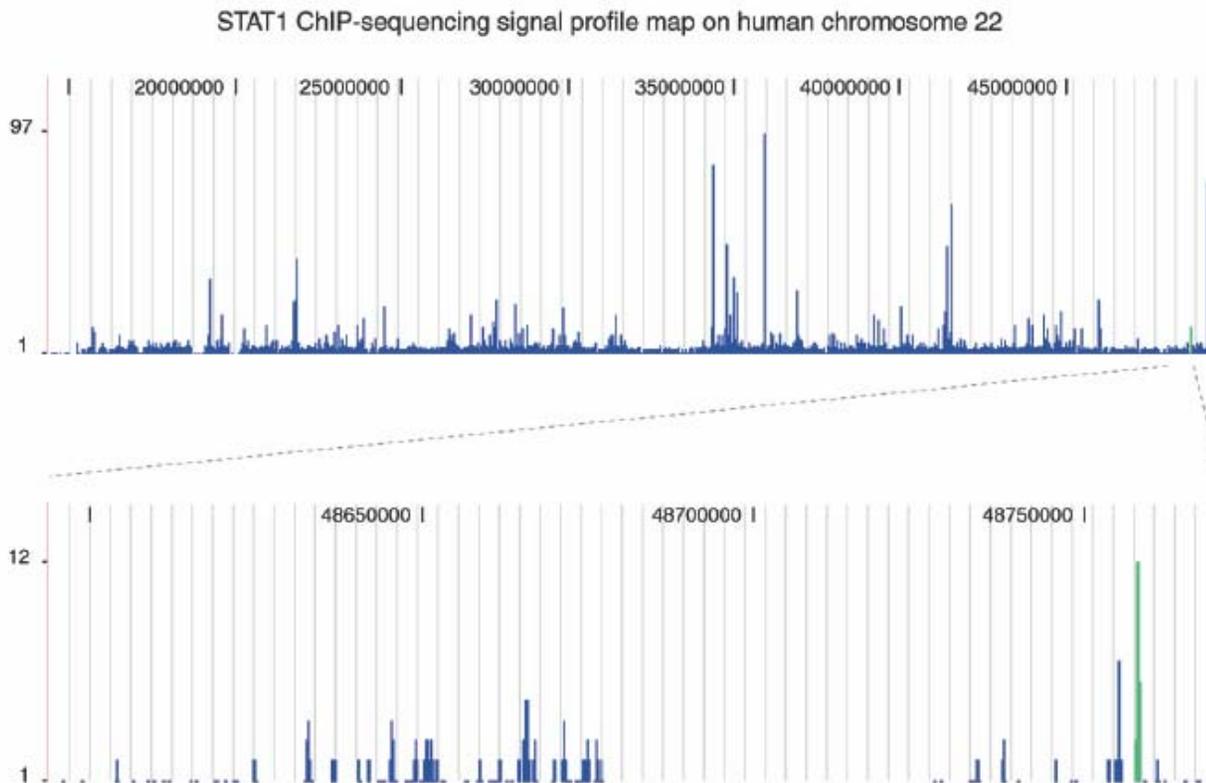
Johnson et al. (2005) TIG, 21, 93-102.



ta_15 ta_16 ta_17 ta_18

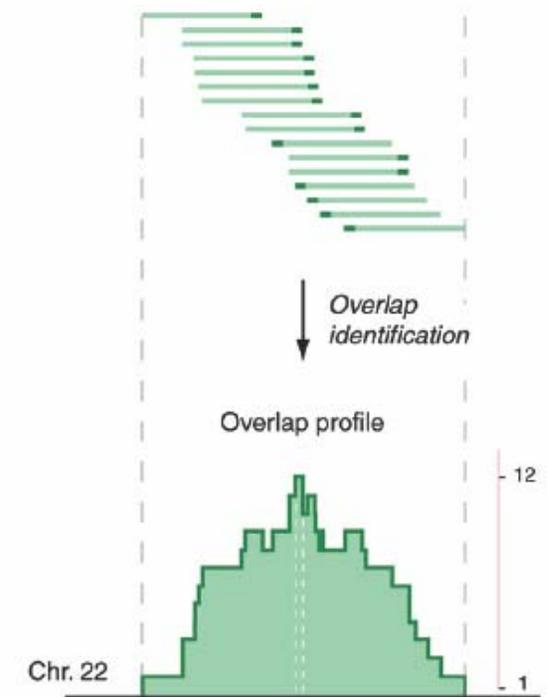


Representative Signal from Chip-Seq



C

16 uniquely mapped sequence reads and their directional extension in a tag cluster



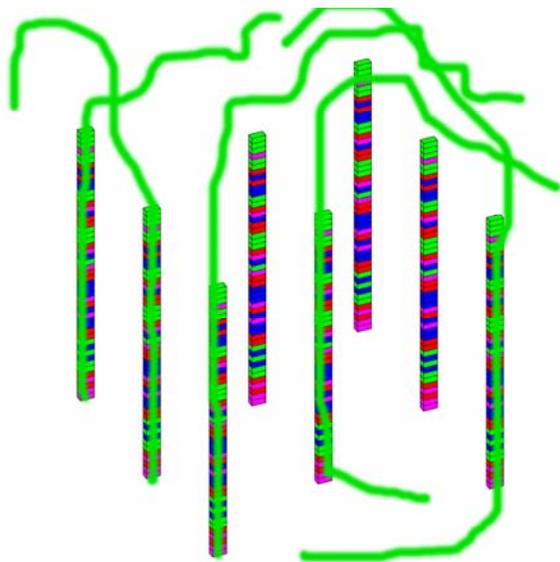
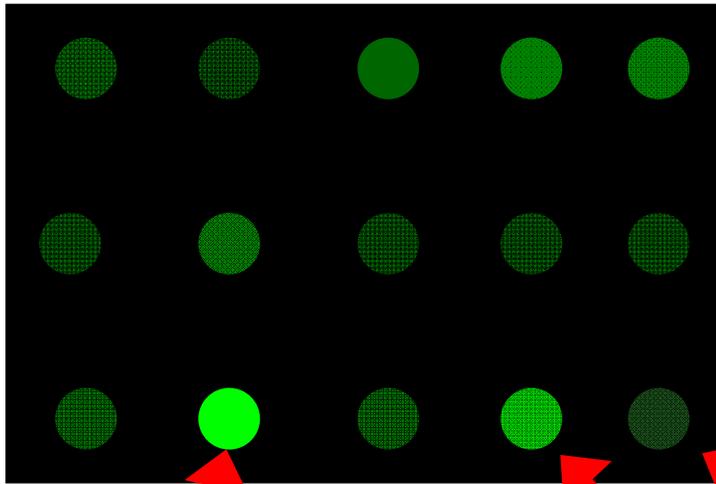
[Robertson et al., Nat. Meth. ('07); Zhang et al. PLOS Comp. Bio. (in revision, '08)]

Signal Processing: Normalizing, Measuring & Correcting for Aspects of Hybridization)

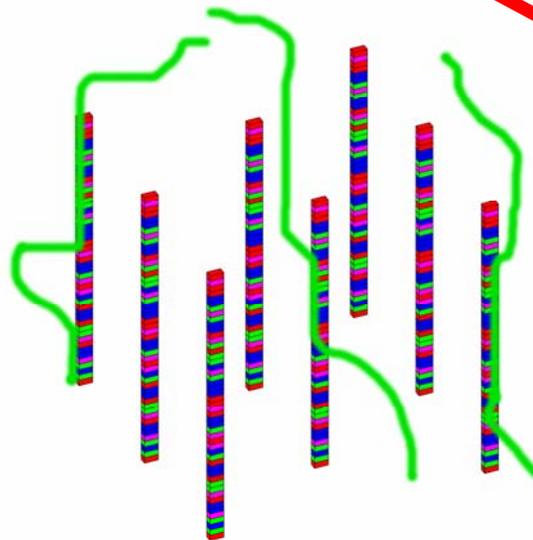


Specific & Non-specific Cross-Hyb.

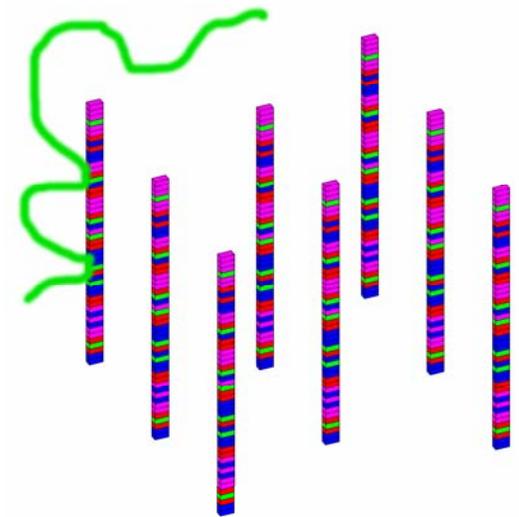
- Perfect match (PM): probe binding intended target
- Specific cross-hyb.: probes binding non-PM targets with a small number of mismatches
- Non-specific cross-hyb.: probes binding targets with many mismatches, due to general stickiness of oligos



Perfect Match



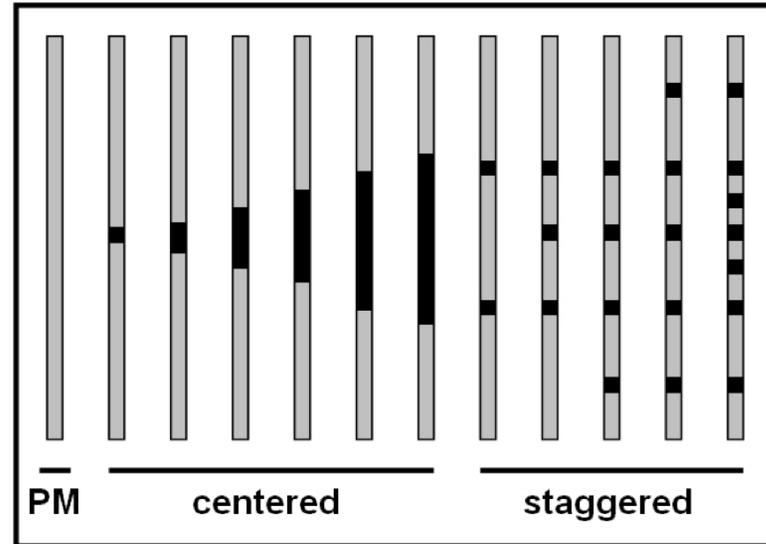
Specific Cross-hyb.



Non-specific Cross-hyb.

Creation of Standardized Datasets for Quantifying Effect of Mismatches

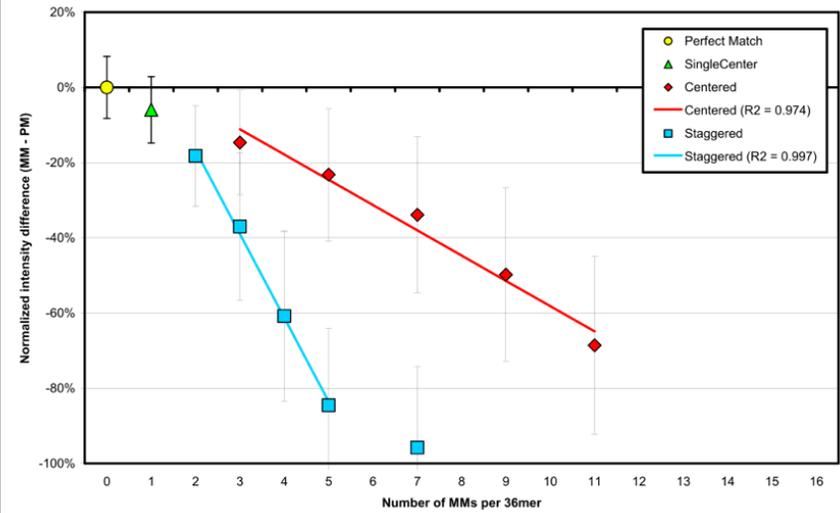
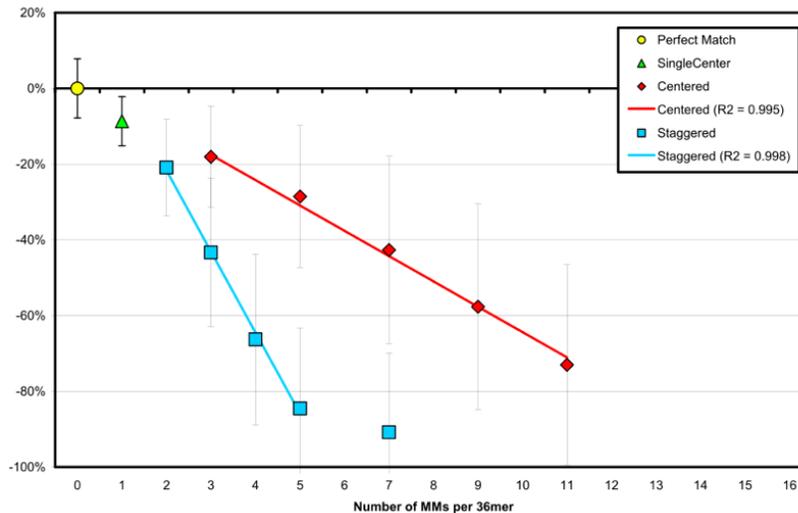
[Serinhaus et al., BMC Genomics (in press)]



Yeast ACT1 Gene

Human HBG2 Gene

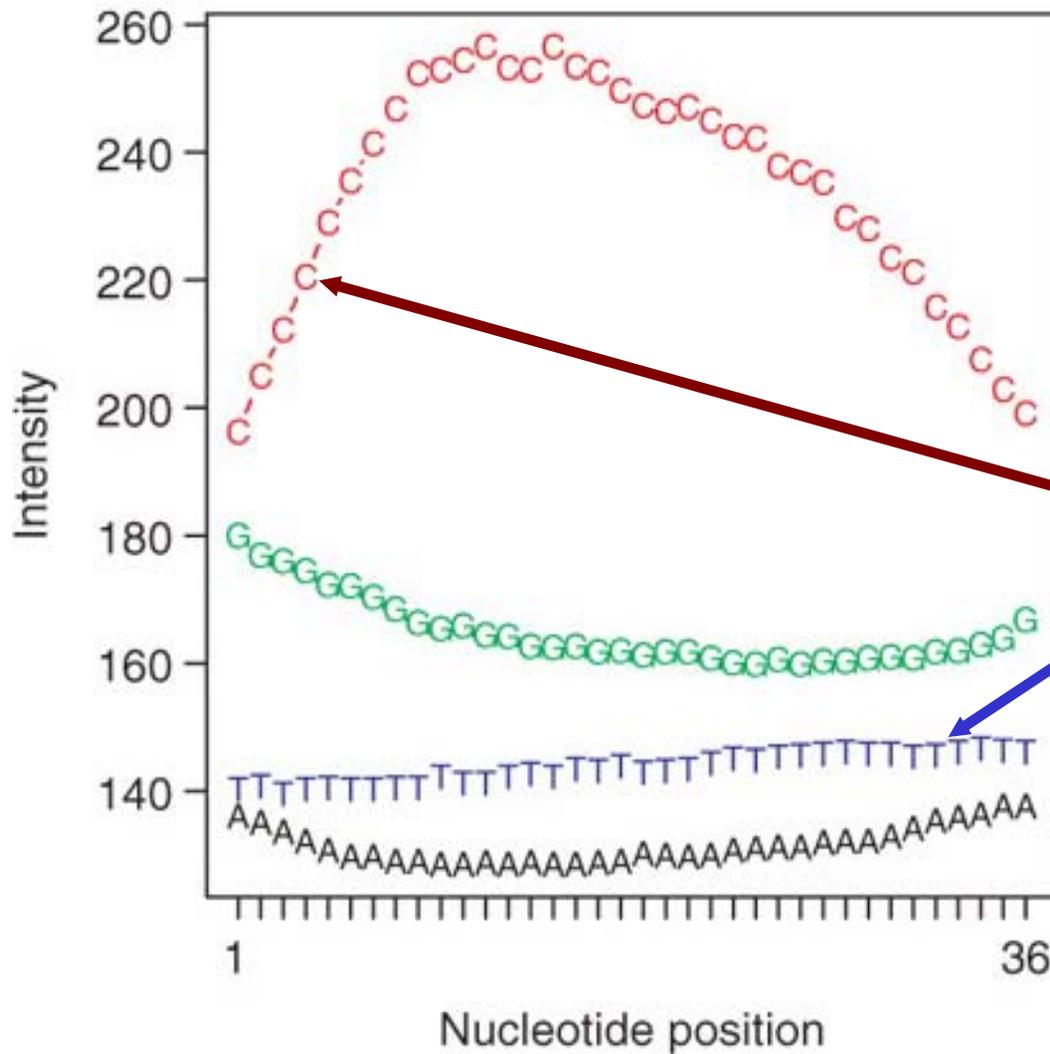
Normalized Intensity
Chg. vs. PM



Number of Mismatches

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

Nimblegen 50th Quantile



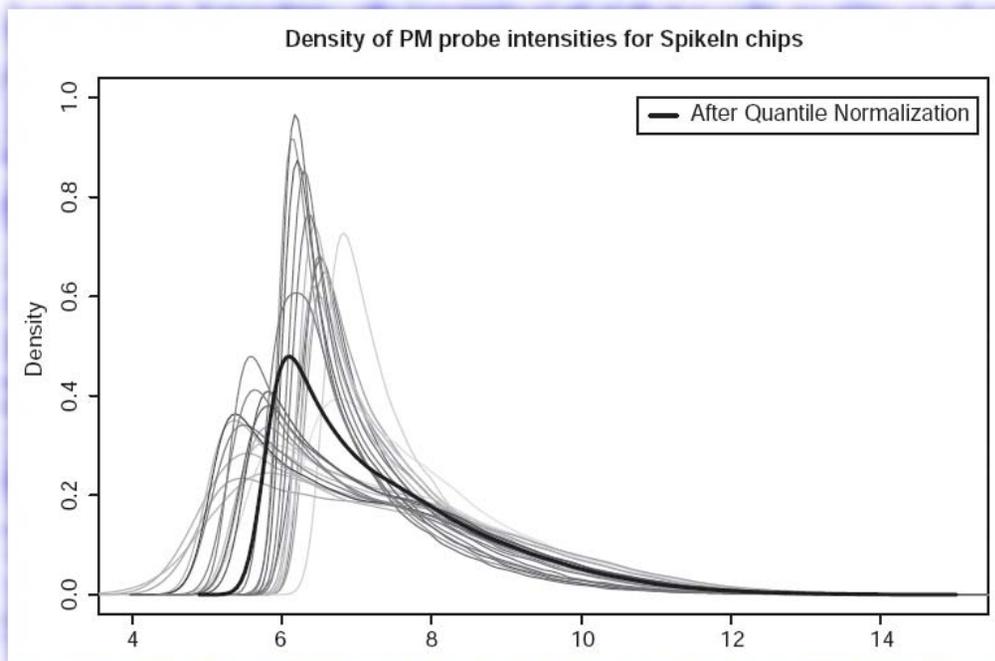
Avg. intensity of all background probes with a C at position 4

Avg. intensity of all background probes with a T at position 33

Quantile Normalization

Gene expression quantile normalization

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



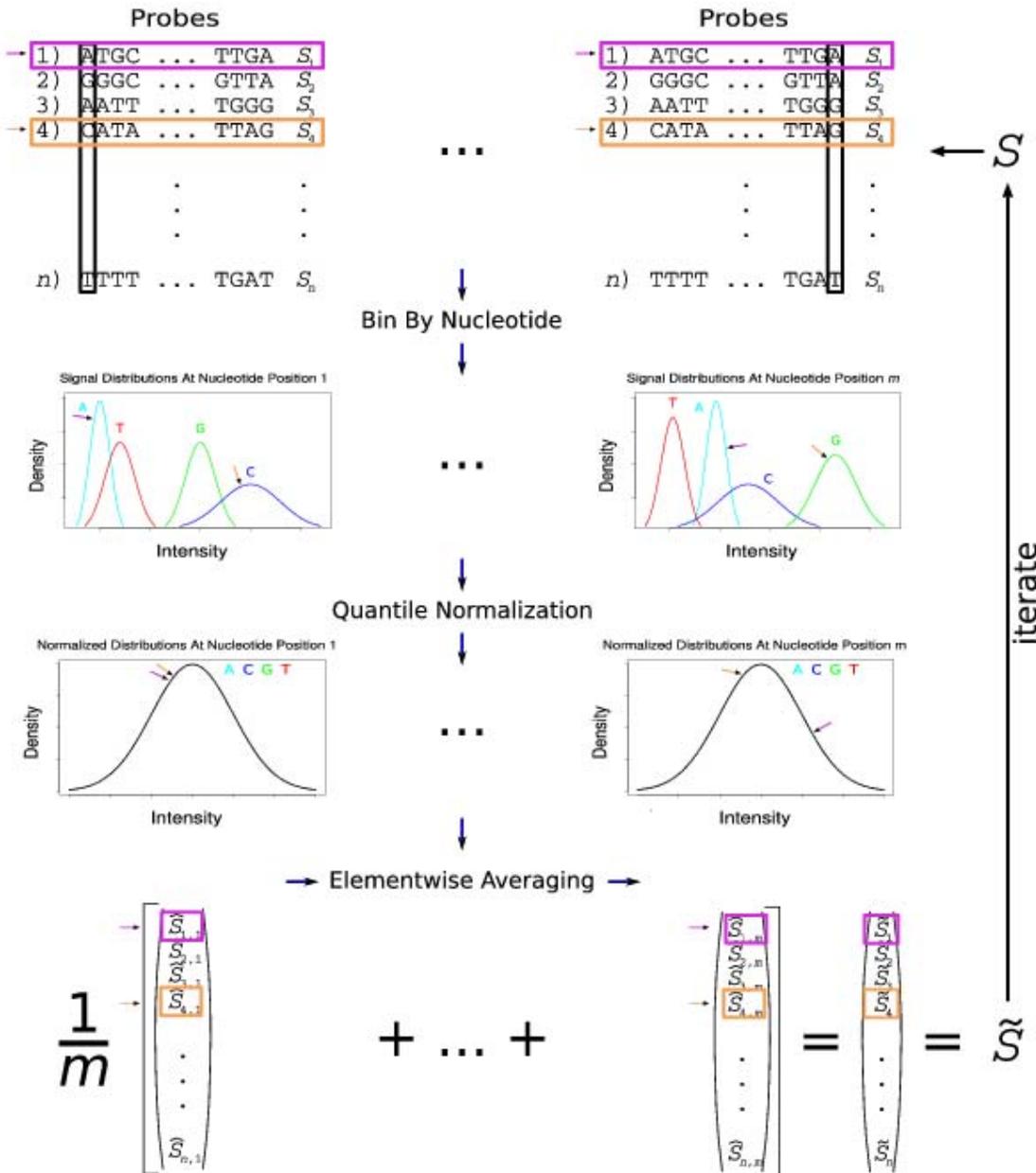
Source: Bolstad, B.M., et al (2003), *Bioinformatics*, **19**, 185-93.

Example

- Distributions that should be the same are forced to be the same (dark line).

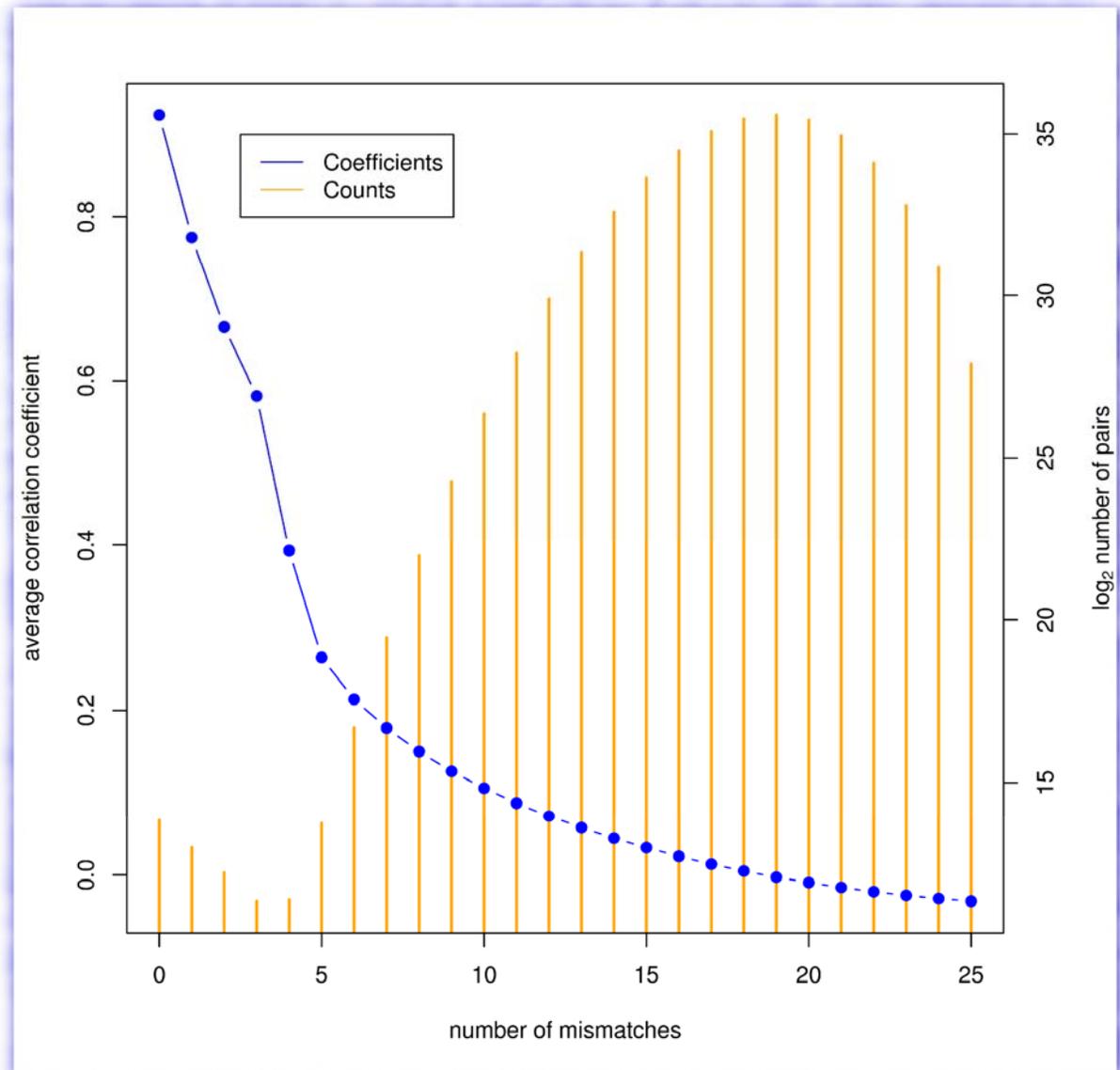
Iterated Quantile Normalization to Correct for Non-specific Cross-hyb.

- 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)



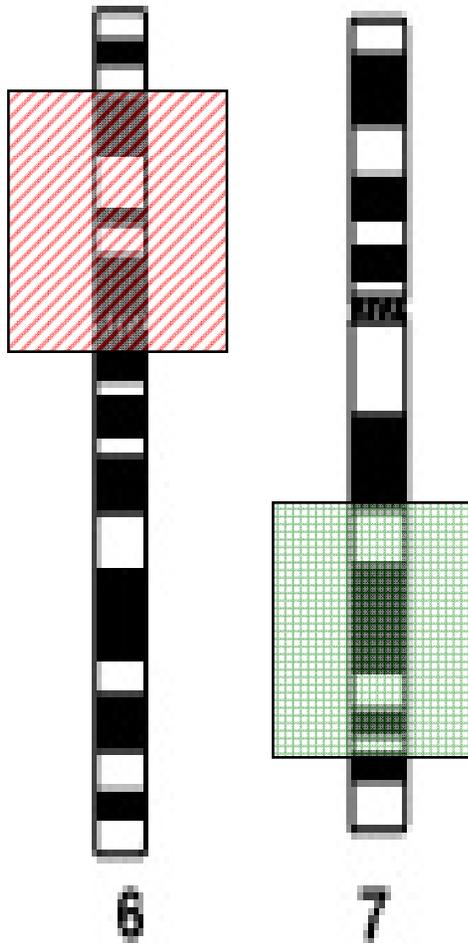
Measuring Specific Cross-Hyb

- 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**, 98-97

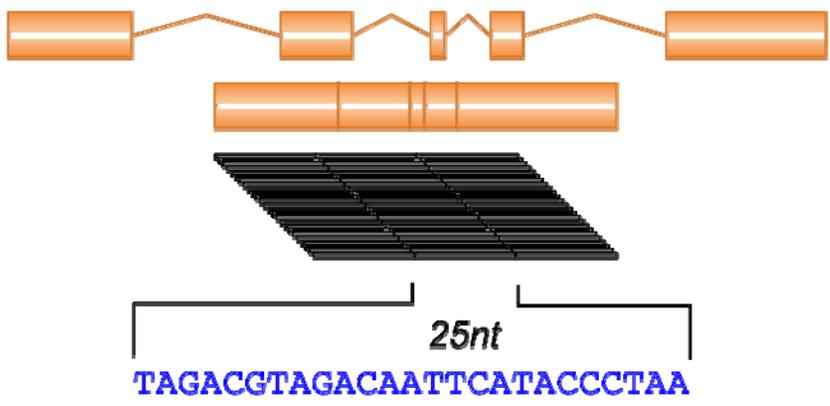
Proof of principle test to “exploit” this



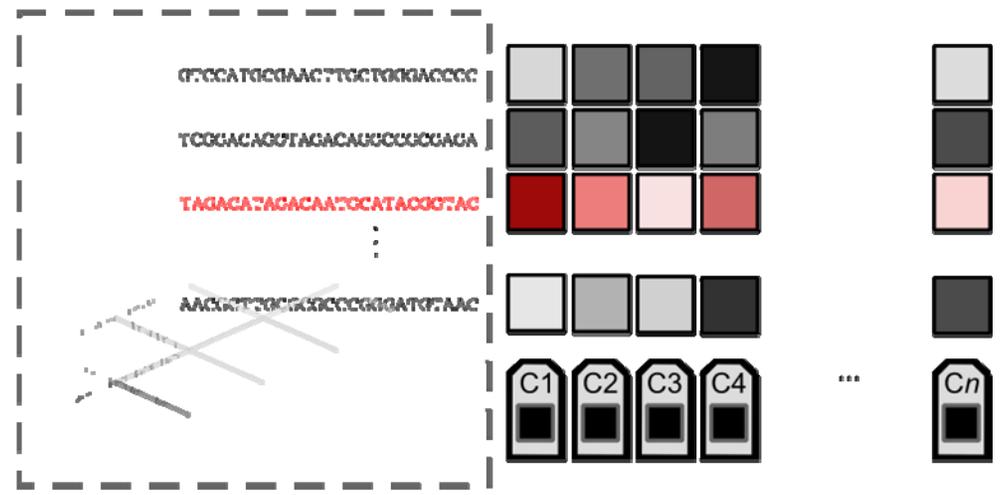
- Using Cheng et al. (2005), predict gene expression levels (and profiles across tissues) for genes on part of chr. #6
- ...Based on closest cross-hybrid tiles on part of chr. #7
- Then compare to measured levels and profile on #6

Nearest Nbr Search on Virtual Tiling

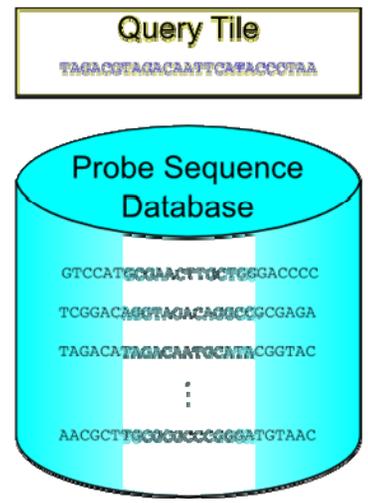
a virtual tiling



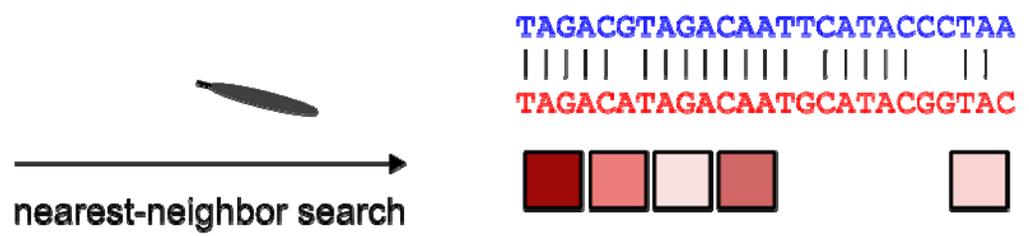
b microarray hybridizations



c similarity search

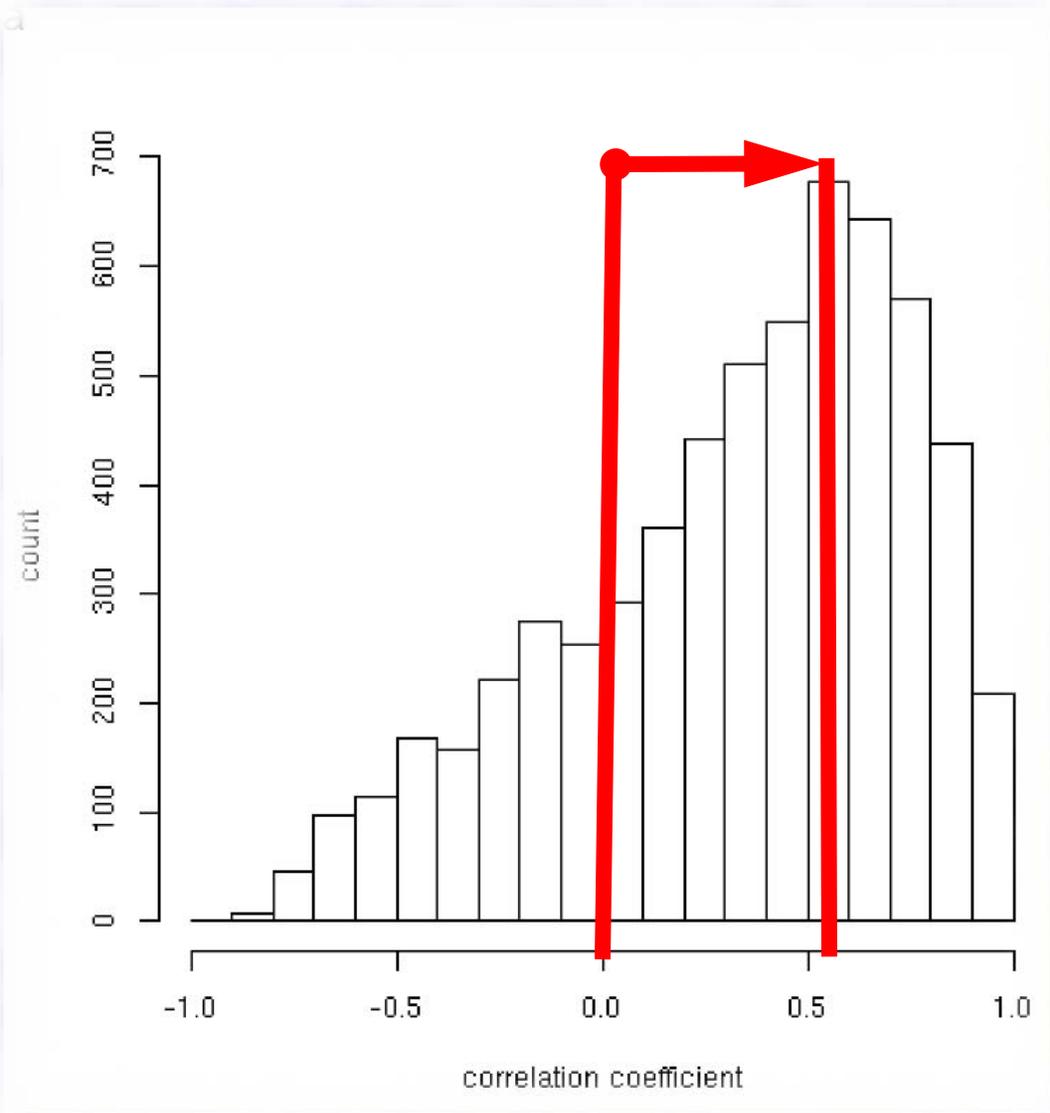


d profile assignment from nearest-neighbor



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)

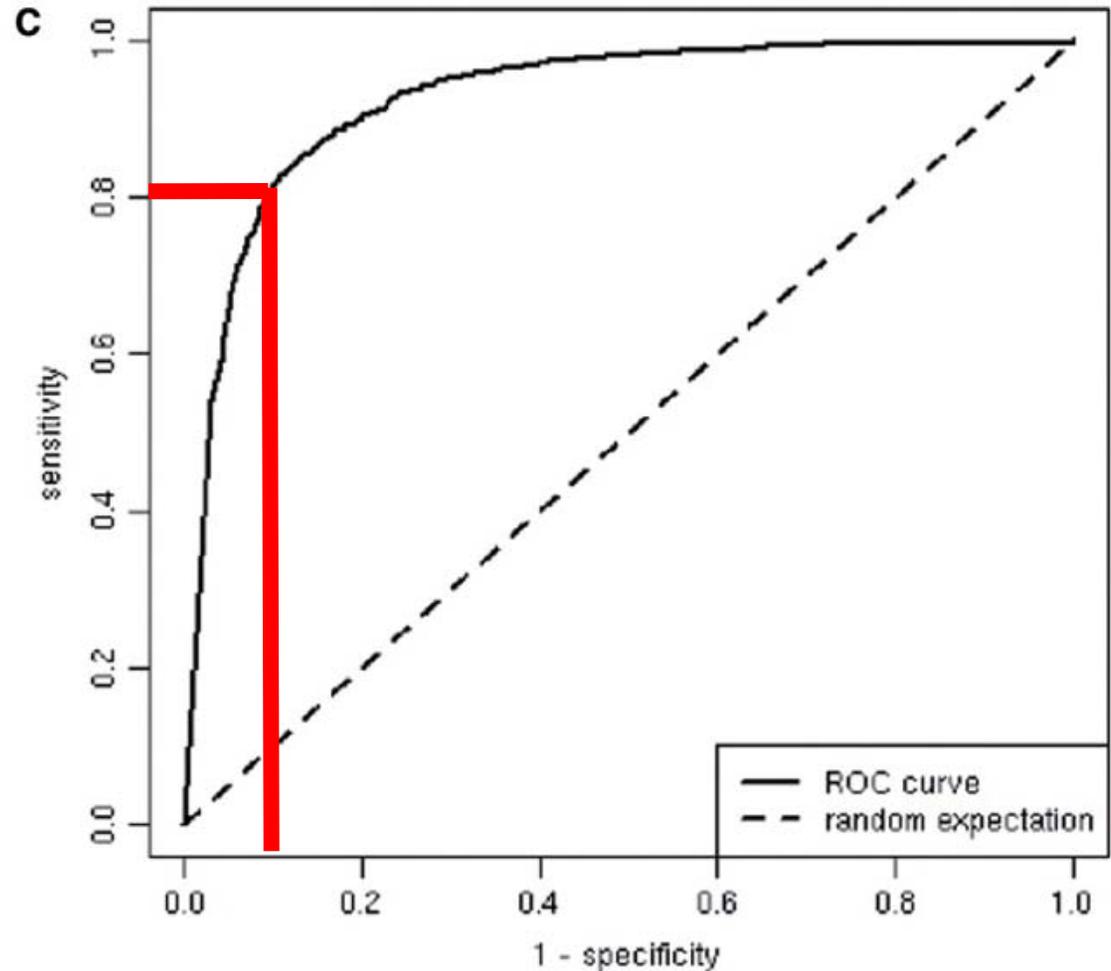


Source: Royce, T.E., et al (2007), *Nucleic Acids Res.*, **23**, 9

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.



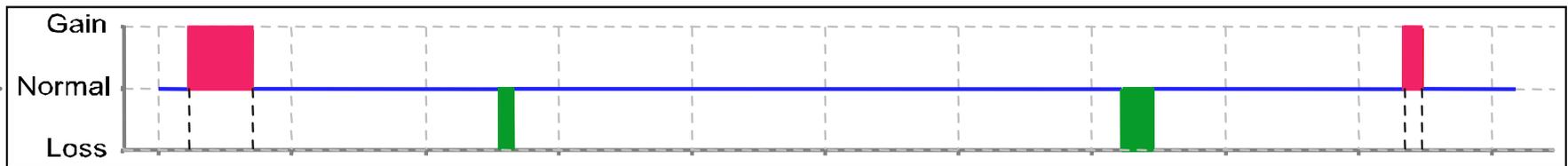
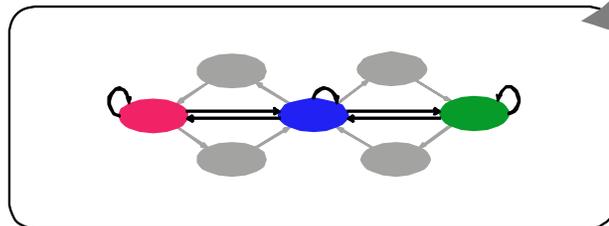
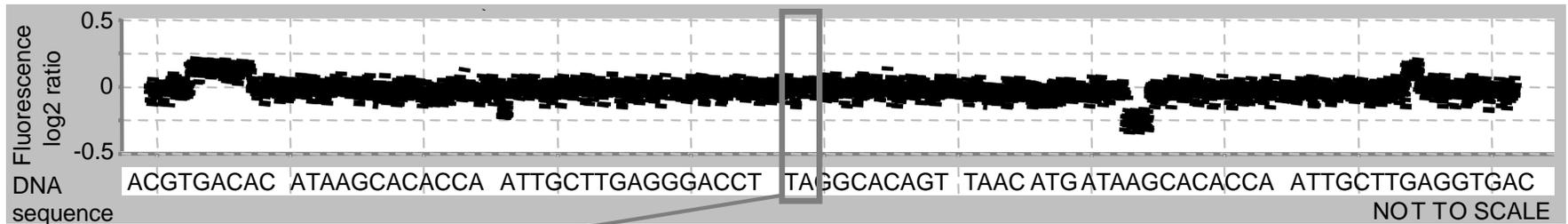
Royce, T. E. et al. Nucl. Acids Res. 2007 35:e99



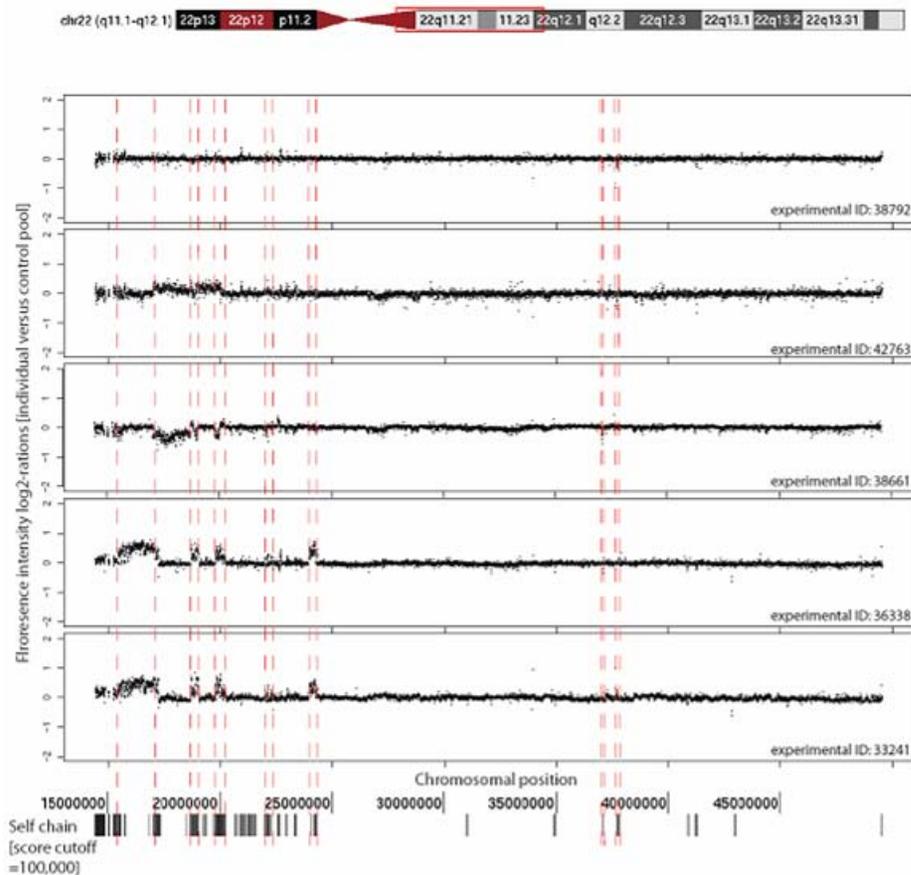
HMM Segmenters,
using "active learning" to find
Annotation Blocks from Raw
Signal

BreakPtr HMM

- To get highest resolution on breakpoints need to smooth & segment the signal
- BreakPtr: prediction of breakpoints, dosage and cross-hybridization using a system based on Hidden Markov Models

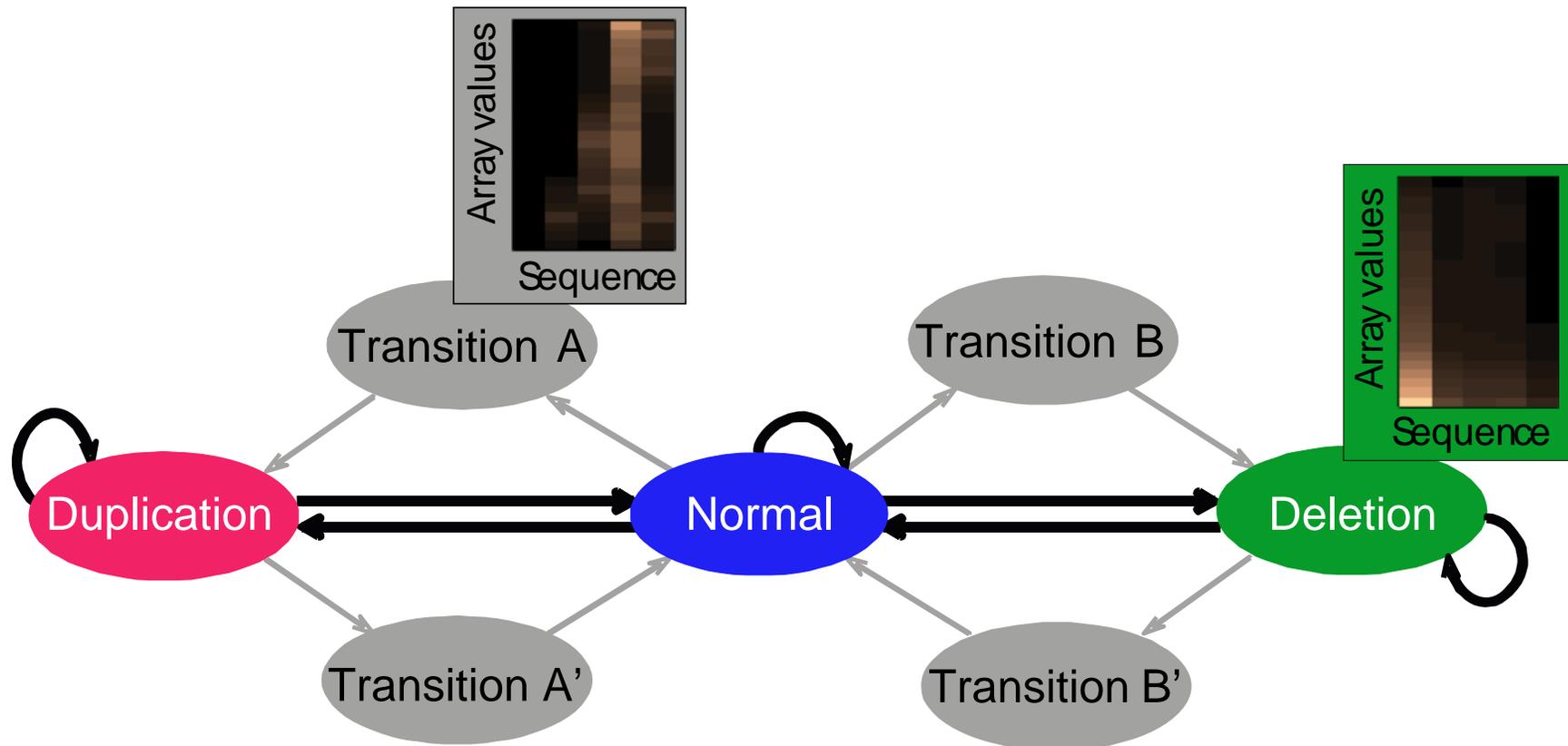


High resolution of tiling arrays allows statistical integration of nucleotide sequence patterns

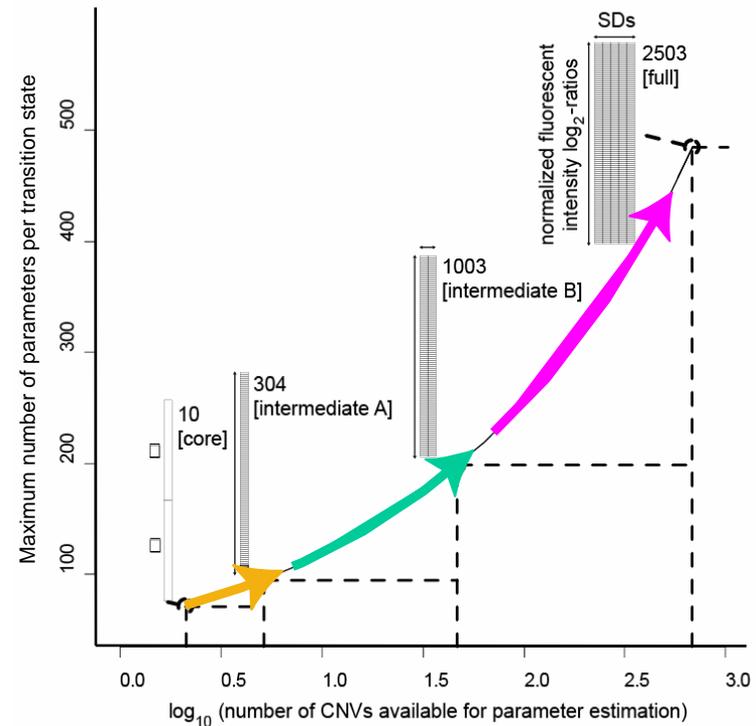
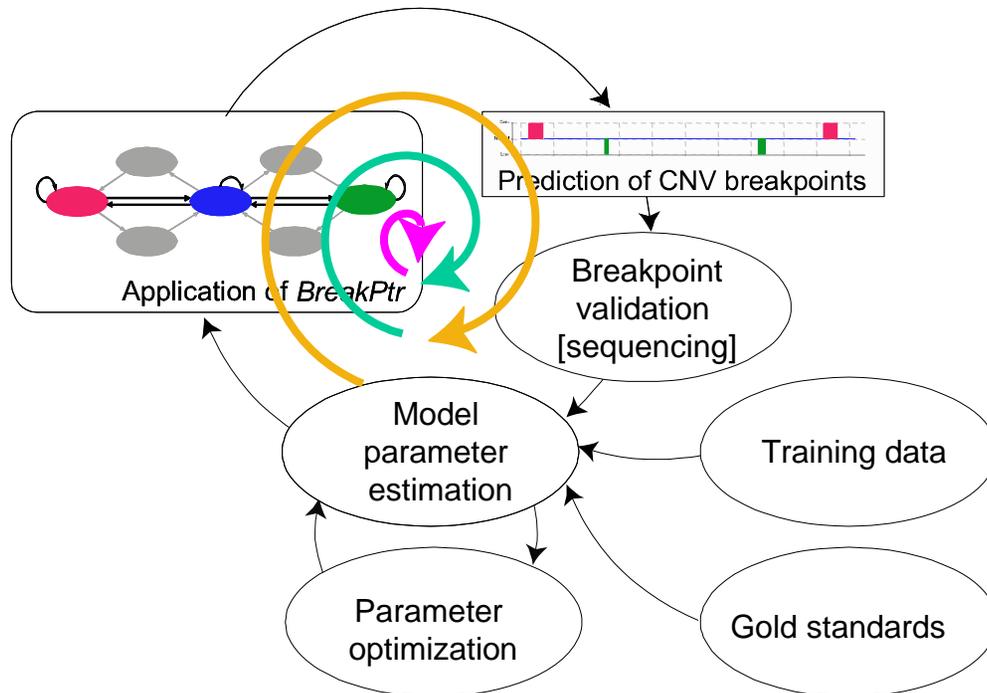


>4-fold enrichment of the breakpoints of copy number variants near segmental duplications (SDs) [e.g. Sharp *et al.*, *Am. J. Hum. Genet.* 2005; 77:78-88].

BreakPtr statistically integrates array signal and DNA sequence signatures
(using a discrete-valued bivariate HMM)



'Active' approach for breakpoint identification: initial scoring with preliminary model, targeted validation (with sequencing), retraining, and rescoring



CNV breakpoints sequenced in ~10 cases following BreakPtr analysis;

Median resolution <300 bp

No improvement in accuracy with higher resolution

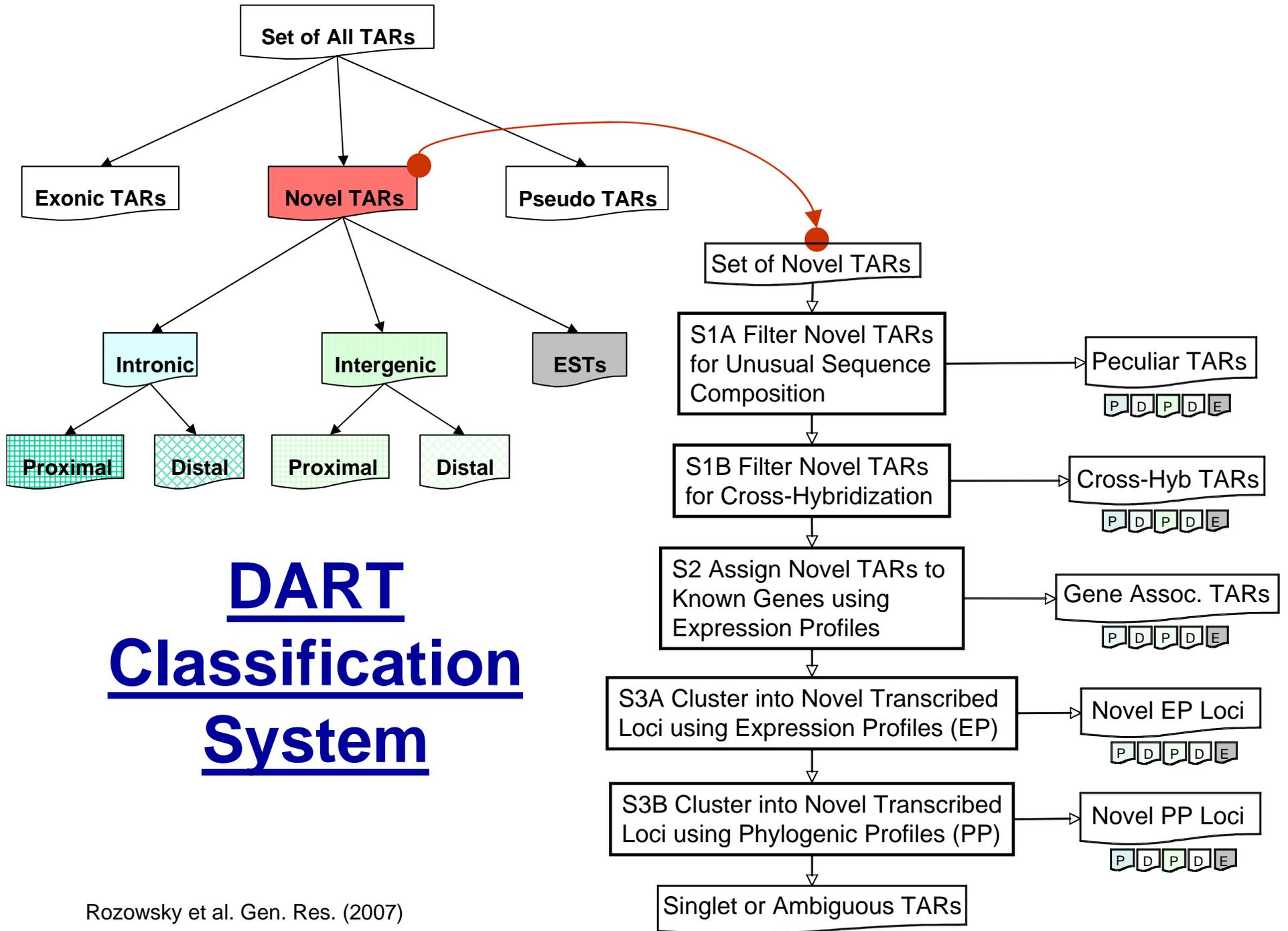
(9nt tiling)

HMM optimized iteratively
(using Expectation Maximization, EM)

Korbel*, Urban* *et al.*, PNAS (2007)

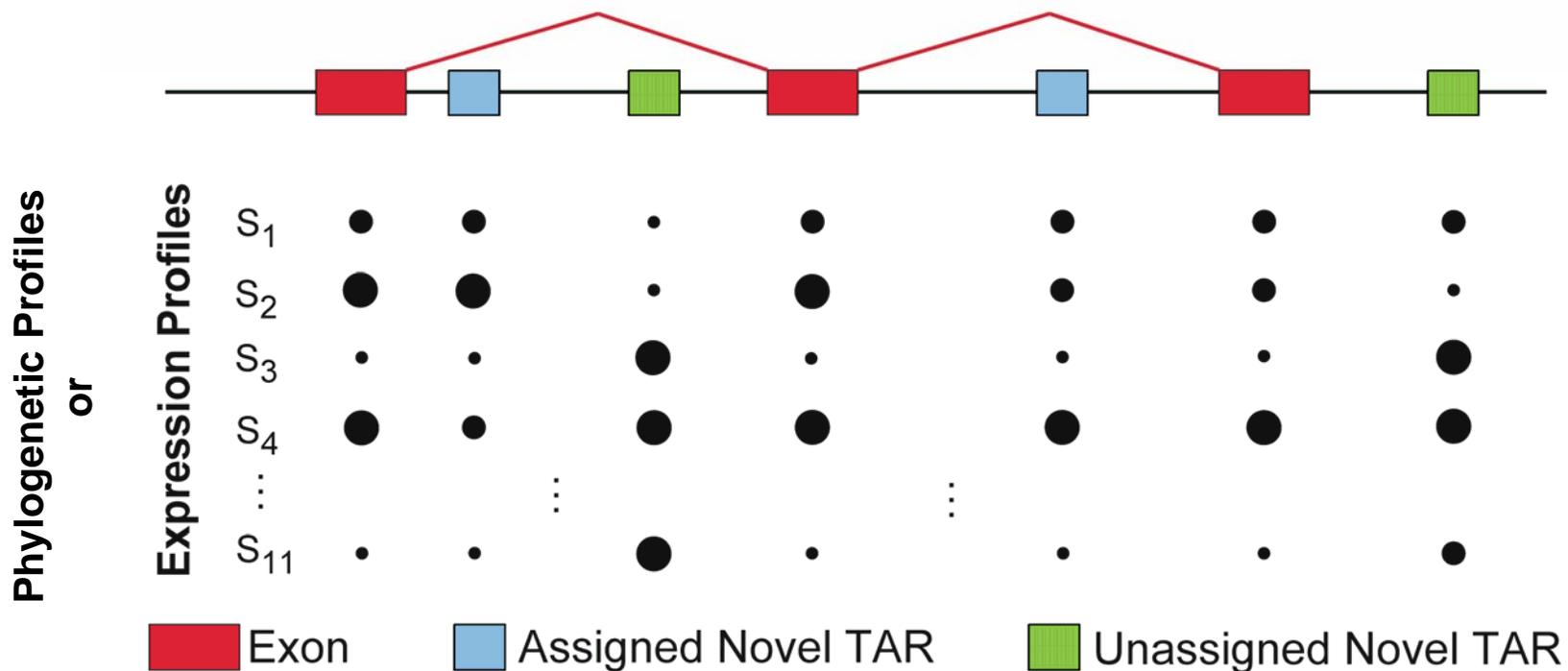


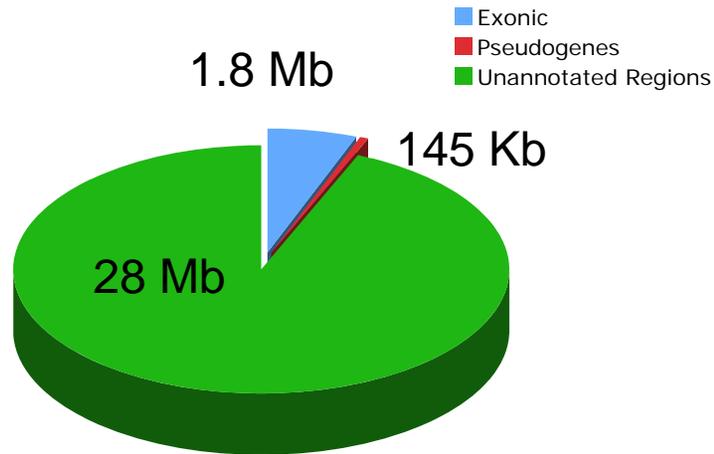
**Annotating a single type of
signal on a large-scale:
Clustering and Classifying Un-
annotated Transcription
(TARs)**



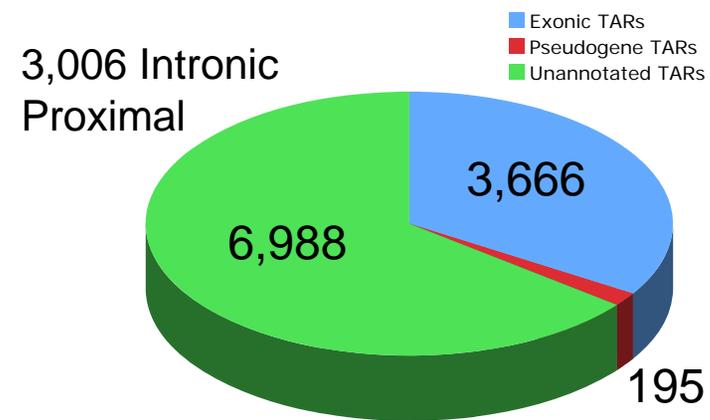
More Developed Annotation: Clustering and Classifying Blocks of Un-annotated Transcription into larger units

Assignment of novel TARs to known gene loci





ENCODE Regions (30 Mb)



Locations of TARs

Of the approx **7,000** Novel TARs

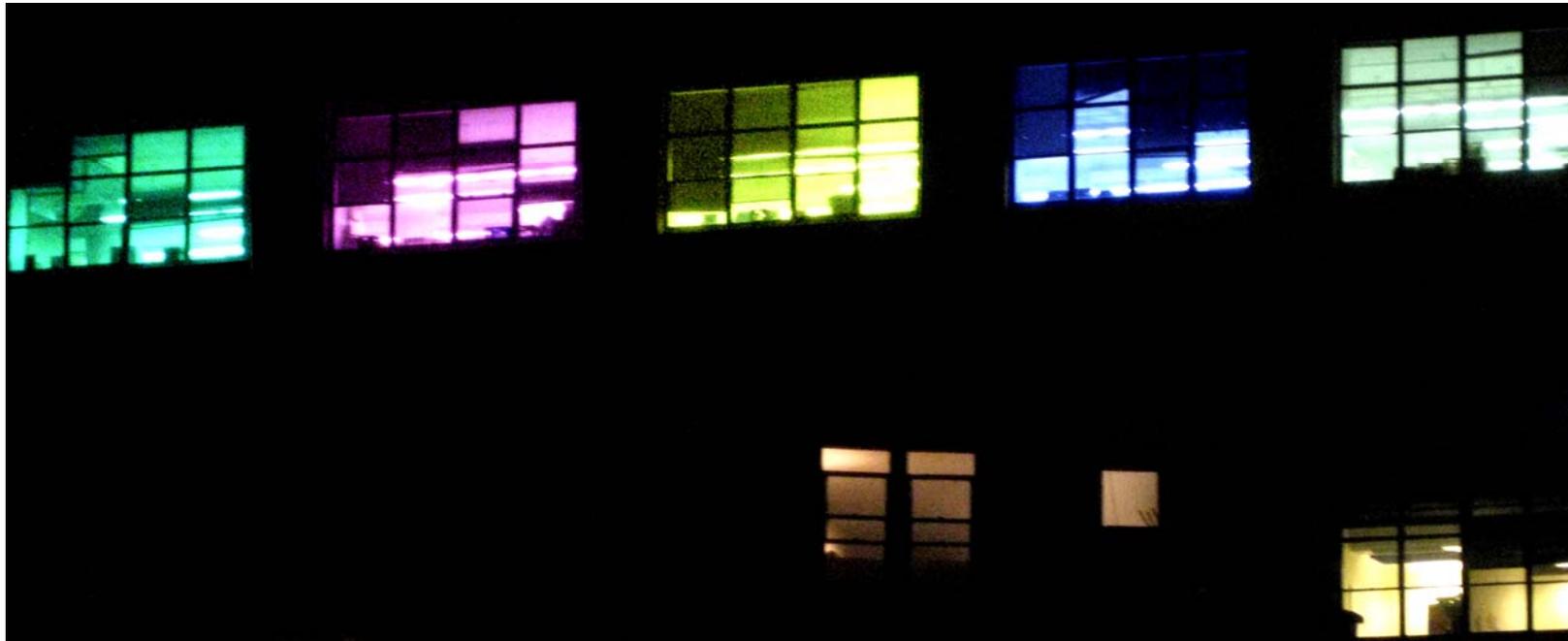
- **955** are assigned to known genes
- **1,463** are clustered into **~200** Novel Loci

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

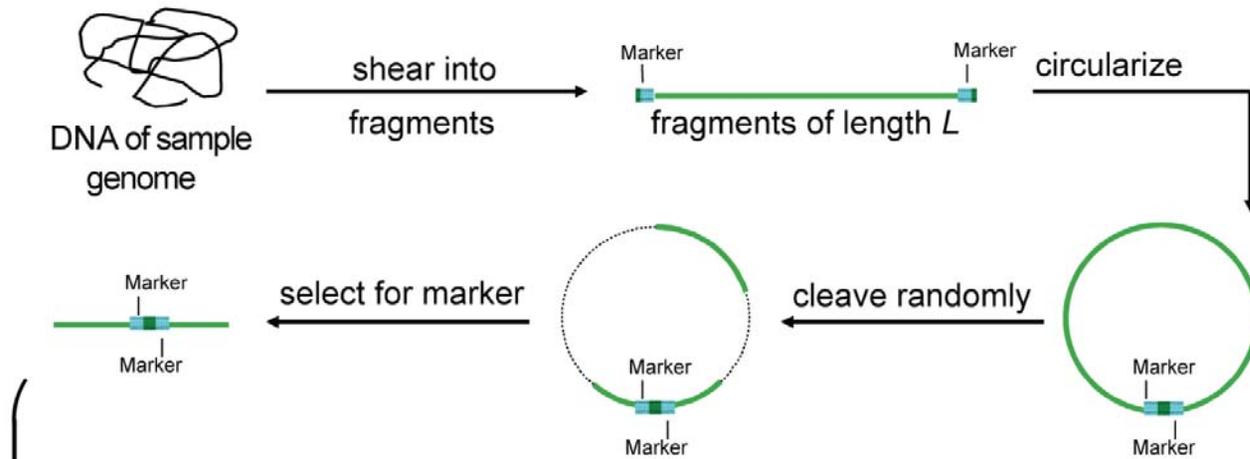
- ◇ RT-PCR & Sequencing
- ◇ 18/46 (39%) confirmed by RT-PCR
- ◇ 4/5 Sequenced Products Map uniquely to correct genomic region

Rozowsky et al. Genome Research (2007)

**Moving Beyond Arrays Next-
Generation Sequencing strategy for
characterizing genomes:
Paired End Mapping
to Find SVs**



Overall Strategy for Analysis of NextGen Seq. Data to Detect Structural Variants

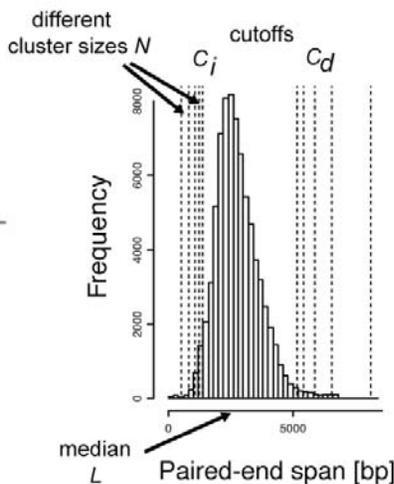
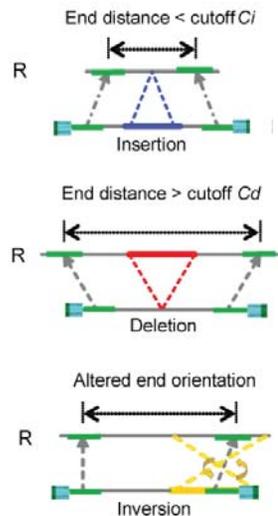
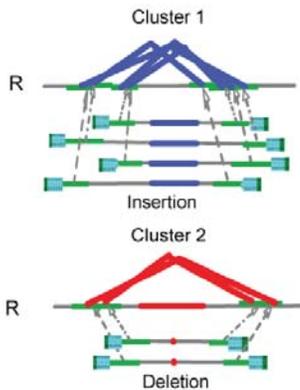


Next generation DNA sequencing, followed by PEMer analysis

- [1] *construct pre-processing*
- [2] *read-alignment*
- [3] *optimal paired-end placement*

[4] *outlier-identification*

[5] *outlier-clustering*



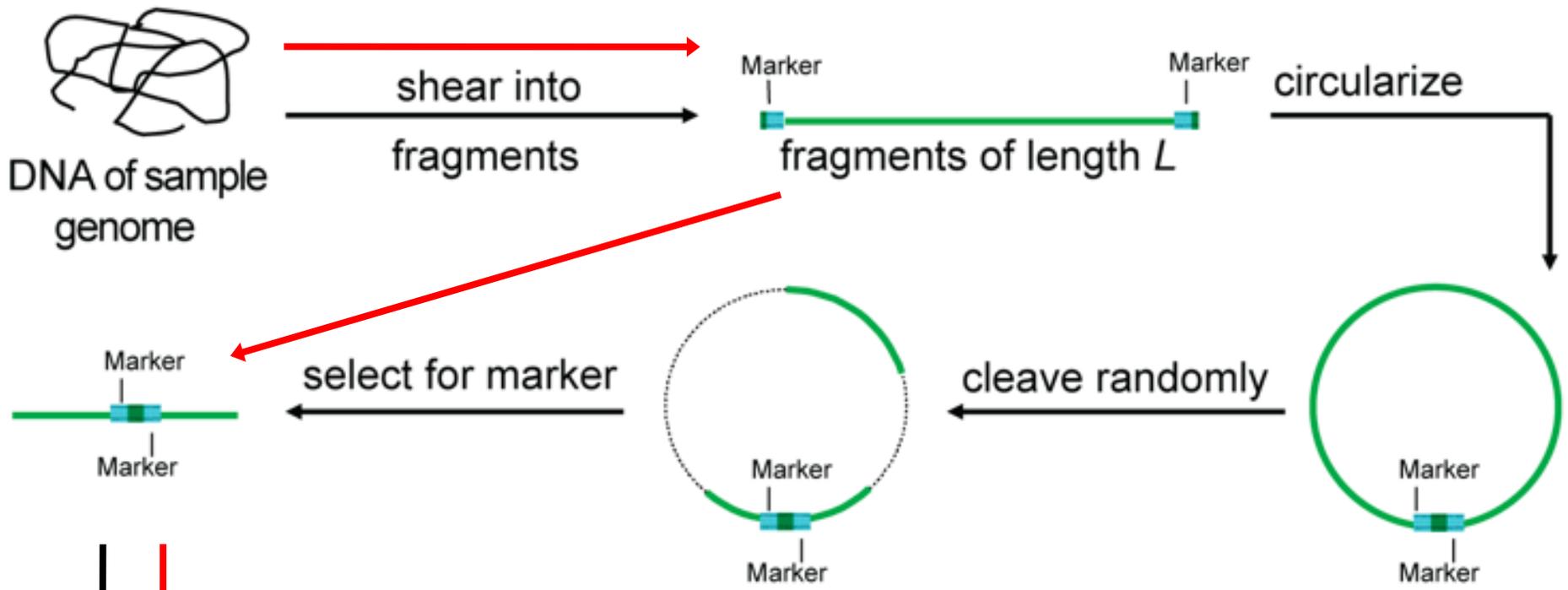
[6] *cluster-merging*

[7] *Display/storage of final SV set*



[Korbel et al., Science ('07); Korbel et al., GenomeBiol. (submitted)]

Simulation strategy



454 sequencing

[Korbel et al.,
GenomeBiol.
(submitted)]

— Simulation

— Experiment

[Korbel et al.,
GenomeBiol.
(submitted)]

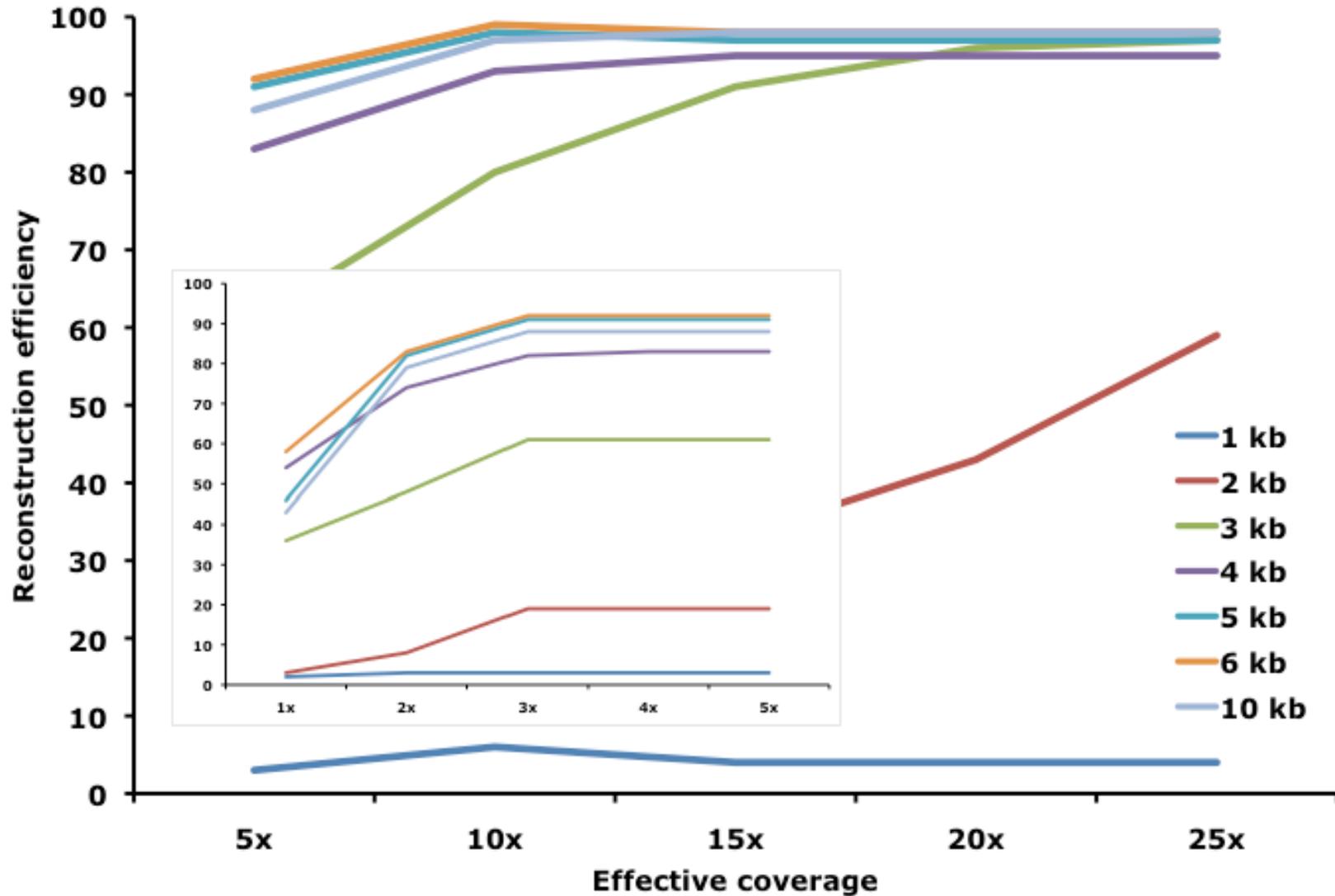
Reconstruction efficiency at 5x coverage

SV size	<i>Single cutoff</i>	<i>Multi-cutoff</i>	<i>Simplified multi-cutoff</i>	<i>Multi-cutoff(*)</i>	<i>Simplified multi-cutoff(*)</i>
1000	3(4)	3(4)	3(4)	3(4)	3(4)
2000	12(13)	23(26)	21(23)	11(13)	6(6)
3000	52(57)	61(68)	61(68)	49(52)	44(46)
4000	84(85)	85(86)	85(86)	80(82)	80(82)
5000	91(93)	91(93)	91(93)	91(93)	91(93)
6000	92(92)	92(92)	92(92)	92(92)	92(92)
10000	88(91)	88(91)	88(91)	88(91)	88(91)
Total	422(435)	443(460)	441(457)	414(427)	404(414)
False positives	31(31)	31(31)	26(31)	5(4)	2(1)

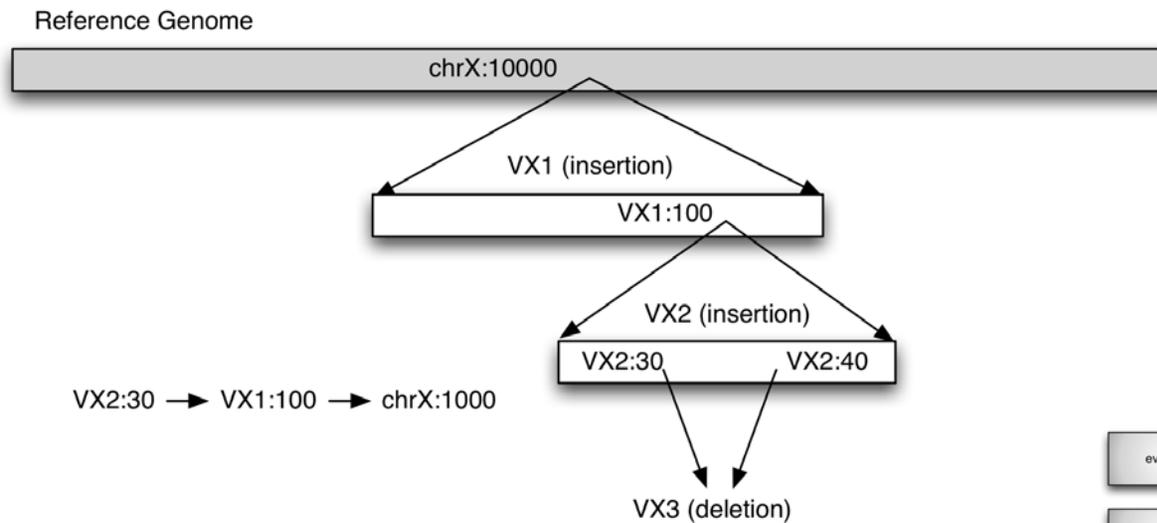
* optimal strategy; Multi-cutoff: overlaps 2-8; Simplified multi-cutoff: overlaps 2, 3, and 4

[Korbel et al.,
GenomeBiol.
(submitted)]

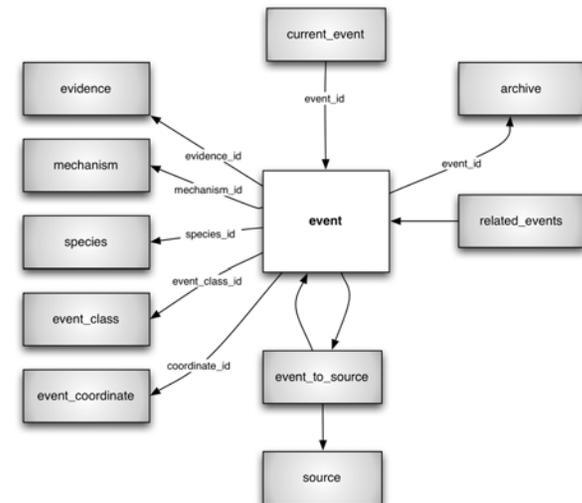
Reconstruction efficiency at different coverage



Building a Database of Variants: Complexities



[Korbel et al.,
GenomeBiol.
(submitted)]



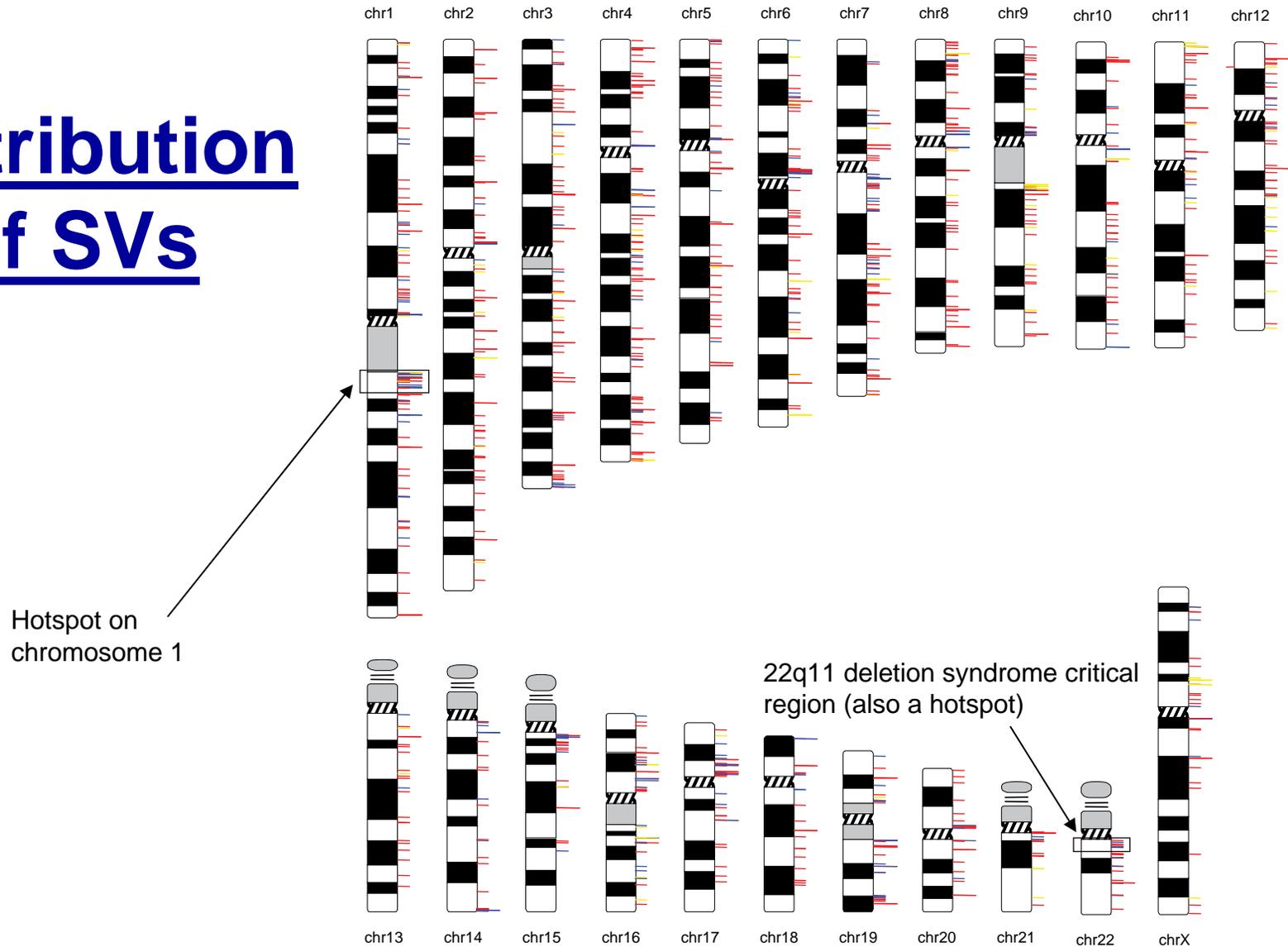
Summary of PEM Results

	NA15510 (Caucasian?, female)	NA18505 (Yoruba, female)
# of sequenced reads	> 10 M.	> 21 M.
Paired ends uniquely mapped	> 4.2 M.	> 8.6 M.
Fold coverage (on 6Gb)	~ 2.1x	~ 4.3x
Predicted Structural Variants*	478	839
<i>Indels</i>	427	758
<i>Inversion breakpoints</i>	51	81
Estimated total variants* with respect to ncbi36, genome-wide	759	902

*at this resolution 95

Korbel et al., 2007 *Science*

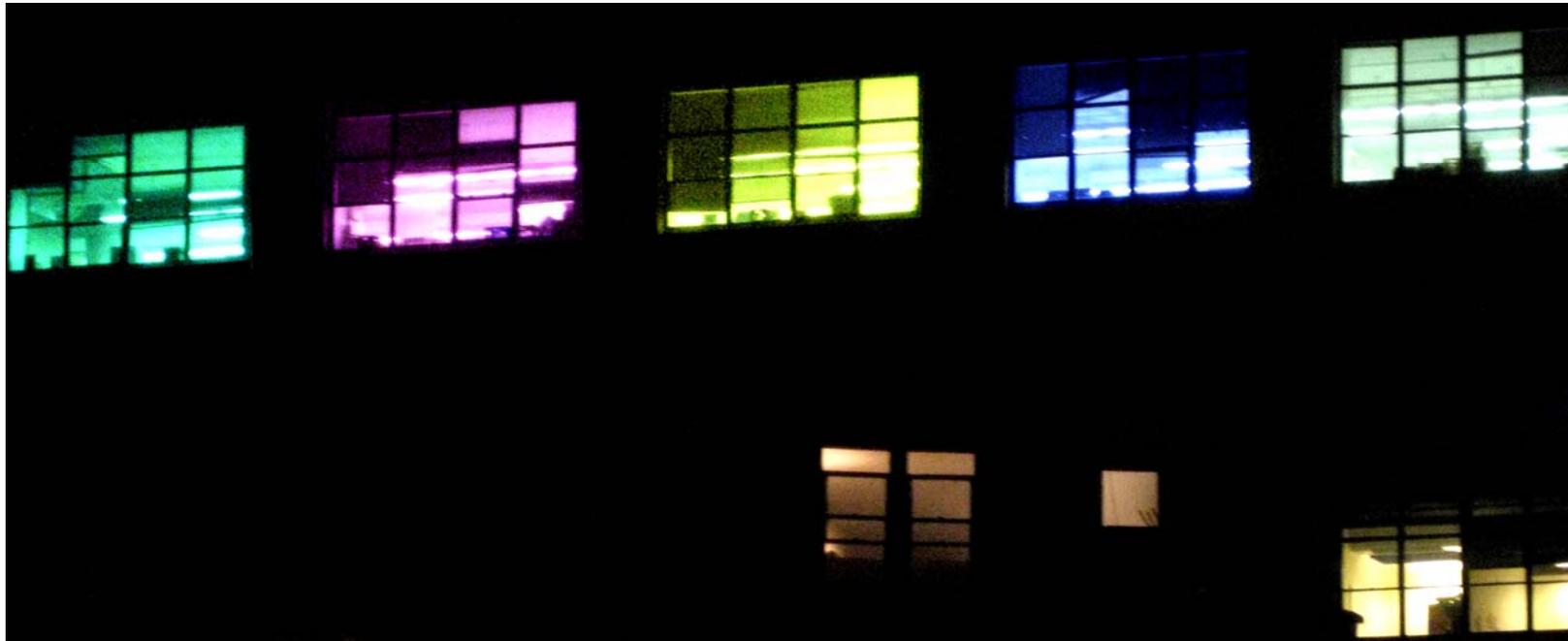
Distribution of SVs



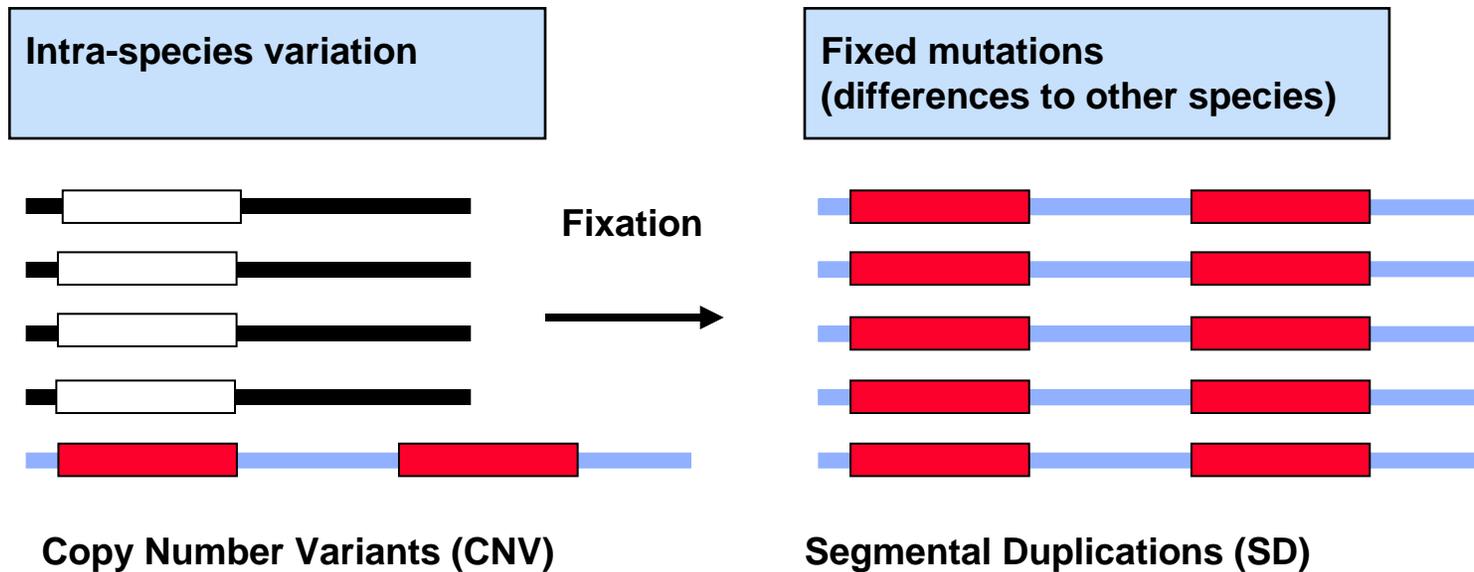
- red deletion; blue insertion; yellow inversion; double line length: same SV in both individuals.

Korbel et al., 2007 *Science*

Analyzing Duplications in the Genome (SDs & CNVs)



SEGMENTAL DUPLICATIONS AND COPY NUMBER VARIANTS ARE RELATED PHENOMENA AND SHOULD HAVE BEEN CREATED BY SIMILAR MECHANISMS



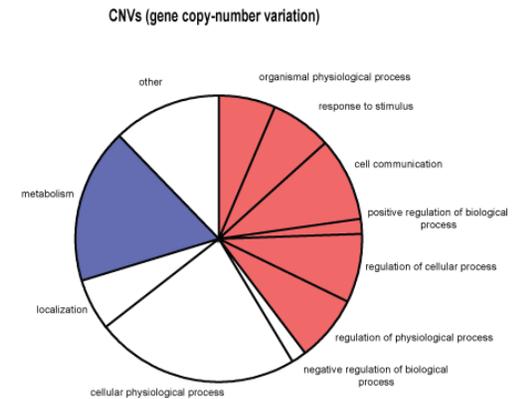
- SDs are the fixed forms of CNVs and arise when a CNV reaches fixation in the population.
- Hence, they should have been created by a similar mechanism

Association of SDs and CNVs with pseudogenes

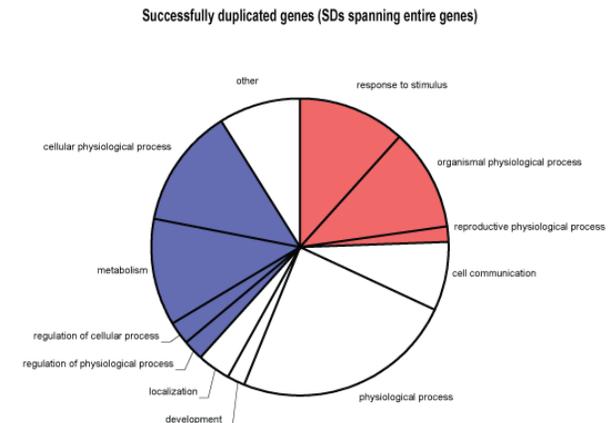
- CNVs are the raw form of variation producing duplicated elements
- Segmental Duplications (SDs) are fixed forms of CNVs/SVs. They give rise to duplicated genes and (eventually) protein families
- Thus, we expect, duplicated pseudogenes (failed duplications) to occur in SDs.
- CNVs and SDs tend to be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

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

A

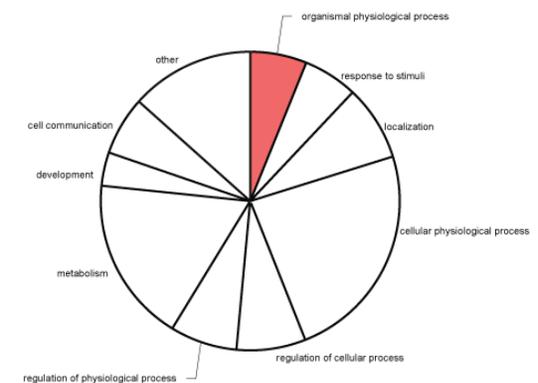


B

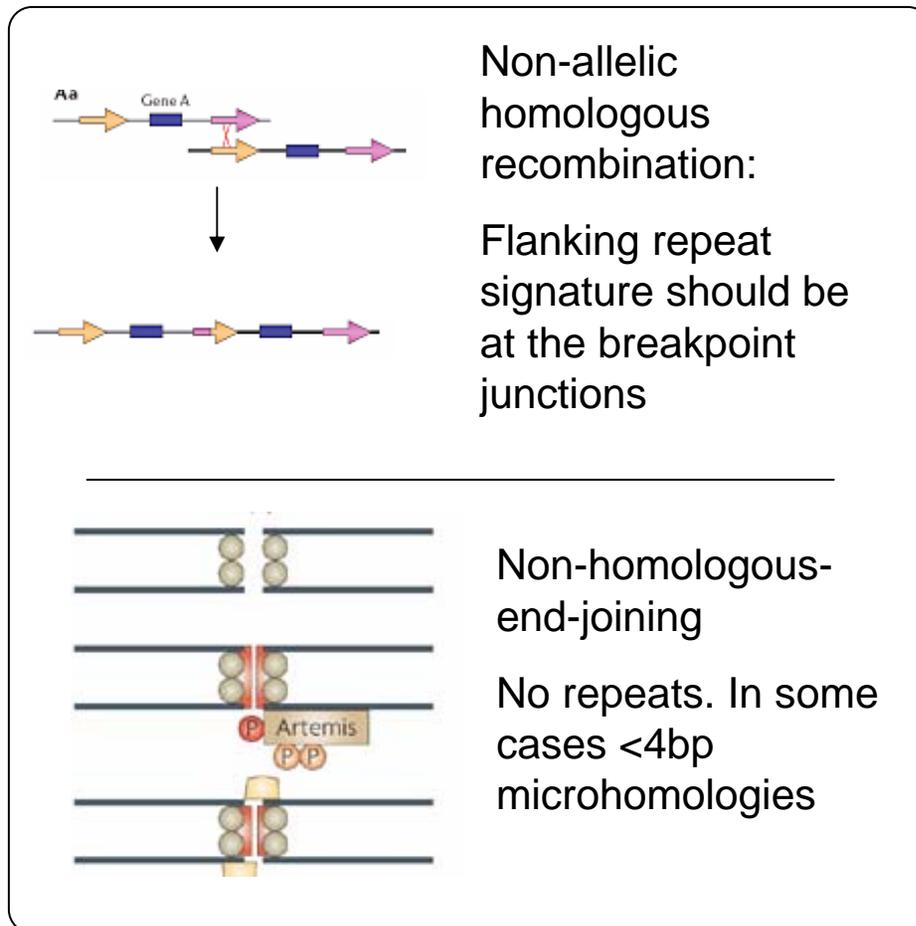


C

Unsuccessful duplicates (duplicated genes inactivated by disruption of coding sequence)



SEVERAL DIFFERENT MECHANISMS HAVE BEEN PROPOSED FOR THE GENESIS OF SDs AND CNVs



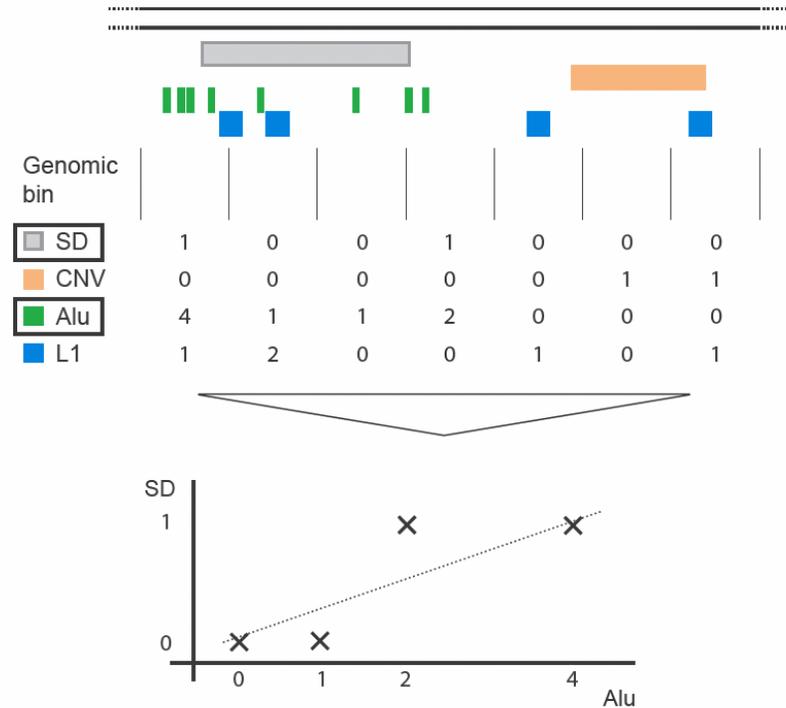
- We can examine breakpoint sequences to determine the mechanism

Problem:

- Most large-scale CNV data is of low resolution (50kb or worse)
- Cannot directly observe repeat signatures in the large-scale data!

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

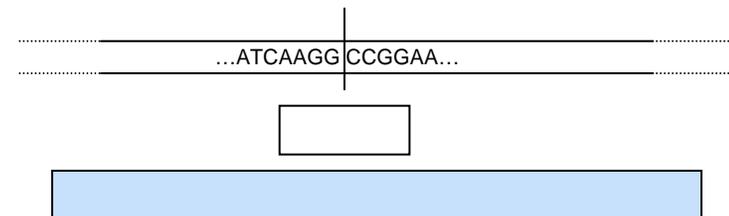
- ① Survey a range of genomic features
- ② Count the number of features in each genomic bin (100kb)
- ③ Calculate correlations / enrichments



If exact CNV breakpoints are known, we can calculate the enrichment of repeat elements relative to the genome or relative to the local environment

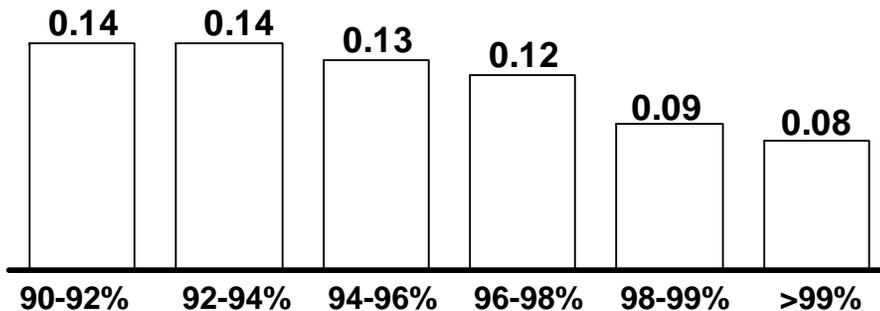
Exact match

Local environment

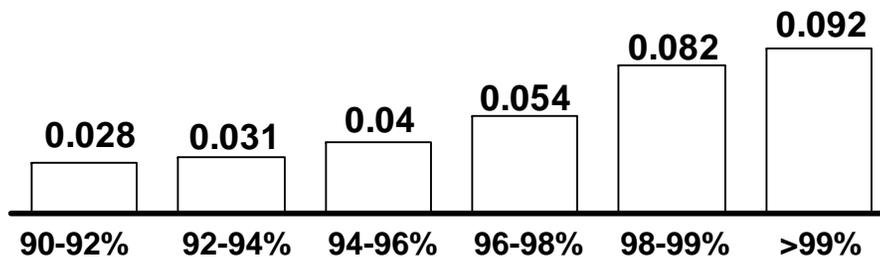


OLDER SDs ARE MUCH MORE LIKELY TO BE FORMED BY ALU ELEMENTS

Alu association with SDs by age



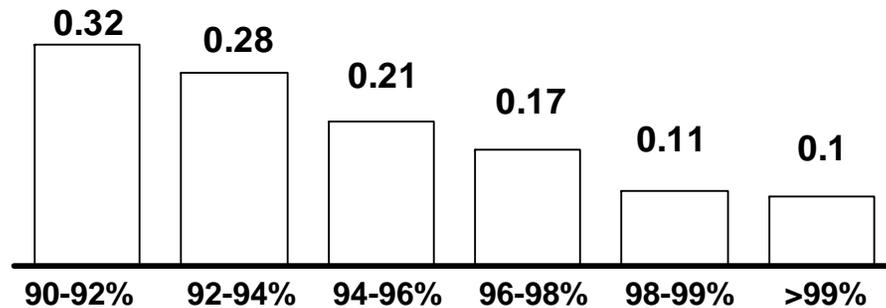
SD association with subtelomeres



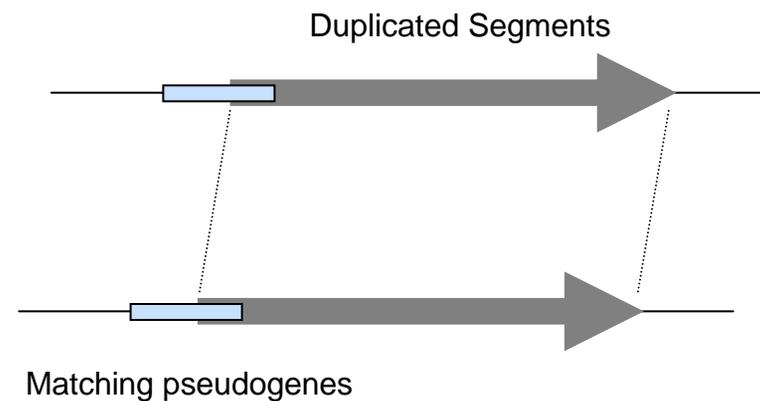
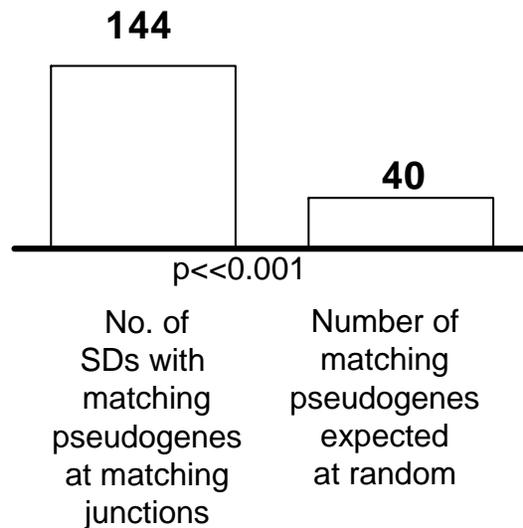
- The co-localization of Alu elements with SDs is highly significant.
- Older SDs have a much higher association with Alus than younger SDs.
- Hence it is likely, that Alu elements were more active in mediating NAHR in the past (consistent with the Alu burst)
- Younger SDs are more likely to be localized in subtelomeres (instable region susceptible to NHEJ)

ANOTHER FUNCTION FOR PSEUDOGENES: SERVING AS REPEATS FOR MEDIATING NAHR

Processed pseudogene association with SDs by age

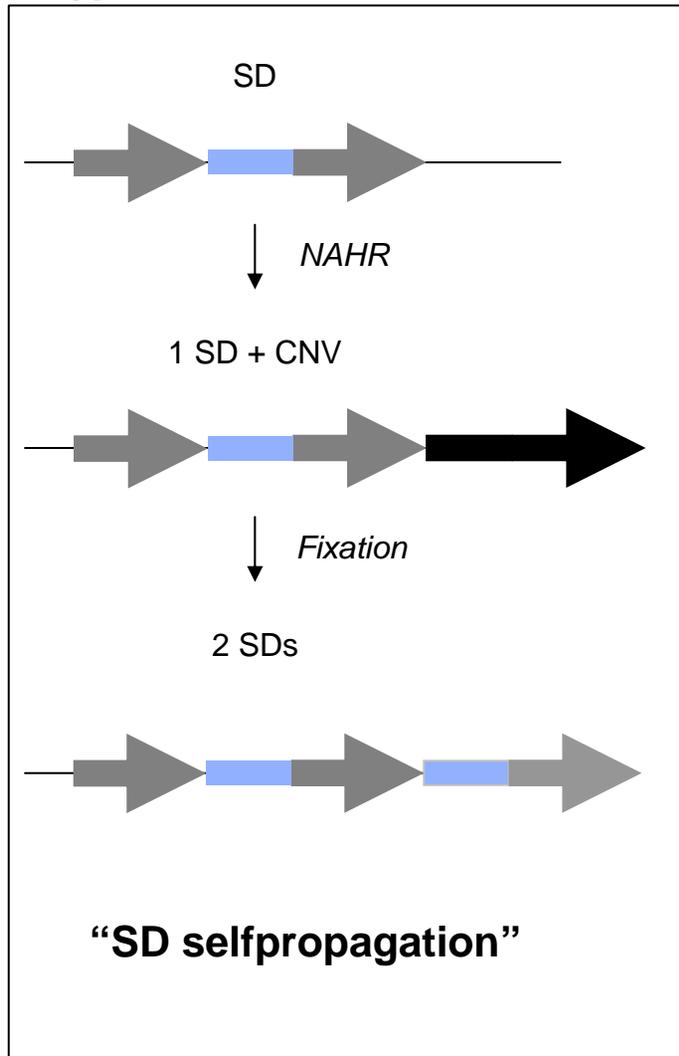


Processed pseudogenes at SD junctions



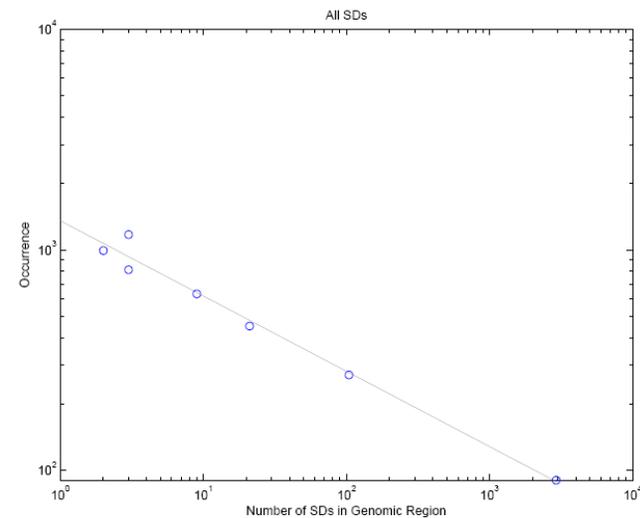
FOCUSSING ON SDS: SDS CAN PROPAGATE THEMSELVES, WHICH LEADS TO A POWER-LAW DISTRIBUTION

Hypothesis



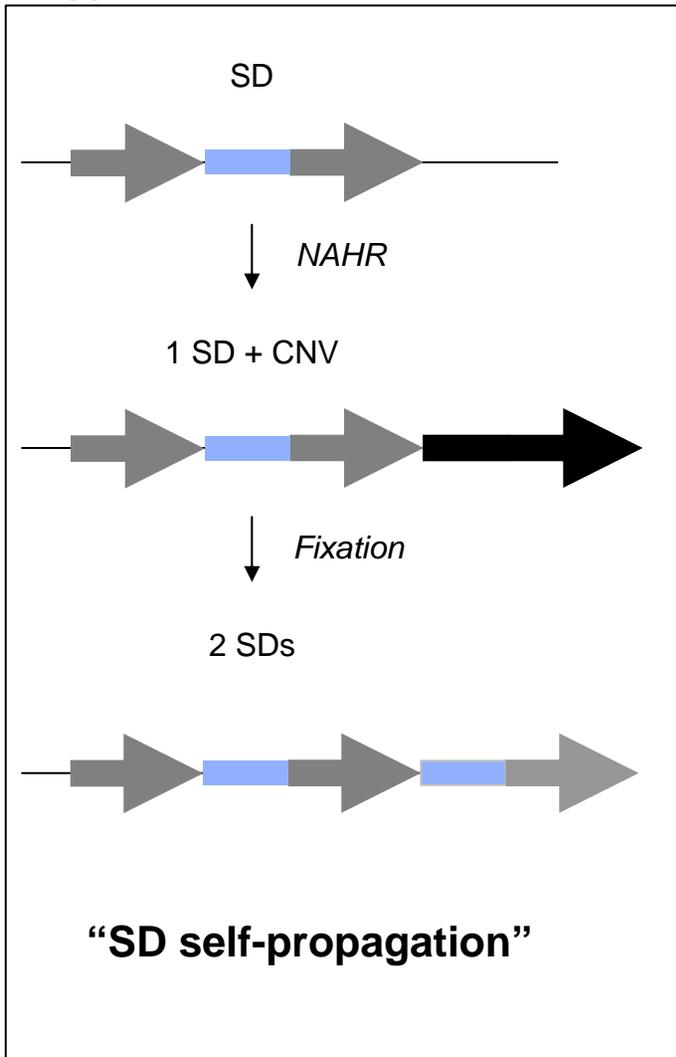
Corollary

- 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 lead to a very skewed (“power-law”) distribution of SDS.



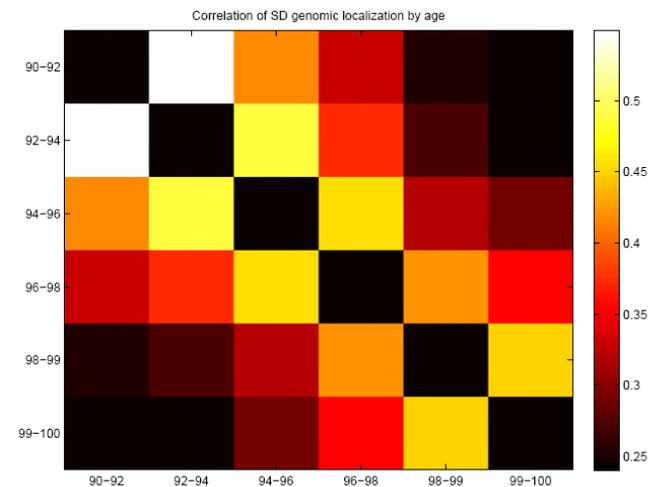
FOCUSSING ON SDS: SDS COLOCALIZE WITH EACH OTHER

Hypothesis



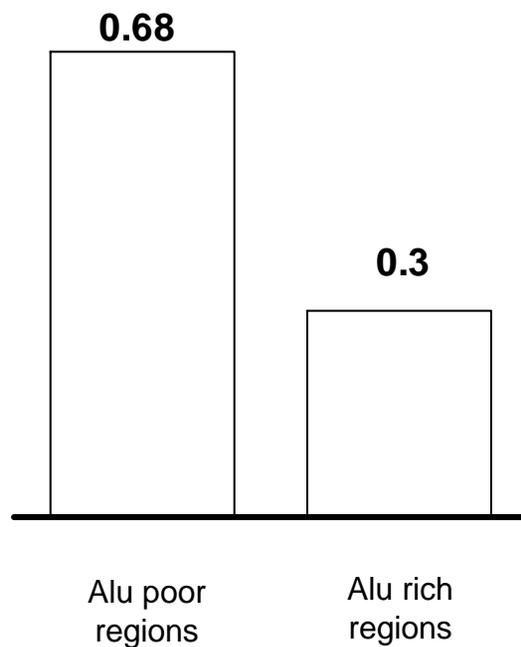
Corollary

- 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:



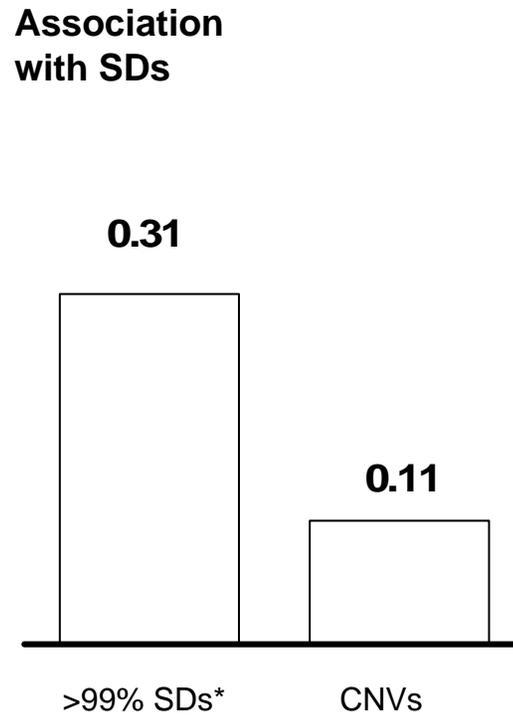
THE ASSOCIATION OF SDs WITH ALU ELEMENTS IS COMPLEMENTARY TO THE ONE WITH SDs

Correlation of young SDs (>99%) with older SDs



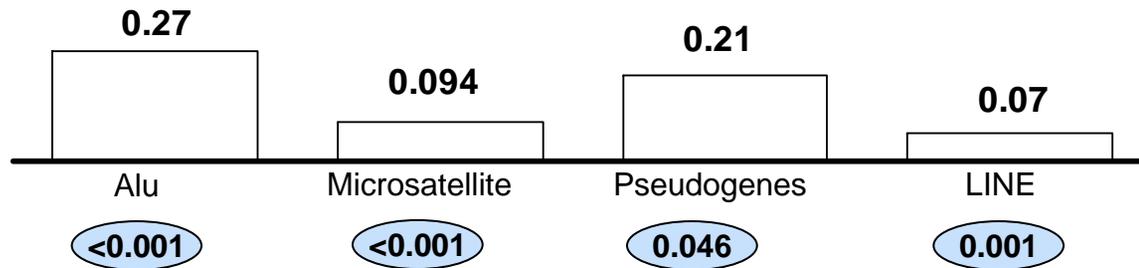
- SDs that lie in Alu rich regions are less likely to be associated with other SDs
- Hence, there is a certain complementarity to SD-mediated NAHR with Alu-mediated NAHR

CNVs ARE LESS ASSOCIATED WITH SDs THAN THE GENERAL SD TREND

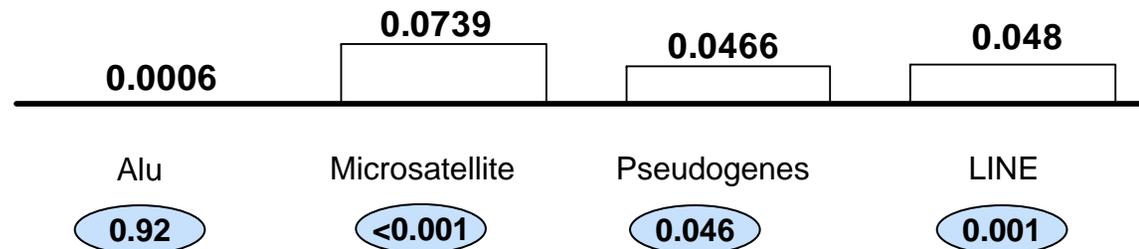


ASSOCIATIONS ARE DIFFERENT FOR SDs AND CNVs

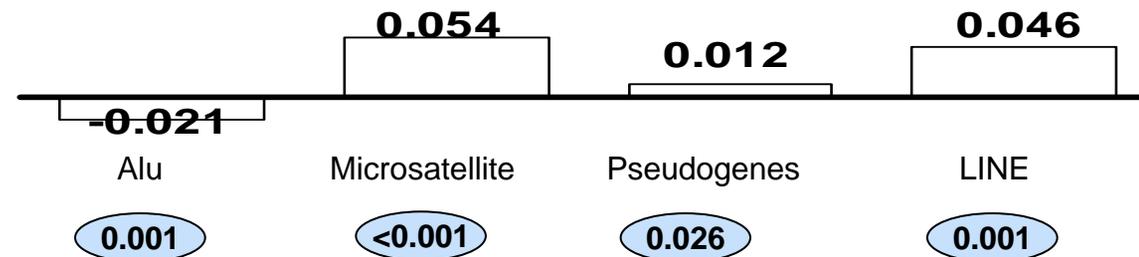
SD association with repeats



CNV association with repeats

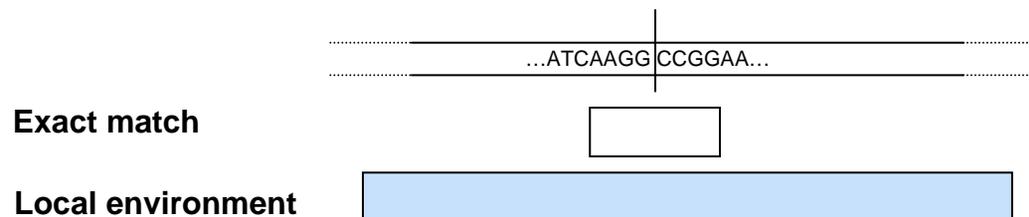


CNV association with repeats after correcting for SD content

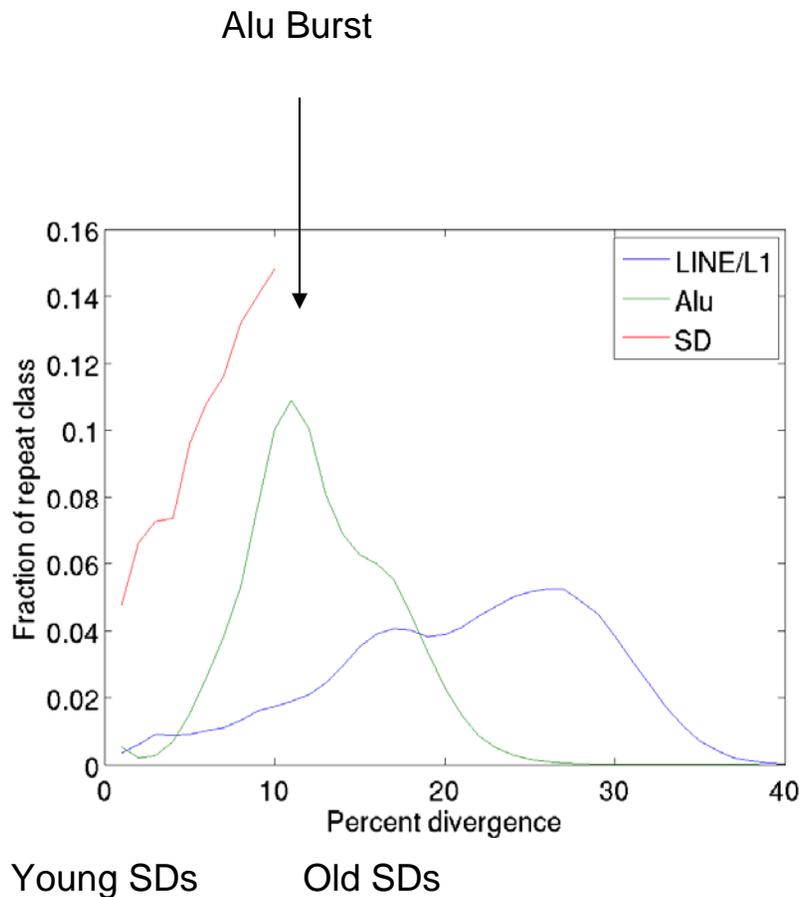


ANALYZING SEQUENCED BREAKPOINTS CONFIRMS THE RESULTS FROM THE COARSE GRAINED ANALYSIS

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

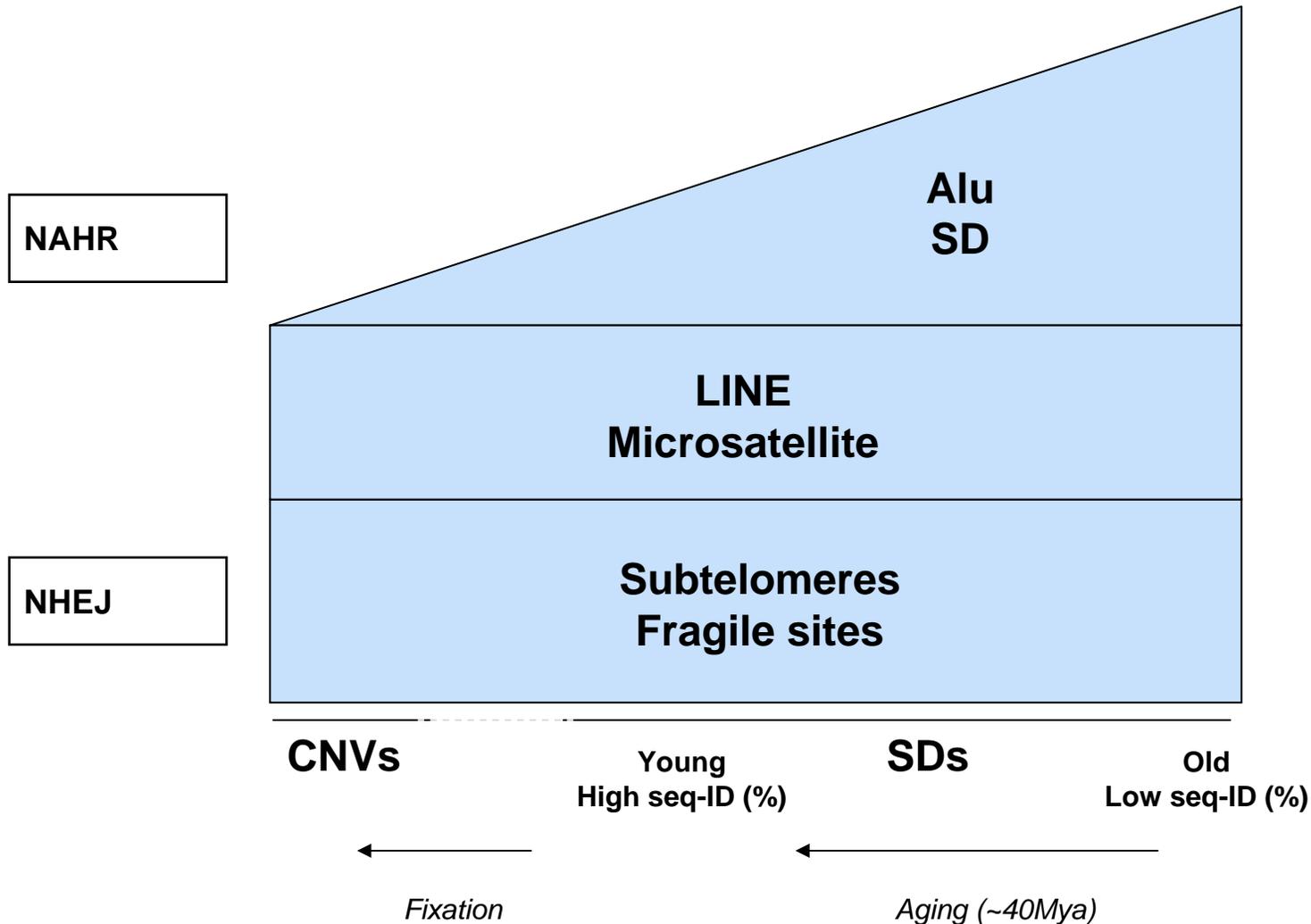


AFTER THE ALU BURST, THE IMPORTANCE OF ALU ELEMENTS FOR GENOME REARRANGEMENT DECLINED RAPIDLY



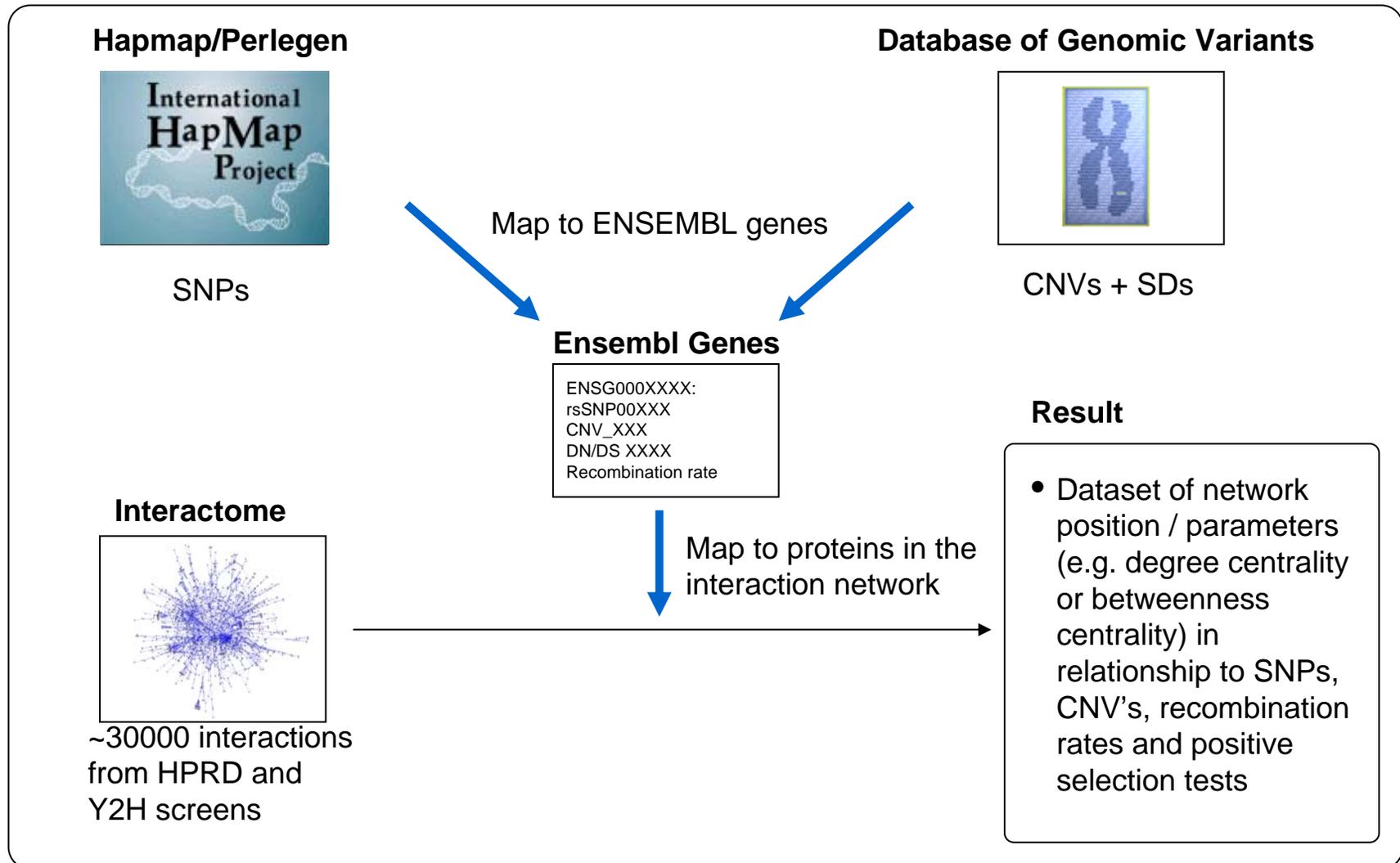
- About 40 million years ago there was a burst in retrotransposon activity
- The majority of Alu elements stem from that time
- This, in turn, led to rapid genome rearrangement via NAHR
- The resulting SDs, could create more SDs, but with Alu activity decaying, their creation slowed

THE MECHANISM DRIVING LARGE SCALE GENOME REARRANGEMENT UNDERWENT A MARKED SHIFT IN THE LAST 40 MYA



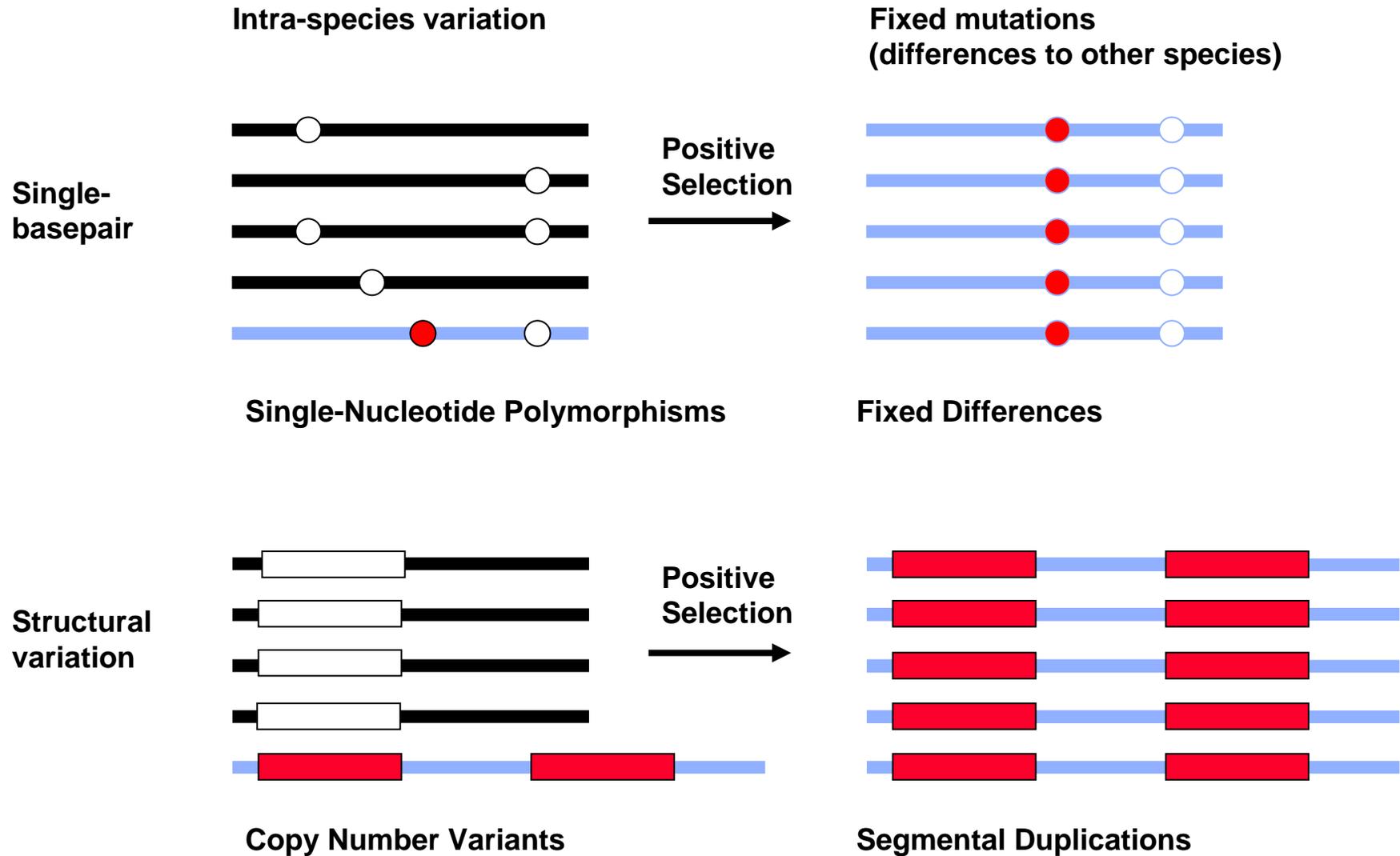
METHODOLOGY: MAP SNP AND CNV DATA ONTO ENSEMBL GENES, AND THEN MAP ENSEMBL GENES TO THE KNOWN INTERACTOME

ILLUSTRATIVE



* From Nielsen et al. *PLoS Biol.* (2005) and Bustamante et al. *Nature* (2005)

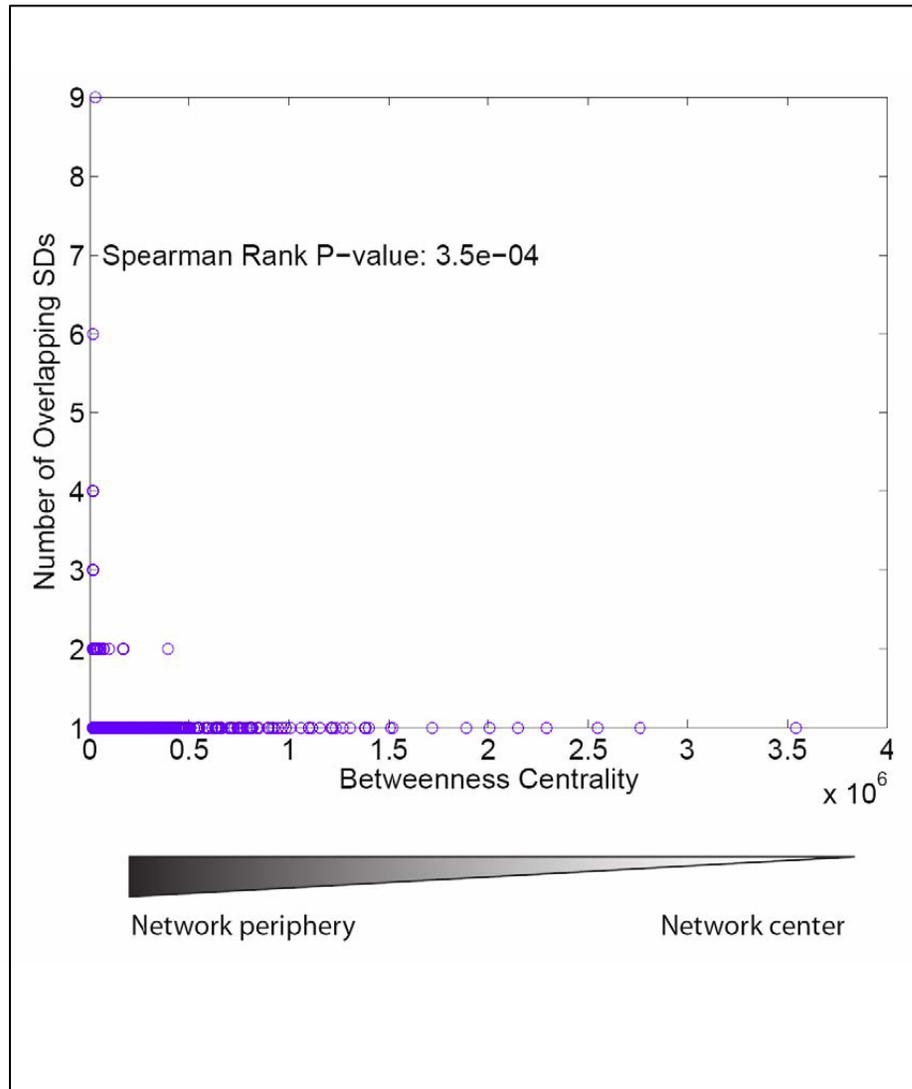
ADAPTIVE EVOLUTION CAN BE SEEN ON TWO DIFFERENT LEVELS



[Often but not always measured by dN/dS]

CENTRAL NODES ARE LESS LIKELY TO LIE INSIDE OF SDs

Centrality vs. SD occurrence



Reasoning

- This result also confirms our initial hypothesis – peripheral nodes tend to lie in regions rich in SDs.
- Since segmental duplications are a different mechanism of ongoing evolution, the less constrained peripheral proteins are enriched in them.
- Note that despite the small size of our dataset for known SD's we get significant correlations. It is to be expected that the correlations will get clearer as more data emerges*

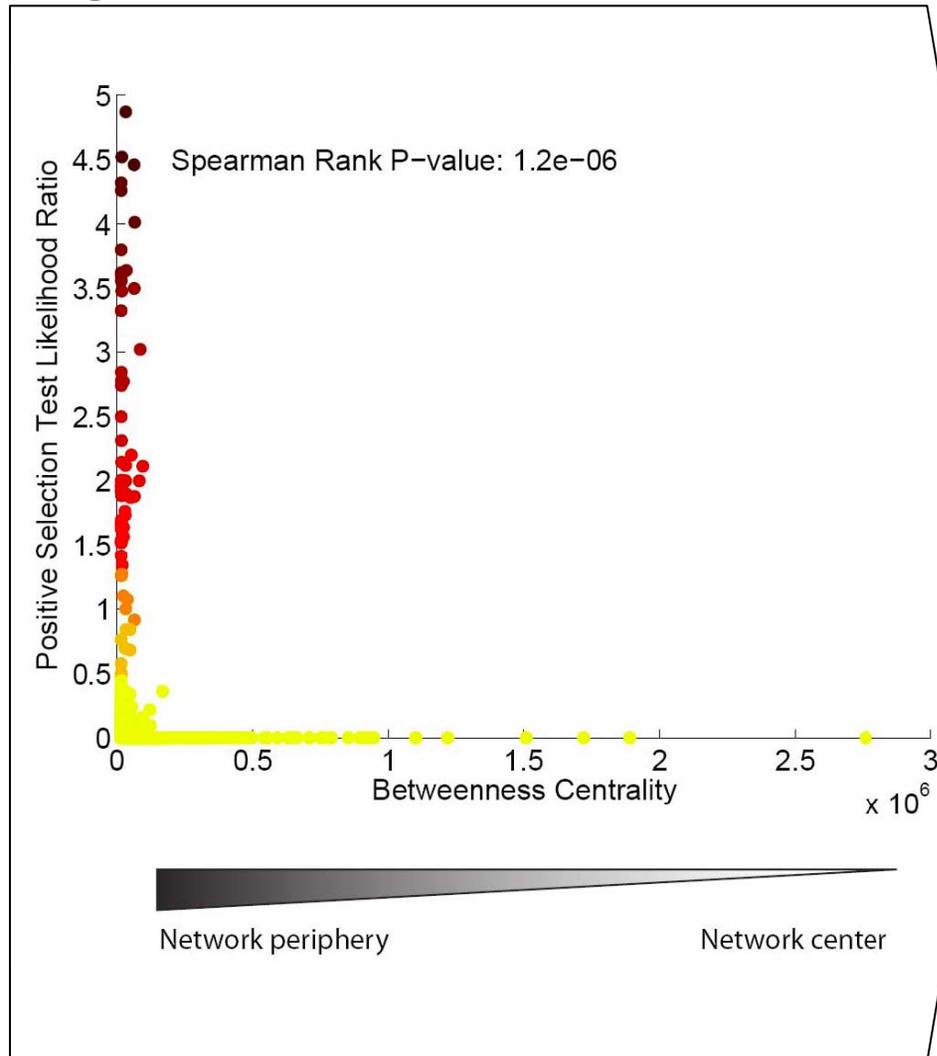
* Specifically, a number of the SDs are likely not fixed, but rather common CNVs in the reference genome

Source: Database of genetic variation, HPRD, Rual et al. *Nature* (2005), and Kim et al. *PNAS* (2007)

CENTRAL PROTEINS ARE LESS LIKELY TO BE UNDER POSITIVE SELECTION

▭ Hubs

Degree vs. Positive Selection



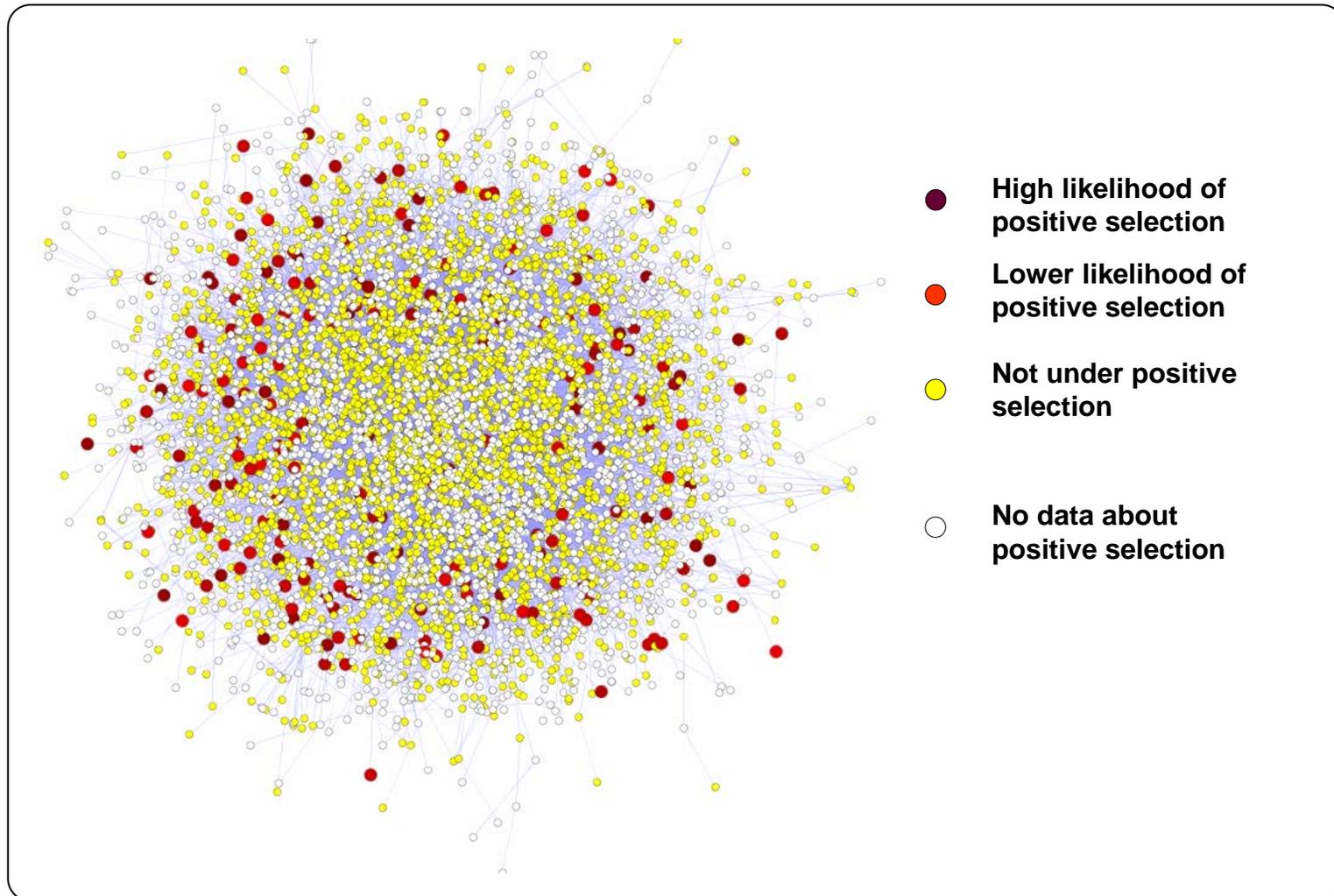
Reasoning

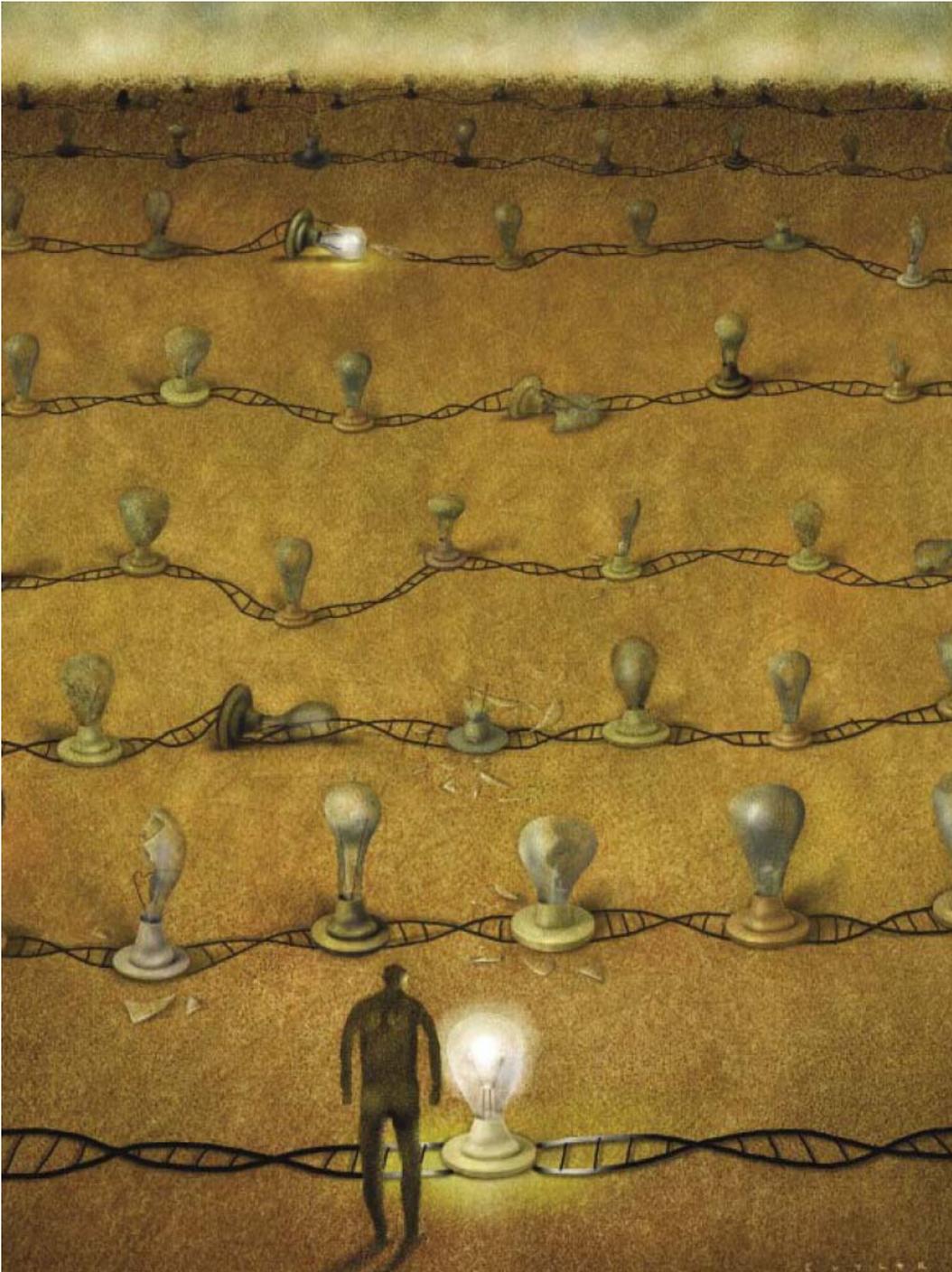
- Peripheral genes are likely to under positive selection, whereas hubs aren't
- This is likely due to the following reasons:
 - Hubs have stronger structural constraints, the network periphery doesn't
 - Most recently evolved functions (e.g. “environmental interaction genes” such as sensory perception genes etc.) would probably lie in the network periphery
- Effect is independent of any bias due to gene expression differences

* With a probability of over 80% to be positively selected as determined by Ka/Ks. Other tests of positive selection (McDonald Kreitmann and LDD) corroborate this result.

POSITIVE SELECTION LARGELY TAKES PLACE AT THE NETWORK PERIPHERY

Positive selection in the human interactome





Integrative Analyses:

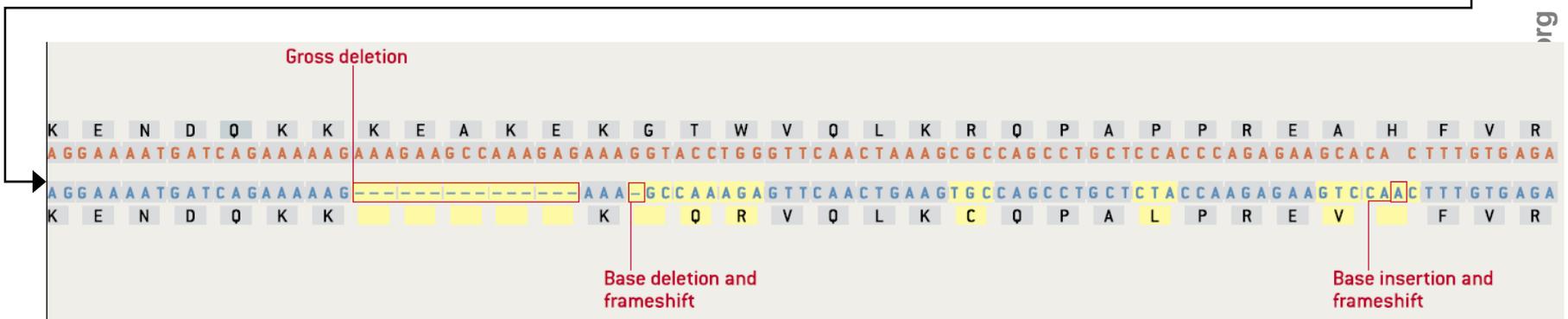
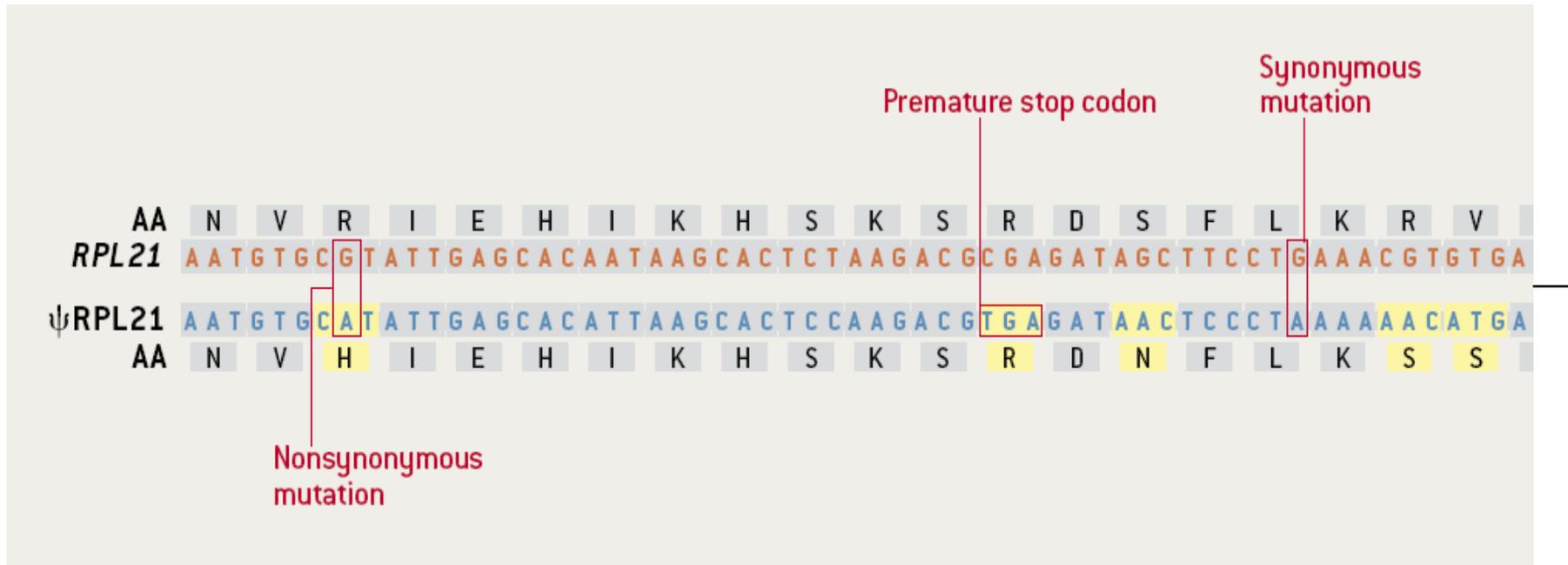
Annotating Pseudogenes and relating them to functional signal and measures of conservation

Illustration from Gerstein & Zheng (2006). Sci Am.

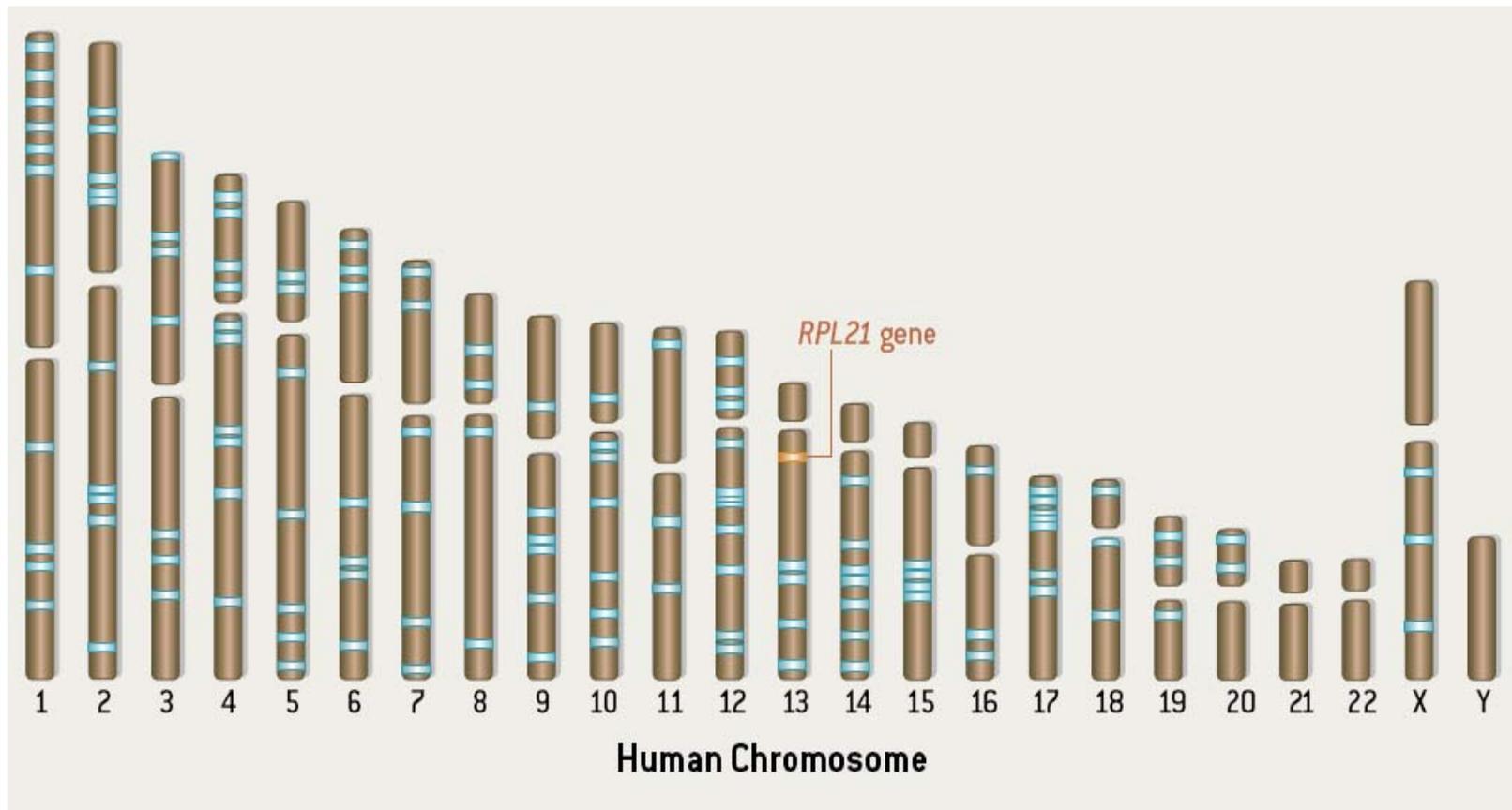
Pseudogenes are among the most interesting intergenic elements

- Formal Properties of Pseudogenes (Ψ G)
 - ◇ Inheritable
 - ◇ 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?...

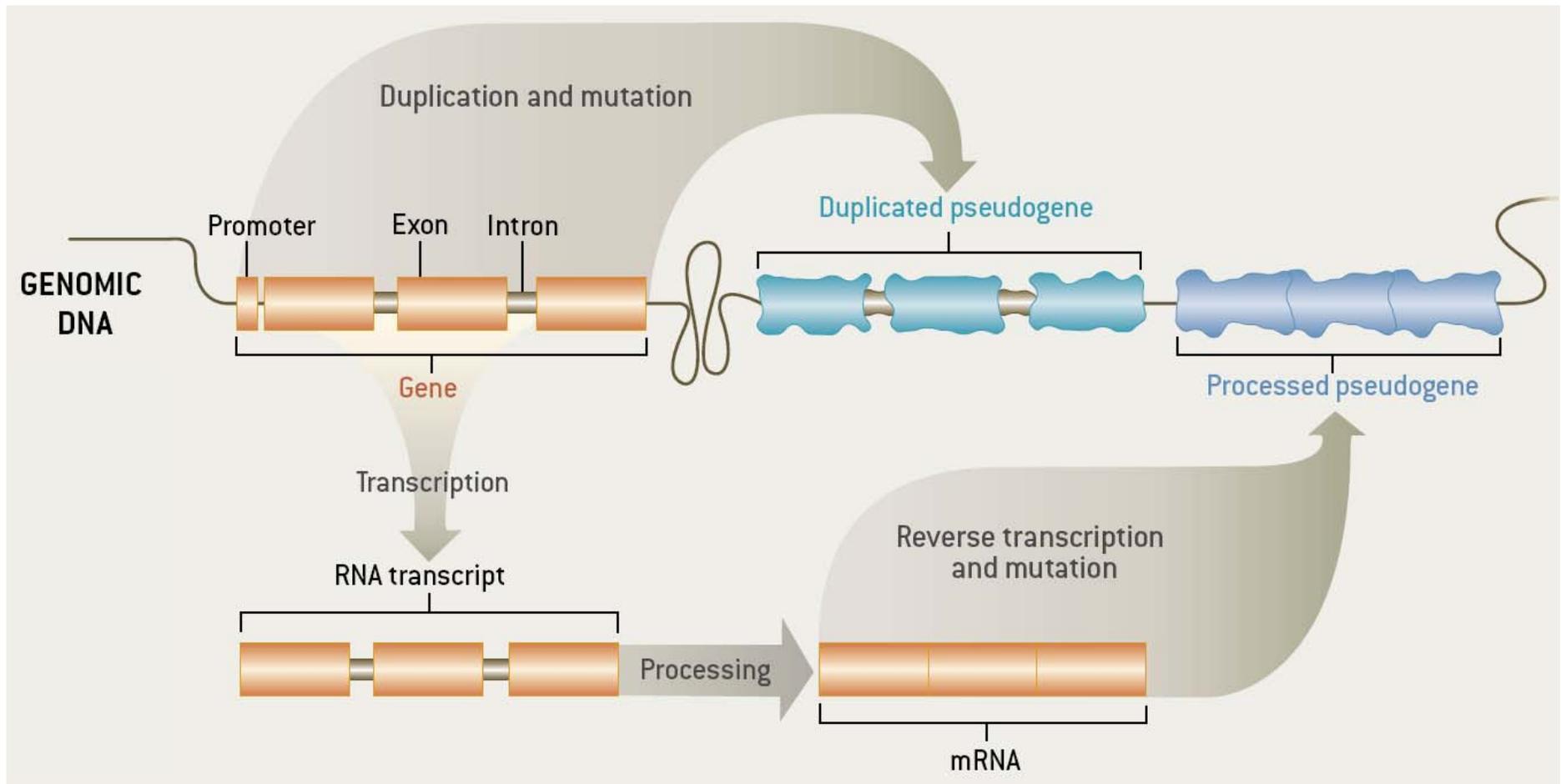
Identifiable Features of a Pseudogene (ψ RPL21)



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?

➤ Important for Doing Accurate Gene Annotation

- Abundant: > 8000 retropseudogenes in human
- High sequence similarity with genes
- 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

[Ruud, *Int. J. Cancer* 1999]

➤ ΨG are “genomic fossils”

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

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 functional genes

- Cross-hybridation in microarray and PCR (Cytokeratin 19, *Int. J. Cancer* 1999)
- Very rarely this gives some pseudogenes regulatory roles

➤ Ψ G are "genomic fossils"

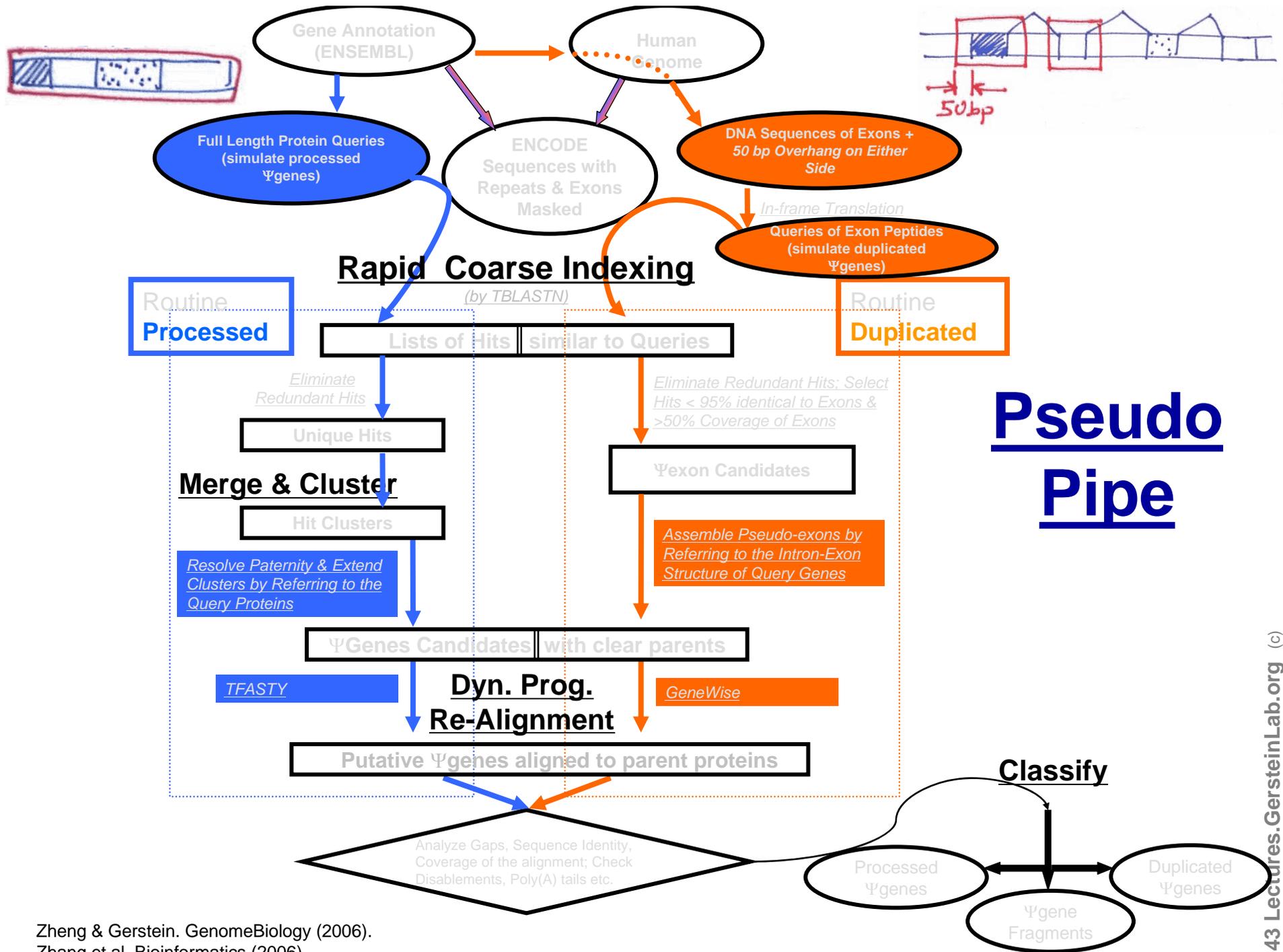
- Study the evolution of genes and genomes
- Measure mutation/insertion rates

In mouse, a pseudogene up-regulates gene expression of *Makorin1* by binding to a transcriptional repressor or an RNA-digesting enzyme [Hirotsume *et al. Nature* **423** 2003]

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. [Ruud, *Int. J. Cancer* 1999]
 - Some pseudogenes have regulatory roles
- Ψ G are “genomic fossils”
 - Study the evolution of genes and genomes
 - Measure mutation/insertion rates





Zheng & Gerstein. GenomeBiology (2006).
 Zhang et al. Bioinformatics (2006)

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

Karro et al., NAR (2007)

DB

Pseudogene.org

ABOUT | PUBLICATION DATA | DATABASE | KNOWLEDGEBASE

Eukaryote Database

CHICKEN	CHIMP
Tax ID: 9031 Build: 2 Pseudogenes: 4179	Tax ID: 9598 BU
Search: Pseudogenes , Sets	Search: Pseudogenes
Download Flatfile: [gtf] [txt] [by chr]	Download Flatfile: [gtf]

DOG	FLY
Tax ID: 9615 Build: 1 Pseudogenes: 2802	Tax ID: 7227 BU
Search: Pseudogenes , Sets	Search: Pseudogenes

Pseudogene Sets

Search	Name	Notes and Links	Reference	Current Build	Current Size	Original Build	Original Size	Download
1	Build 28 pseudogenes	7888 pseudogenes found by Zhaolei Zhang ...	Zhang et. al.	36	6186	28	7888	[gtf] [txt]
2	Pseudogene.org pipeline output	13979 pseudogenes were identified by our...	www.pseudogene.org	36	13111	34	13979	[gtf] [txt]
3	Human chr. 22	522 pseudogenes were identified on chrom...	Zheng et. al.	36	402	34	522	[gtf] [txt]
4	Transcribed Pseudogenes	201 transcribed pseudogenes were identified	urn:lsid:pseudogene.org:9606.PFSet:3 522 pseudogenes were identified on chromosome 22 by Deyou Zheng based on human build 34. Later this set was mapped to the current human build 35, 509 pseudogenes remaining. Original Analysis is accessible by clicking the link.	36	200	34	201	[gtf] [txt]
5	Bork Pseudogenes	19629 pseudogenes were identified by D.T...		36	17493	34	19629	[gtf] [txt]
6	Hopsigen pseudogenes	5791 processed pseudogenes	Khelifi et. al.	36	4080	34	5791	[gtf] [txt]

Links	Pseudogene Accession Number (tsid format)	Name	Chromosome	Start	Stop	Strand	Type	Protein
	urn:lsid:pseudogene.org:9606.Pseudogene:48953	ENSP00000321017.Human.chr22.mb14	22	14456056	14456172	-	FP	ENSP00000
	urn:lsid:pseudogene.org:9606.Pseudogene:5542	ENSP00000347298.Human.chr22.mb14	22	14464249	14464833	-	Processed	ENSP00000
	urn:lsid:pseudogene.org:9606.Pseudogene:7182	ENSP00000252487.Human.chr22.mb14	22	14502720	14503765	-	Processed	ENSP00000

5 Methods of Assignment

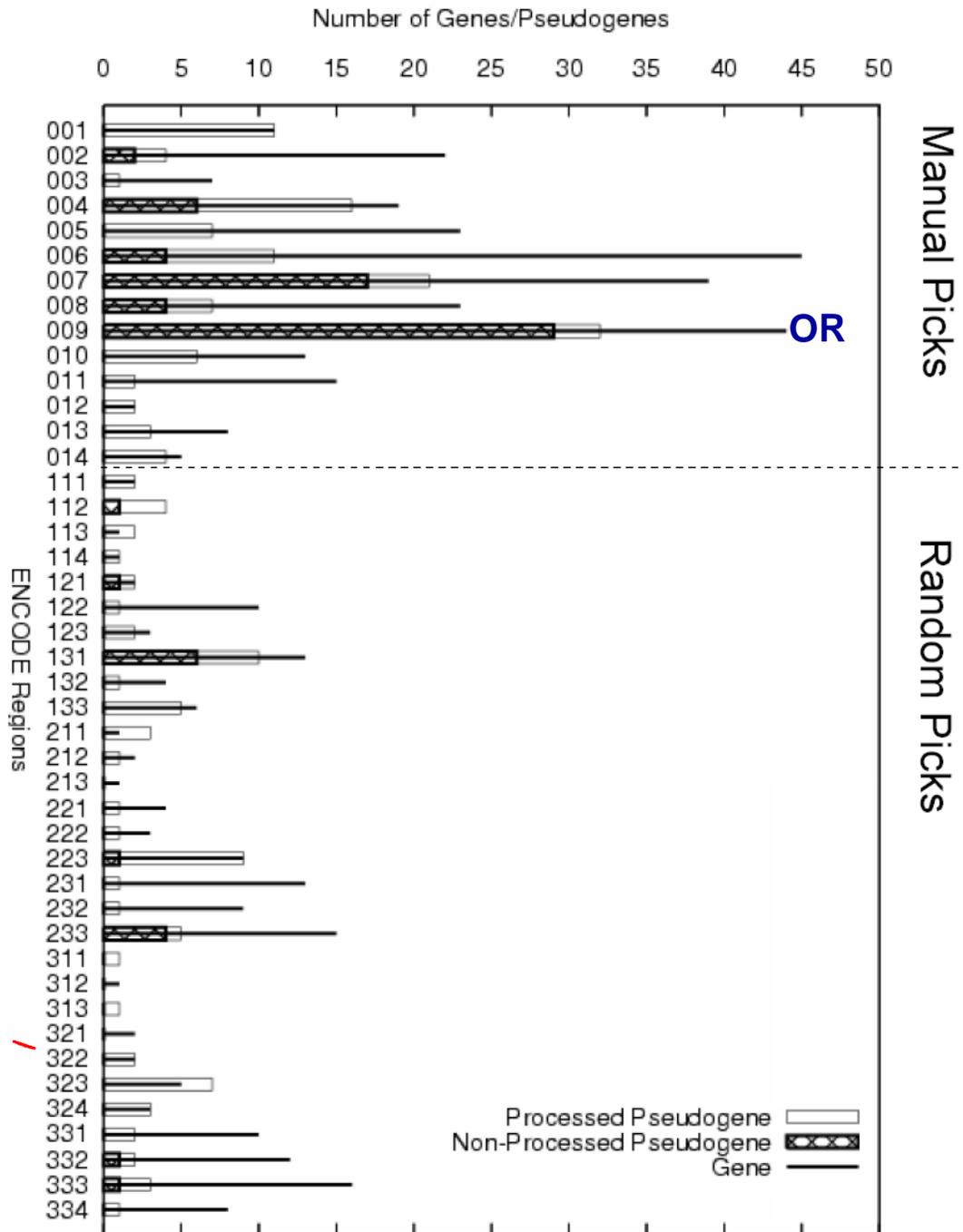
- **4 automatic pipelines**
 - ◇ retroFinder+pseudoFinder (UCSC), PseudoPipe (Yale), GIS
 - ◇ Comparing protein or transcript v genomic DNA, filtering, application of rules
- **HAVANA manual**
- **What is a pseudogene?**
 - ◇ Different criteria
- **Conservative approach here**
 - ◇ Can't overlap gene annotation
 - ◇ Need to have a protein alignment
 - ◇ 201 pseudogenes vs ~400 genes

Overall Results: Regional Distribution

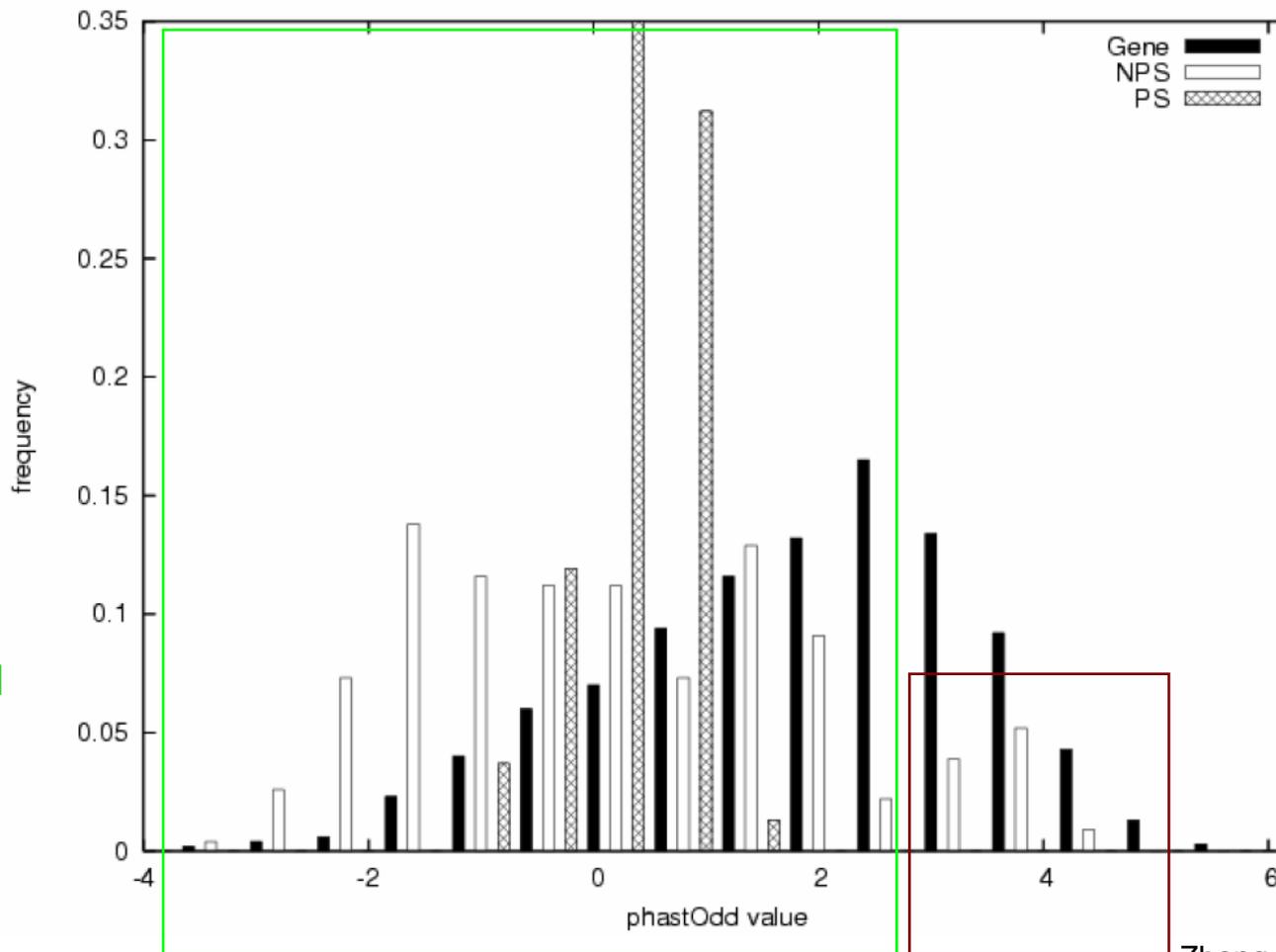
**201 pseudogenes
77 non-processed
124 processed**

Zheng et al. (2007) Gen. Res.

**browser
+
pseudogene.org/ENCODE**



Using phastOdd value to examine neutral evolution of pseudogenes



most good candidates for studying mutational processes

a few non-proc. ψ G under constraint

Zheng et al. (2007) Gen. Res.

representative pseudogenes drawn from 201 total

	A	B	C	D	E	F
human -	⊗	⊗	⊗	⊗	⊗	⊗
chimp -	⊗	■	⊗	⊗	■	■
baboon -	⊗	⊗	⊗	⊗	⊗	■
macaque -	⊗	⊗	⊗	⊗	⊗	■
marmoset -	⊗	○	⊗	■	■	■
galago -	⊗	○	⊗	⊗	■	■
rat -	○	○	⊗	■	⊗	■
mouse -	⊗	○	⊗	■	⊗	■
rabbit -	○	○	○	■	⊗	■
cow -	○	○	○	⊗	○	■
dog -	⊗	○	○	■	⊗	⊗
rfbat -	⊗	○	○	⊗	■	■
shrew -	⊗	○	○	⊗	■	■
armadillo -	⊗	○	○	⊗	○	■
elephant -	⊗	○	○	■	⊗	■
tenrec -	○	○	○	■	⊗	⊗
monodelphis -	○	○	○	■	⊗	■
platypus -	○	○	○	■	⊗	■
chicken -	○	○	○	■	○	■
xenopus -	○	○	○	○	○	⊗
tetraodon -	○	○	⊗	■	○	⊗
zebrafish -	○	○	○	■	⊗	■

History of Pseudogene Preservation

Based on
alignment from
ENCODE MSA
group

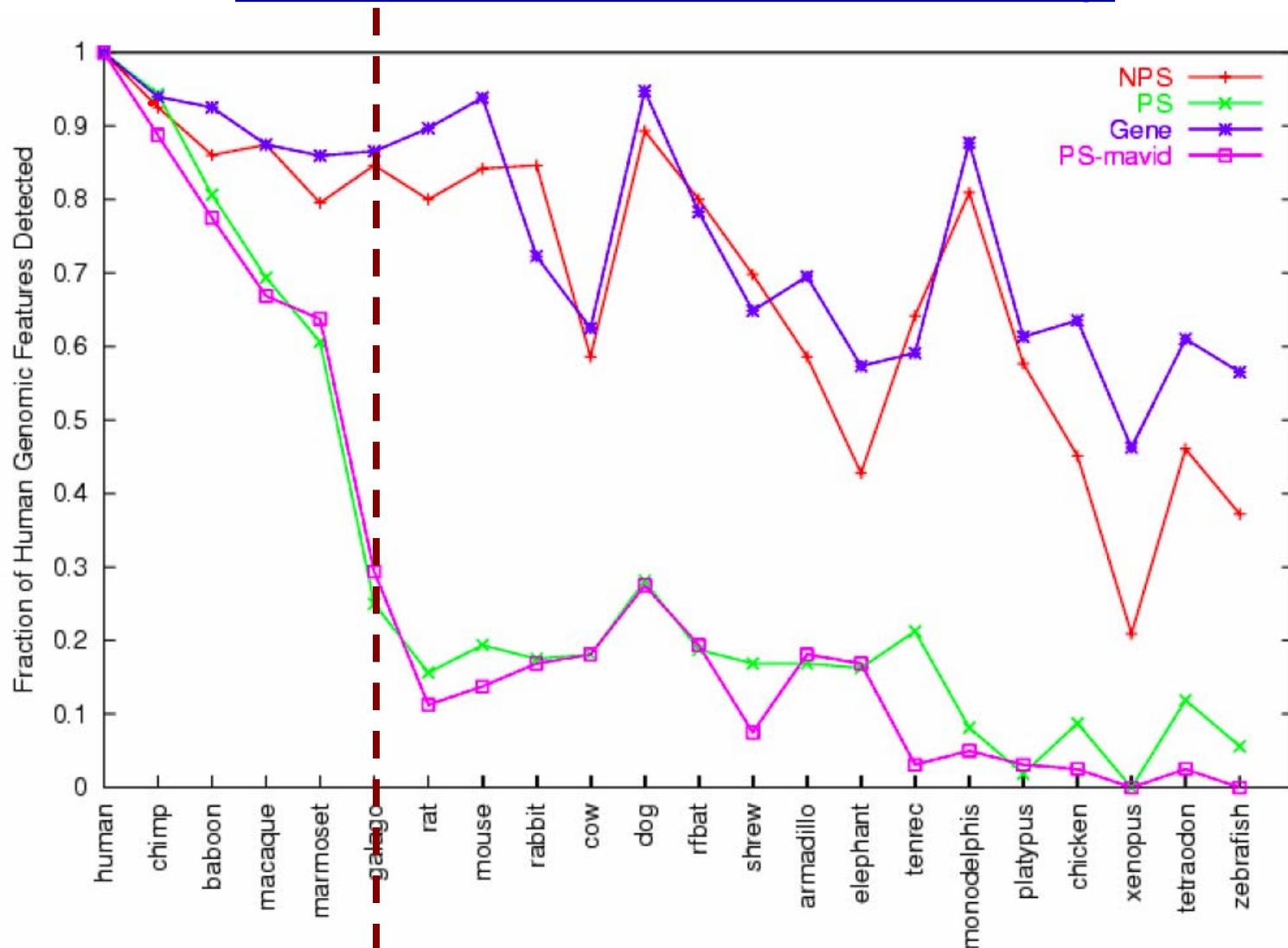
Zheng et al. (2007) Gen. Res.

Absent ○

Present with Disablement ⊗

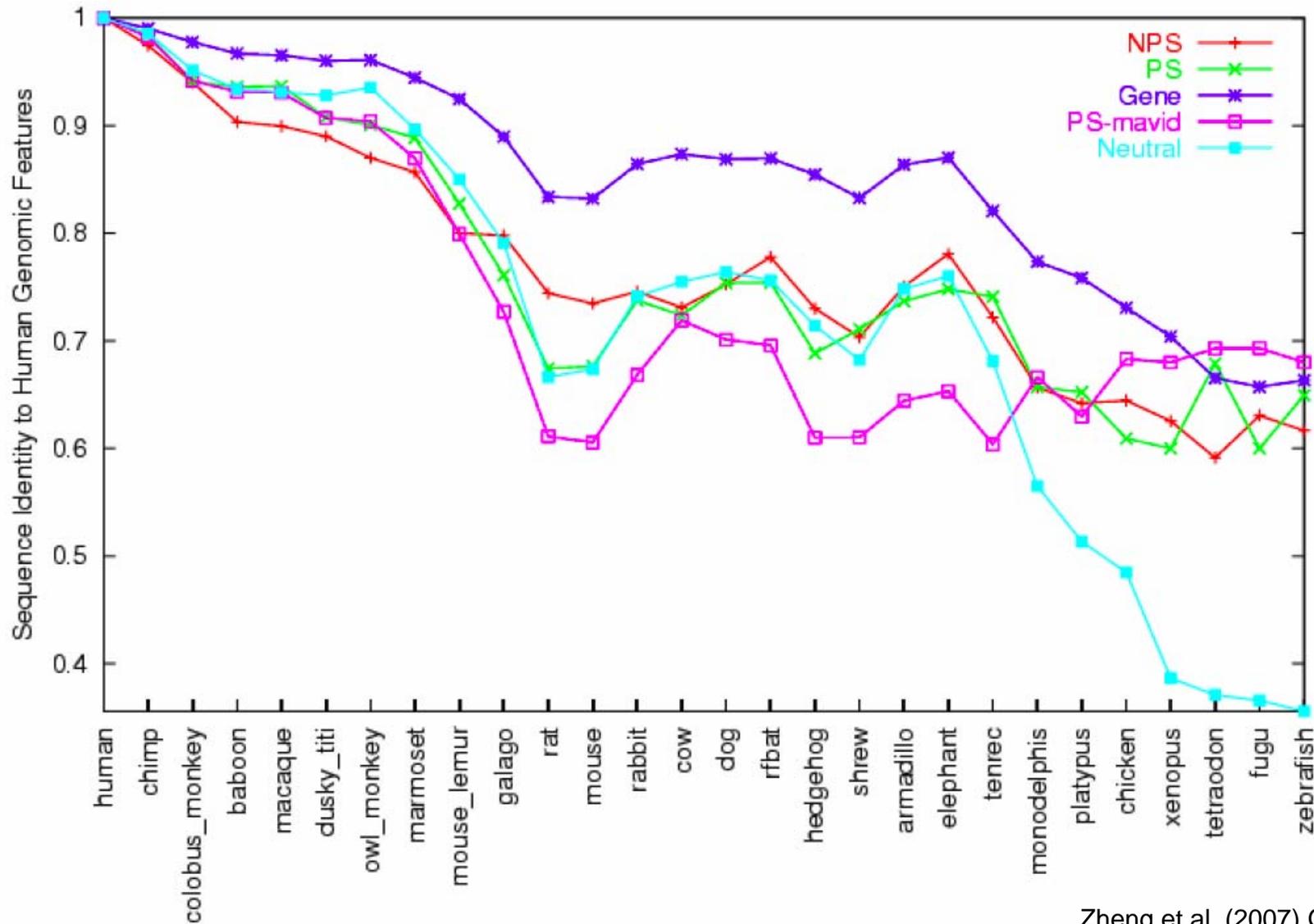
Present without Disablement ■

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

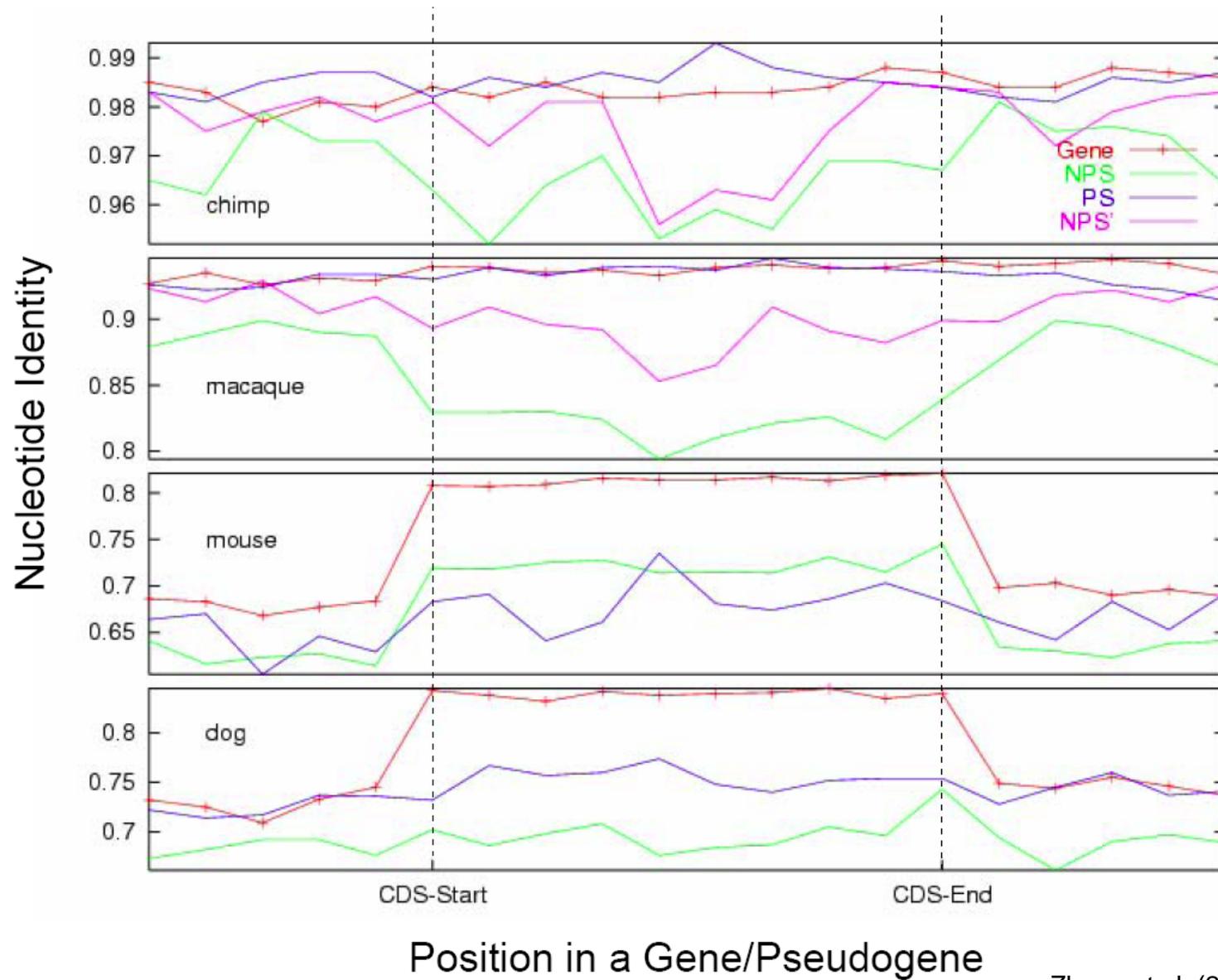


Zheng et al. (2007) Gen. Res.

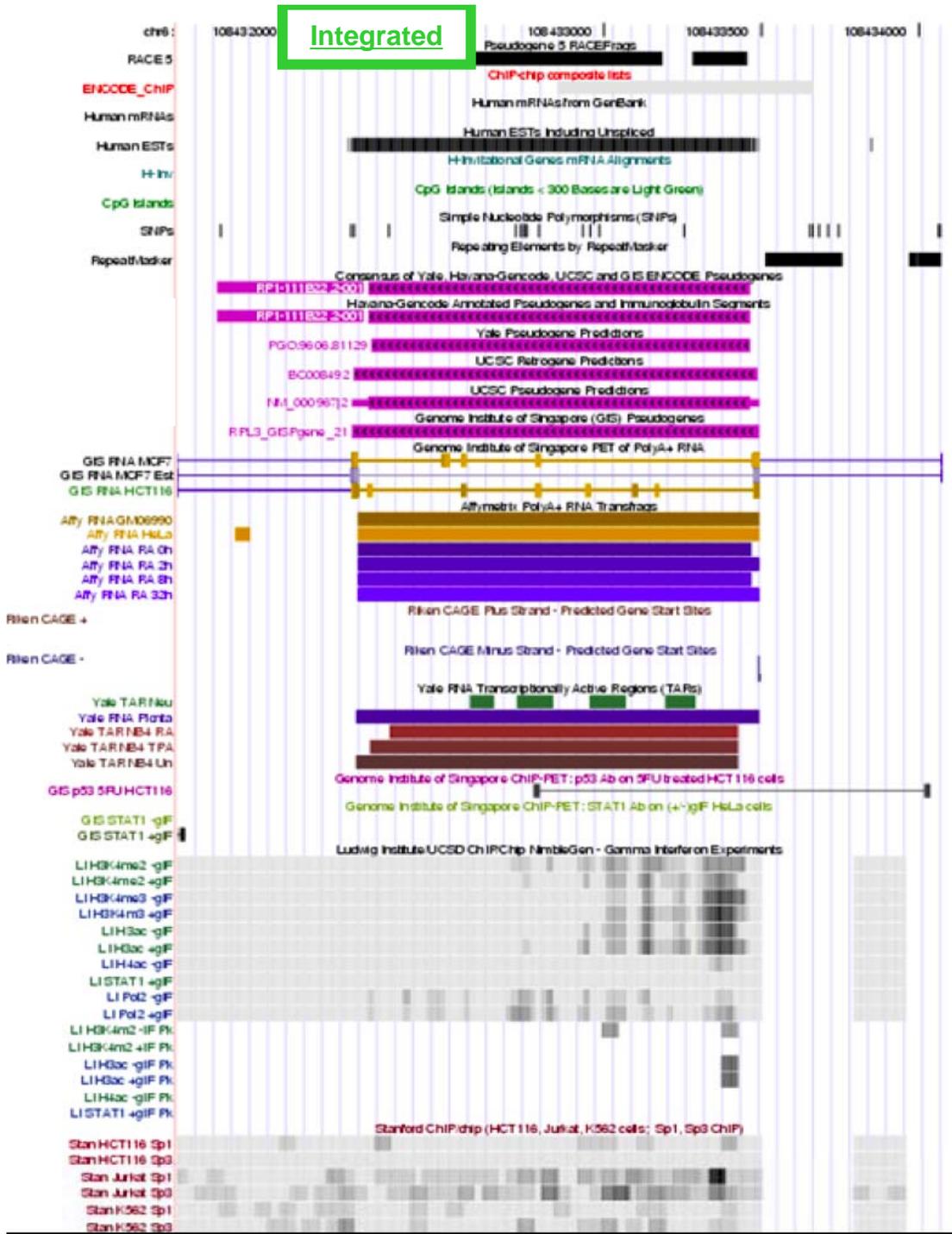
Sequence Decay of Pseudogenes, Approximately Neutral



Sequence Decay of Pseudogenes Relative to their Immediate Genomic Context



Connecting Intergenic Activity to Pseudogenes



Composite
ChIP
hit

Special
ψG
tracks in
browser

diTAG

CAGE

TARs

ChIP-
chip

Connecting TARs (TxFragments) 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.

Zheng et al. (2007) Gen. Res.

Intersection of Pseudogenes with Transcriptional Evidence

	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 

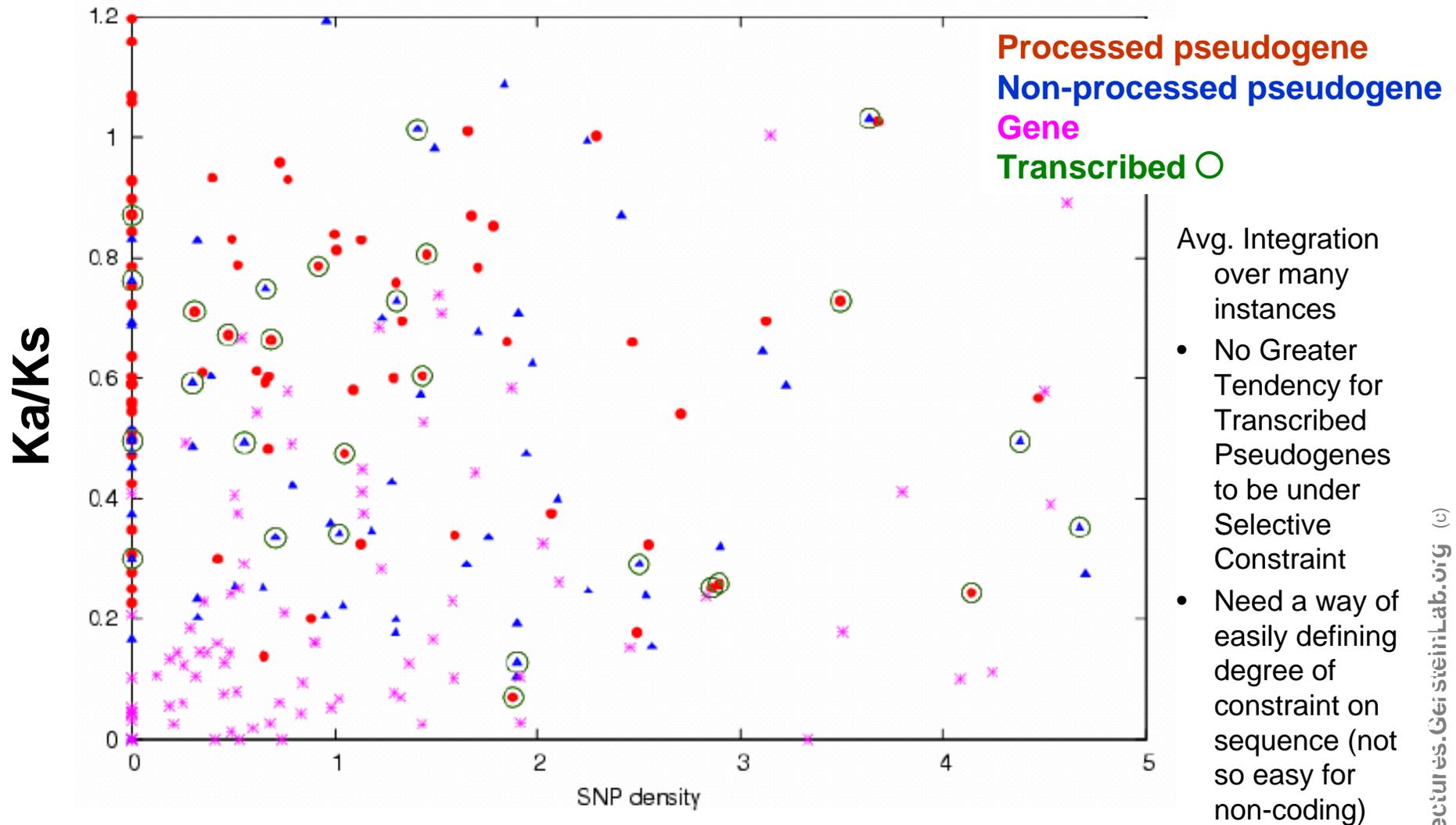
Excluding TARs (due to cross-hyb issues)

Targeted RACE expts to 160 pseudogenes, gives 14

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

Integrated

Integrating Transcriptional Evidence with Gene Annotation and Sequence Constraints



Measurement of Short-time variation (pN+pS)

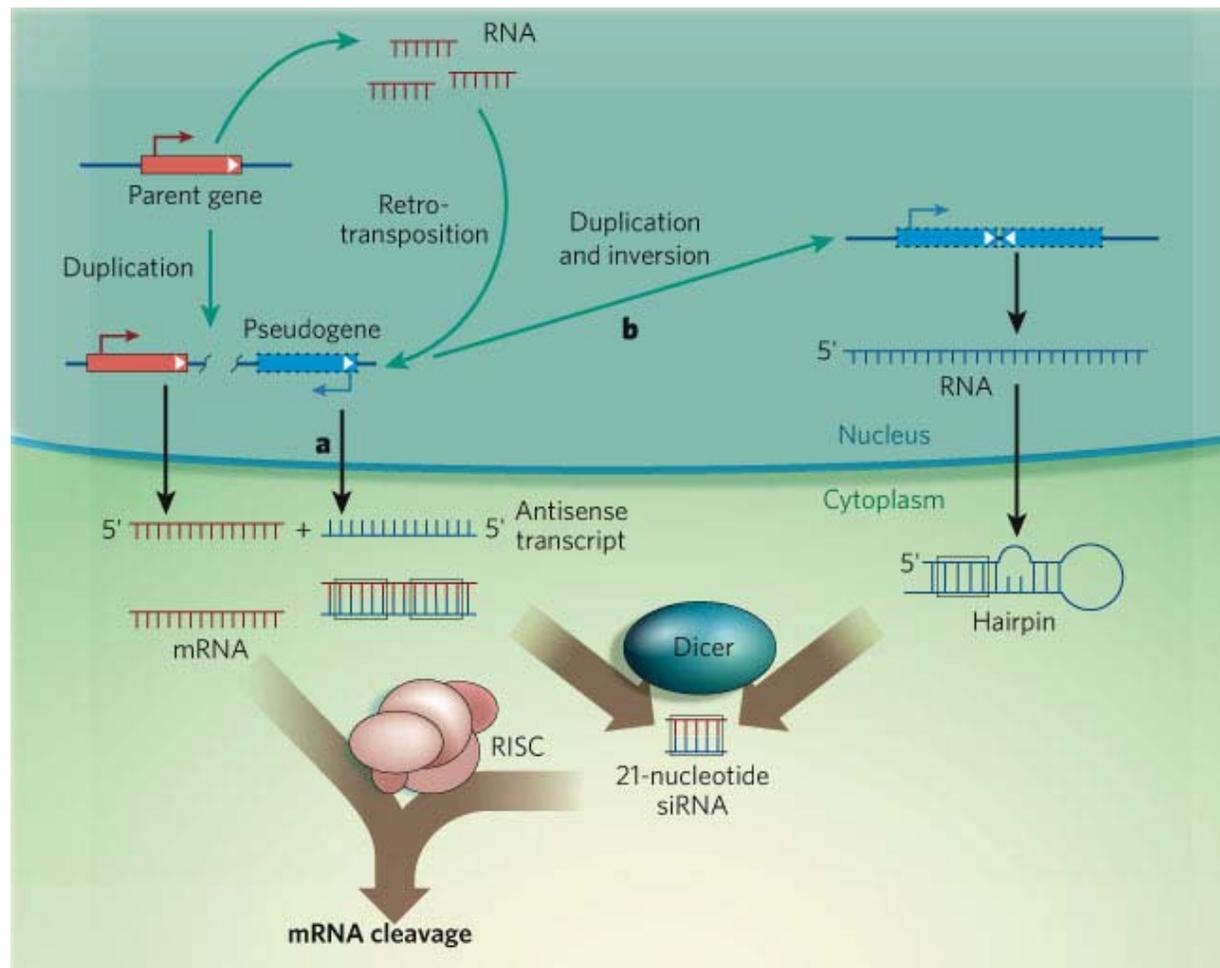
Zheng et al. (2007) Gen. Res.

Conclusion:
The distinction
between gene and
non-gene is
becoming less
clearcut

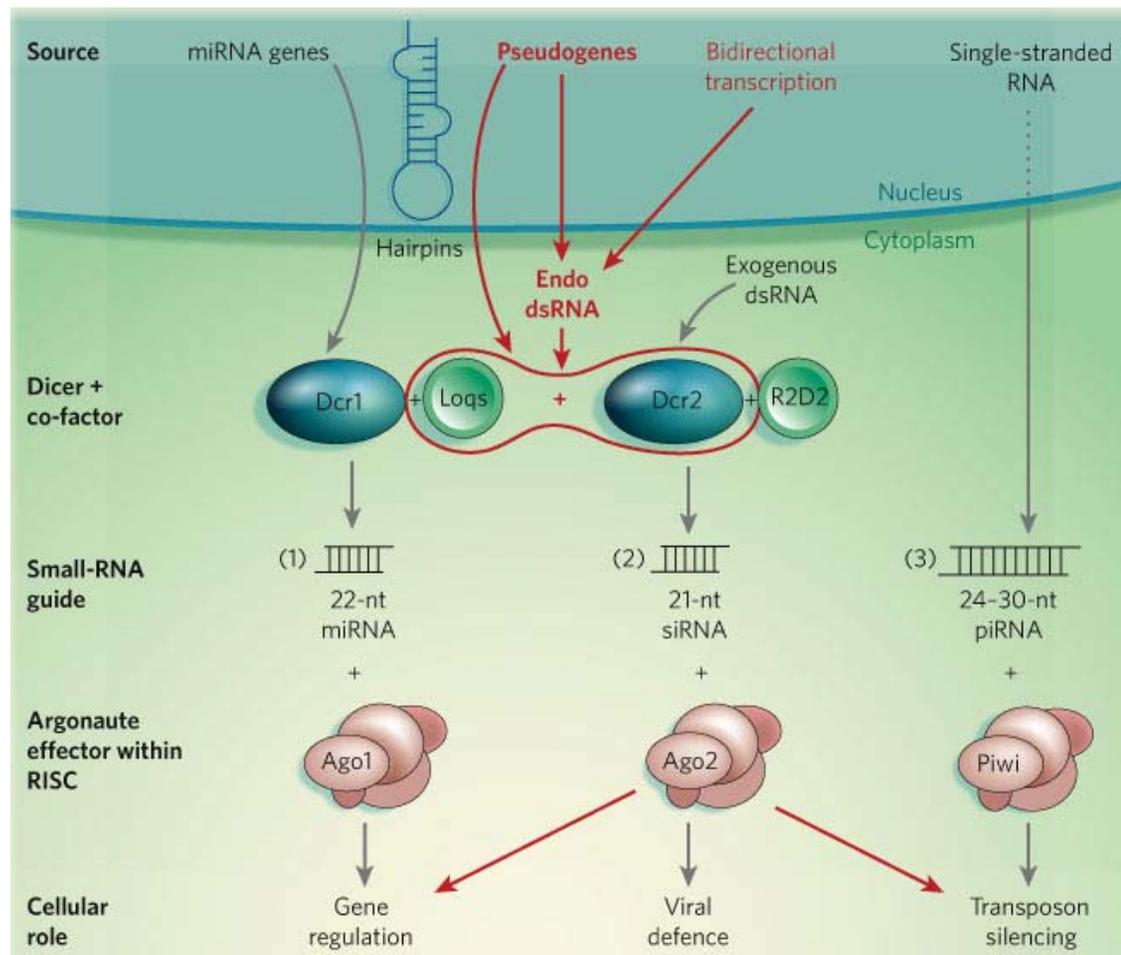
What are Active Pseudogenes Doing? Potential for Gene Regulation via endo-siRNA

- Recent Discovery in Mouse and 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).

How could a pseudogene be involved in RNAi?

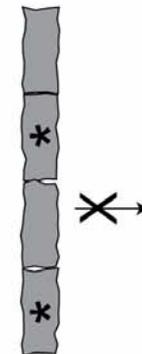
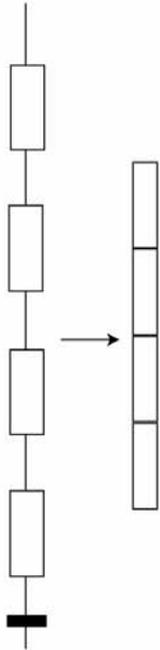


Very Speculatively, Papers Blur Boundaries betw. siRNAs and miRNAs



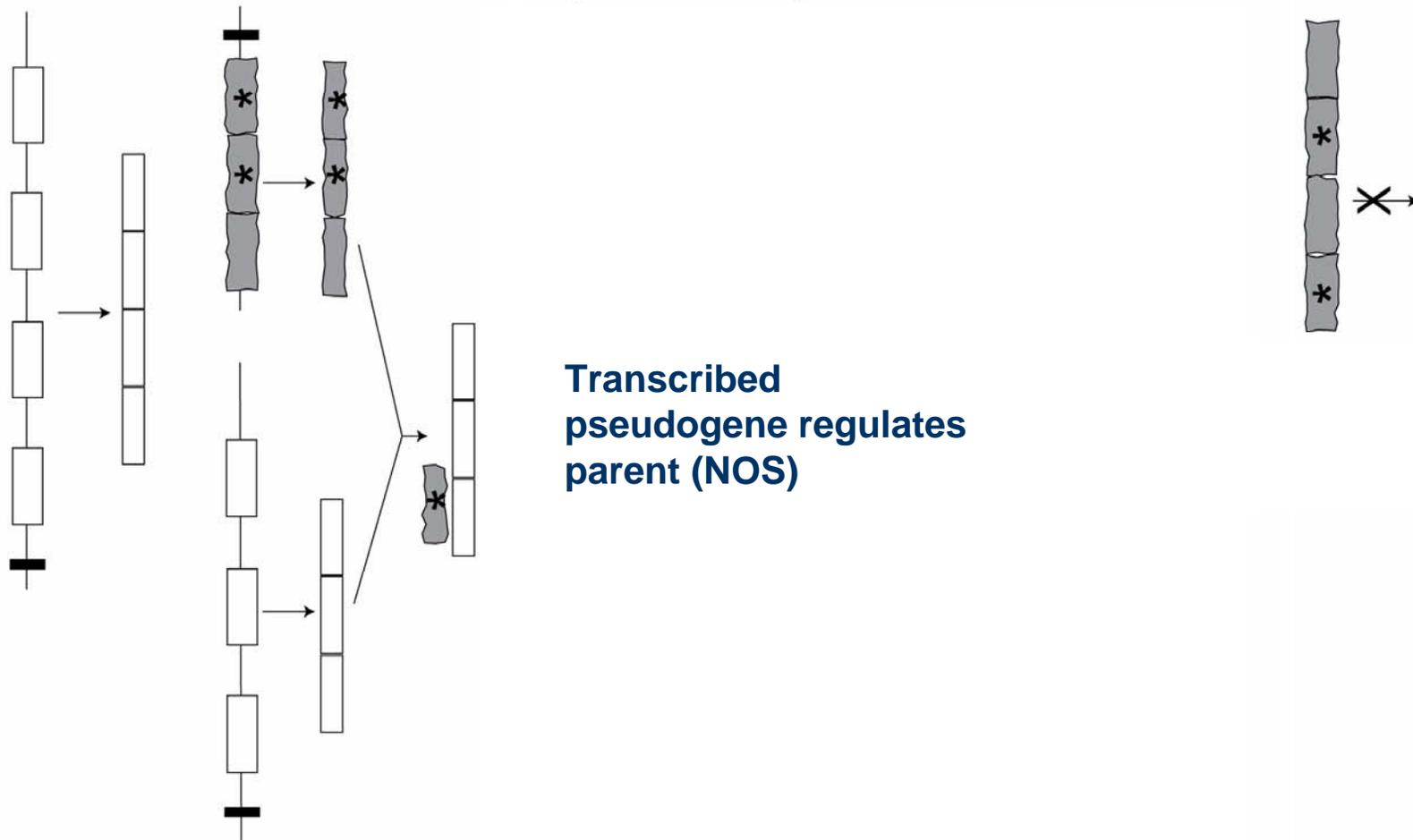
Genes & Pseudogenes

(b) Dead Pseudogene



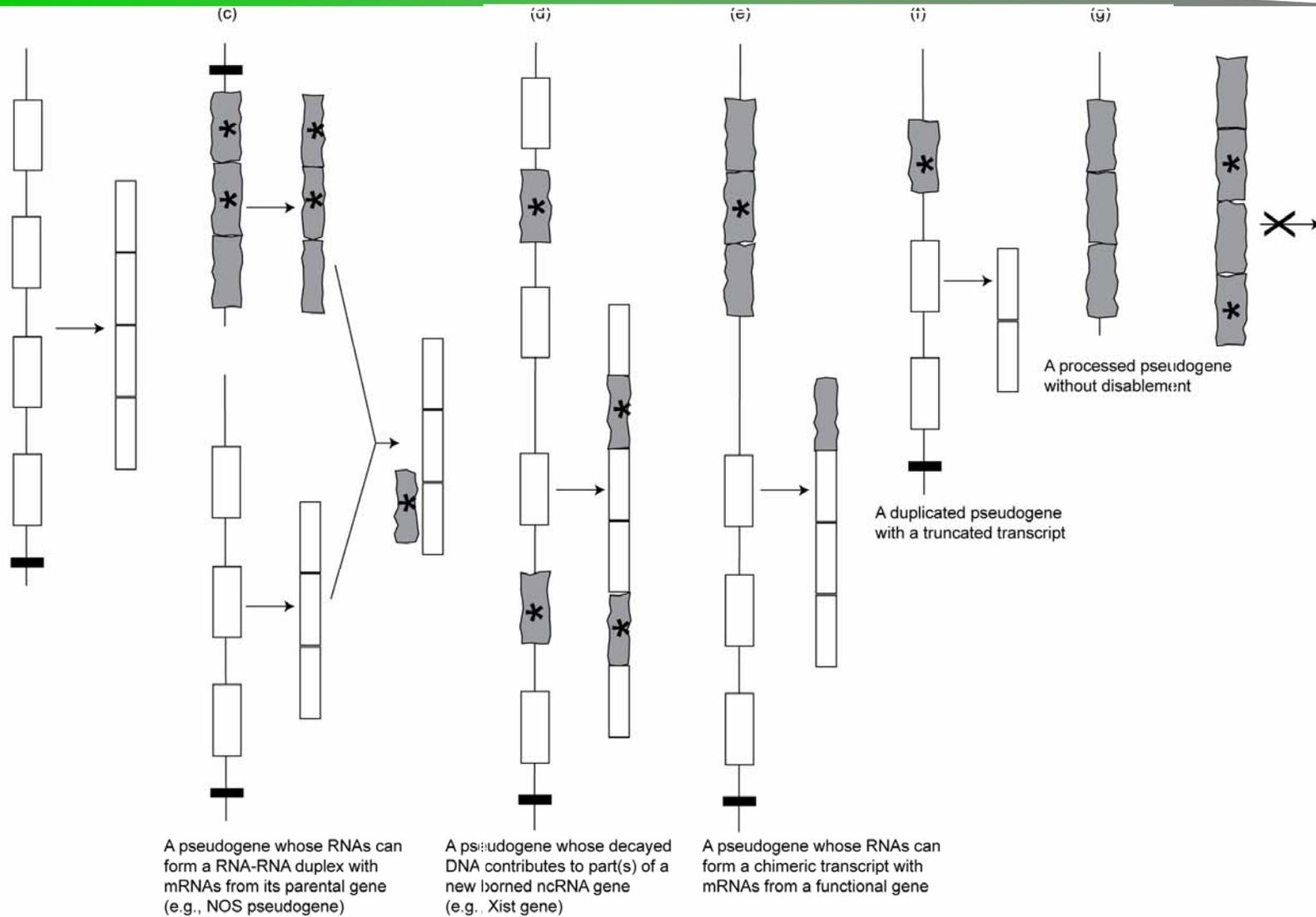
Genes or Pseudogenes?

(b) Dead Pseudogene

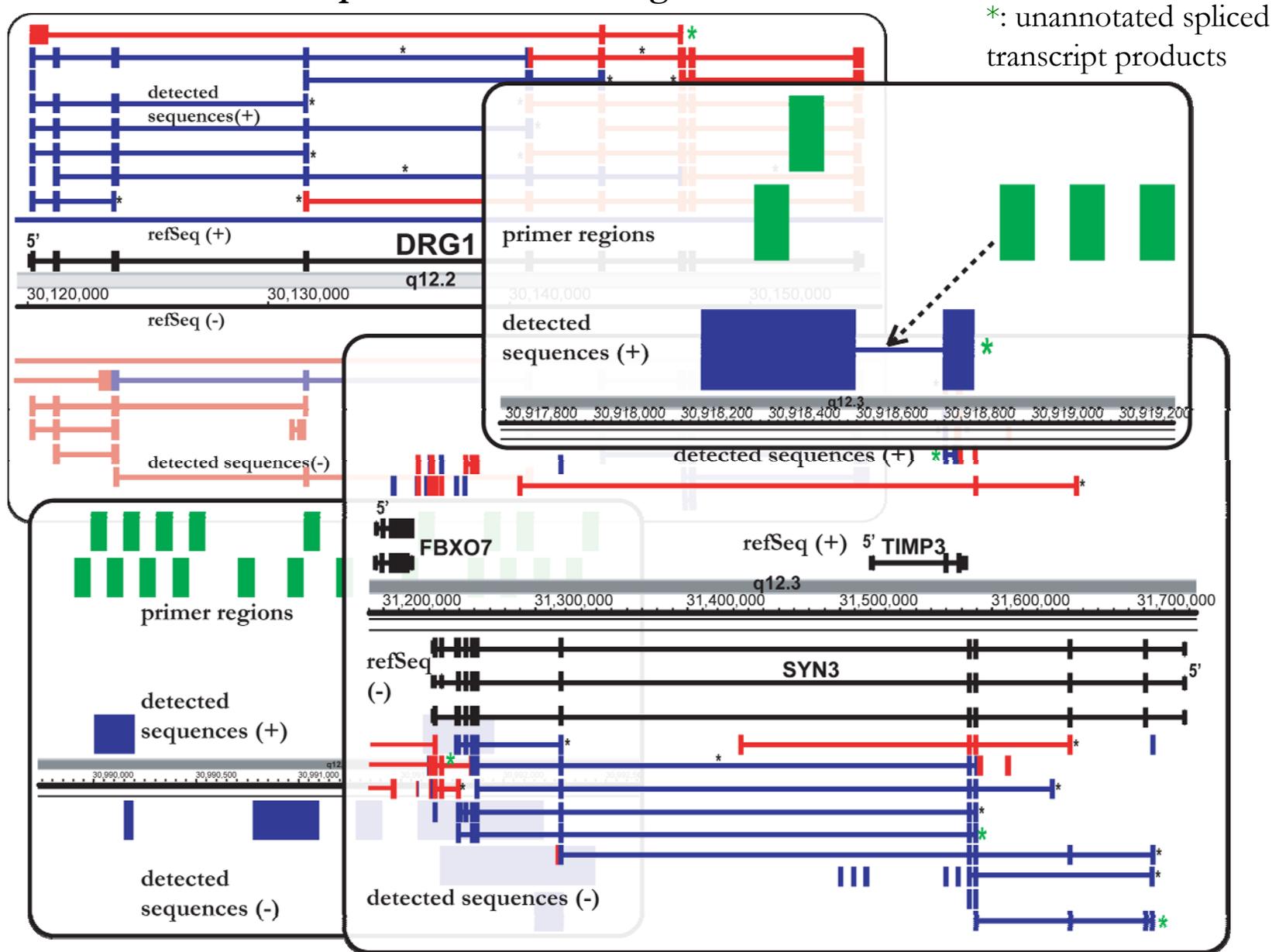


Genes or Pseudogenes?

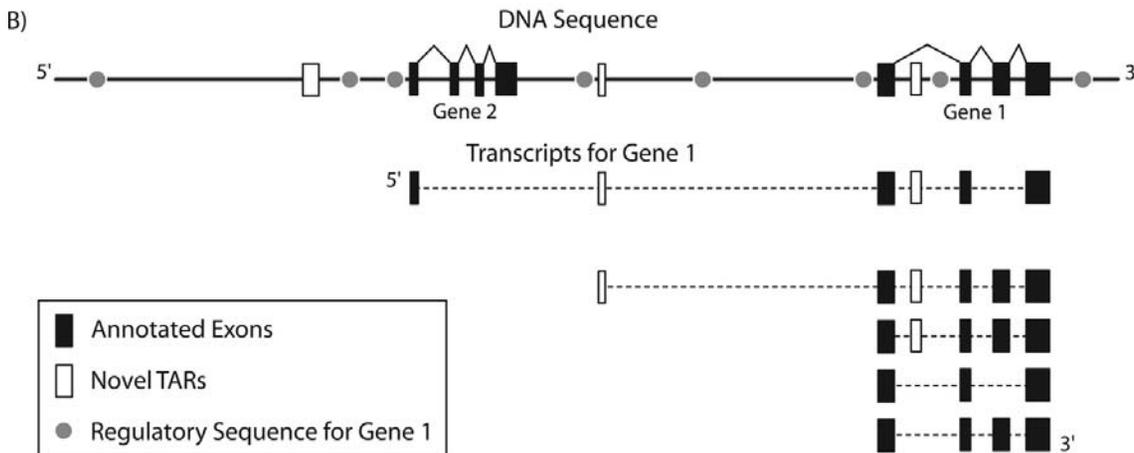
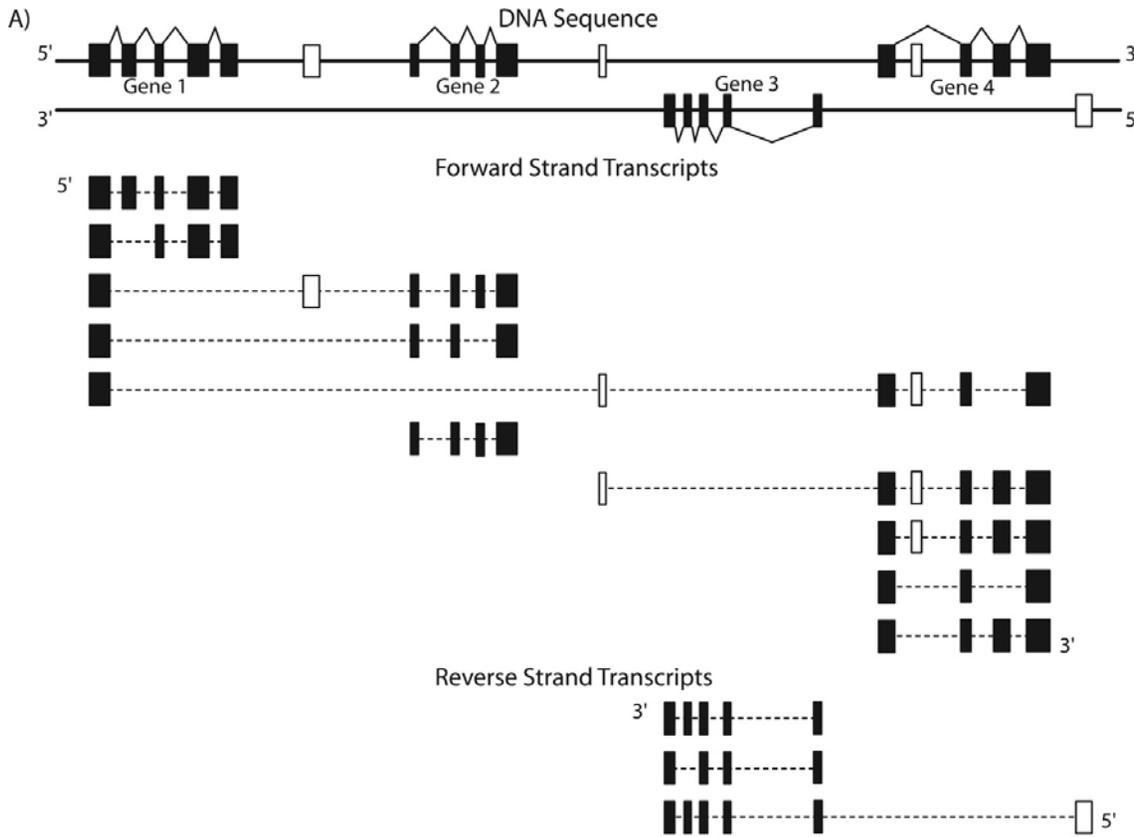
(b) Dead Pseudogene



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



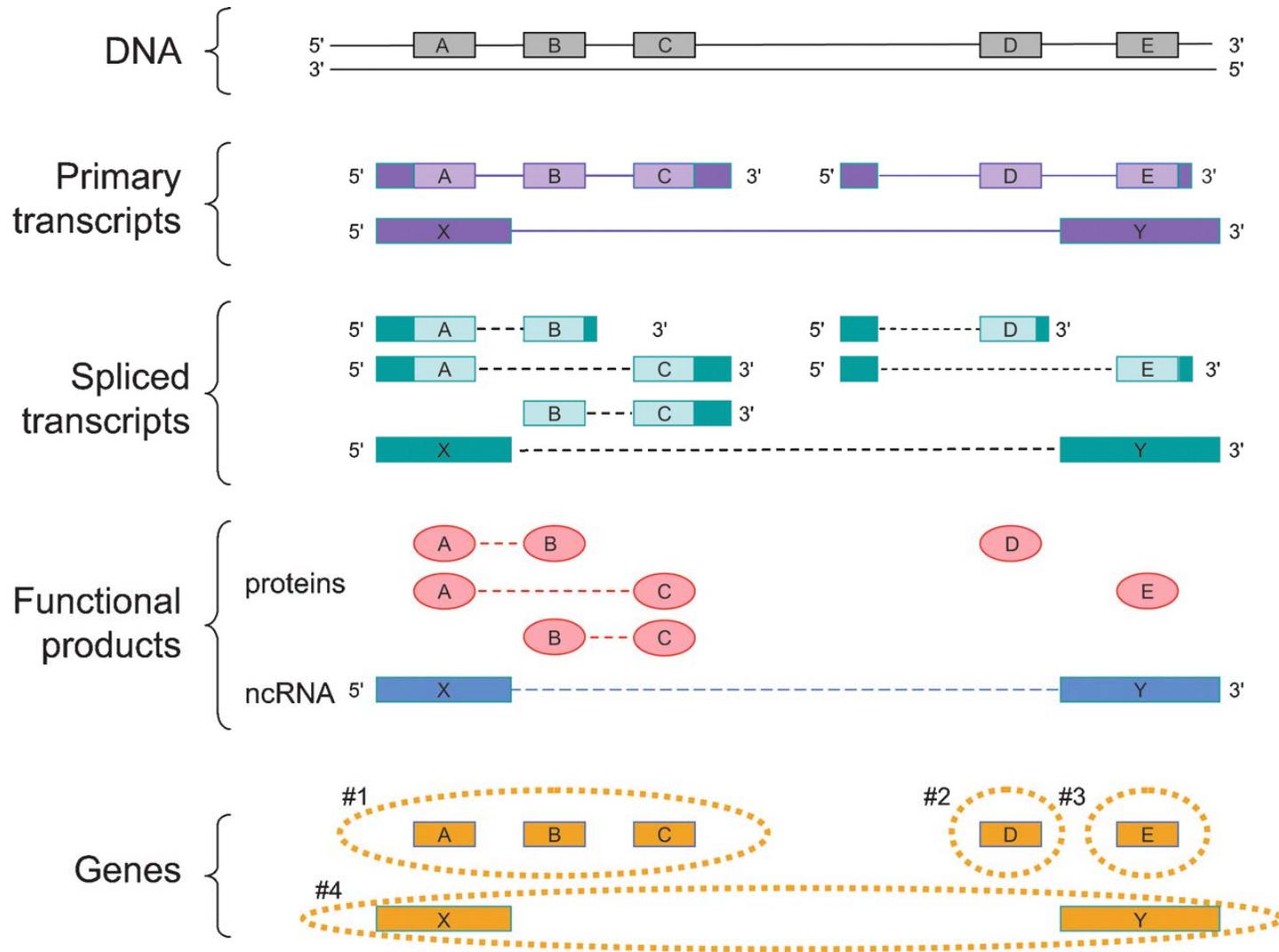
Biological complexity revealed by ENCODE: Long Interleaved Transcripts and Distributed Regulation

What is a Gene? and What is not a Gene?

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

Proposed Re-definition of a Gene: “Gene is a union of genomic sequences encoding a coherent set of potentially overlapping functional products.”

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

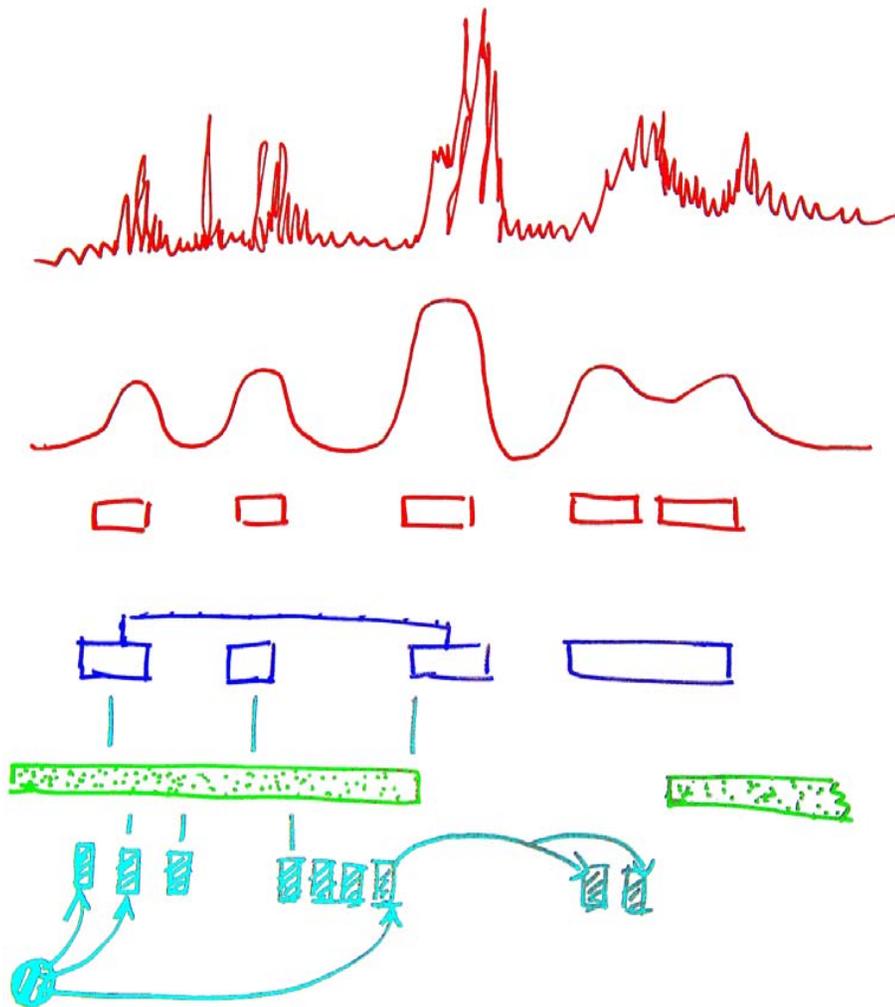




Summary

Overview of Annotation Process

- Doing large-scale similarity comparison, looking for repeated or deleted regions
- **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

Processing the Raw Experimental Signal, developing scoring technology

- Simulating to correct for non-uniform coverage of the genome in Chip-seq experiments and using this to better score the experiments

Large-scale analysis of a single type of
"signal" :
First-Pass Annotation Clustering and
Characterizing Novel Transcribed Regions
and Groups of Binding Sites

- DART classification of TARs
 - ◇ 1300 TARs in ~200 novel ENCODE loci
 - based on expression and phylogenetic clustering
- Deserts and Forests of Binding Activity
 - ◇ on ~50kb scale
 - ◇ Biplot gives broad separation of seq. specific and non-specific factors and associated genomic bins

Integrative Annotation: Relating Pseudogenes to Conservation & Transcription

- Annotation: Pseudogene Assignment
 - ◇ Consensus annotation from automatic pipelines & manual curation gives 201 in ENCODE
 - ~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

ENCODE Acknowledgements

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

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

+

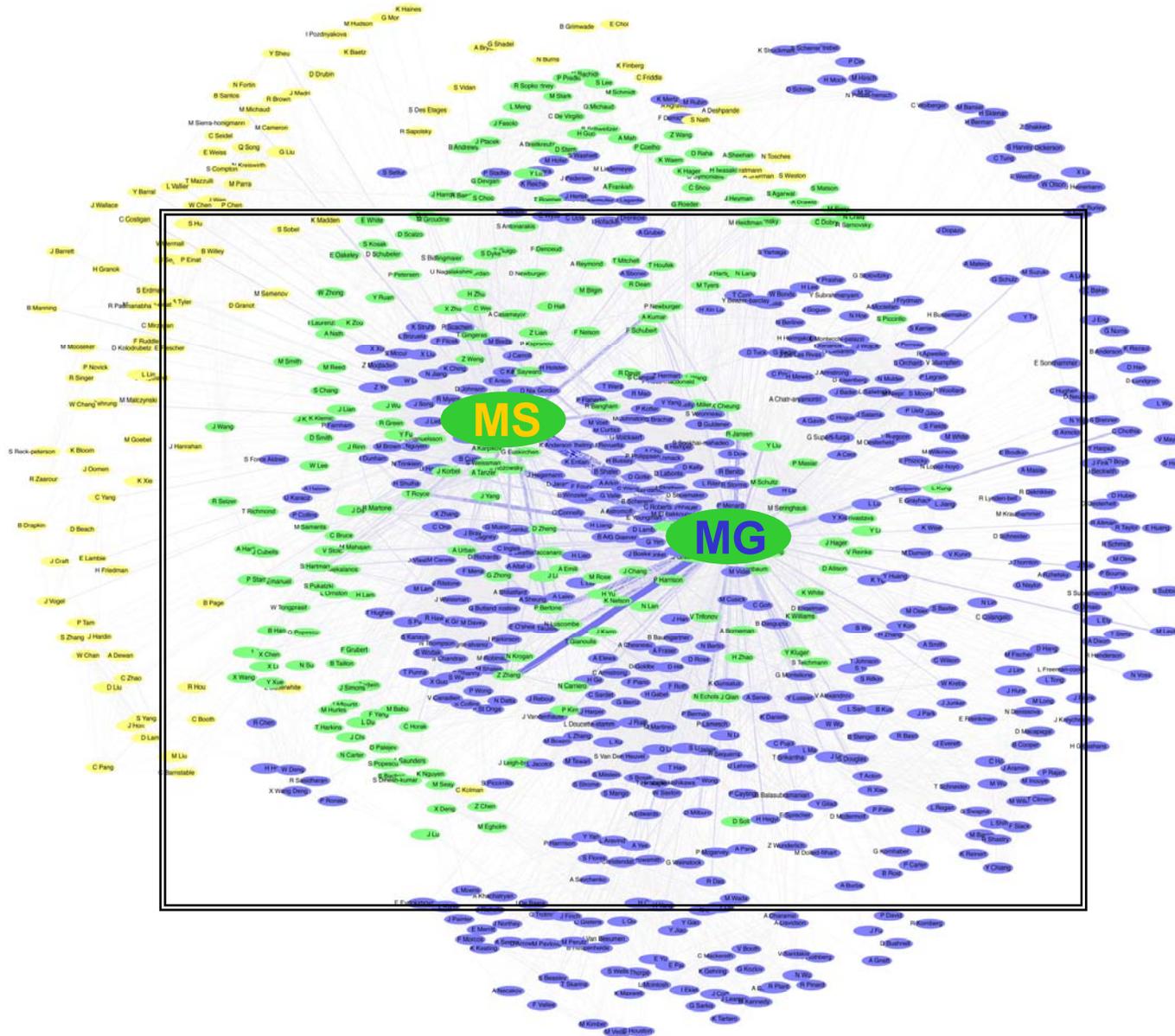
N Trinklein, U Karaöz, A Halees, SF Aldred, PJ Collins, RM Myers

+

Consortium

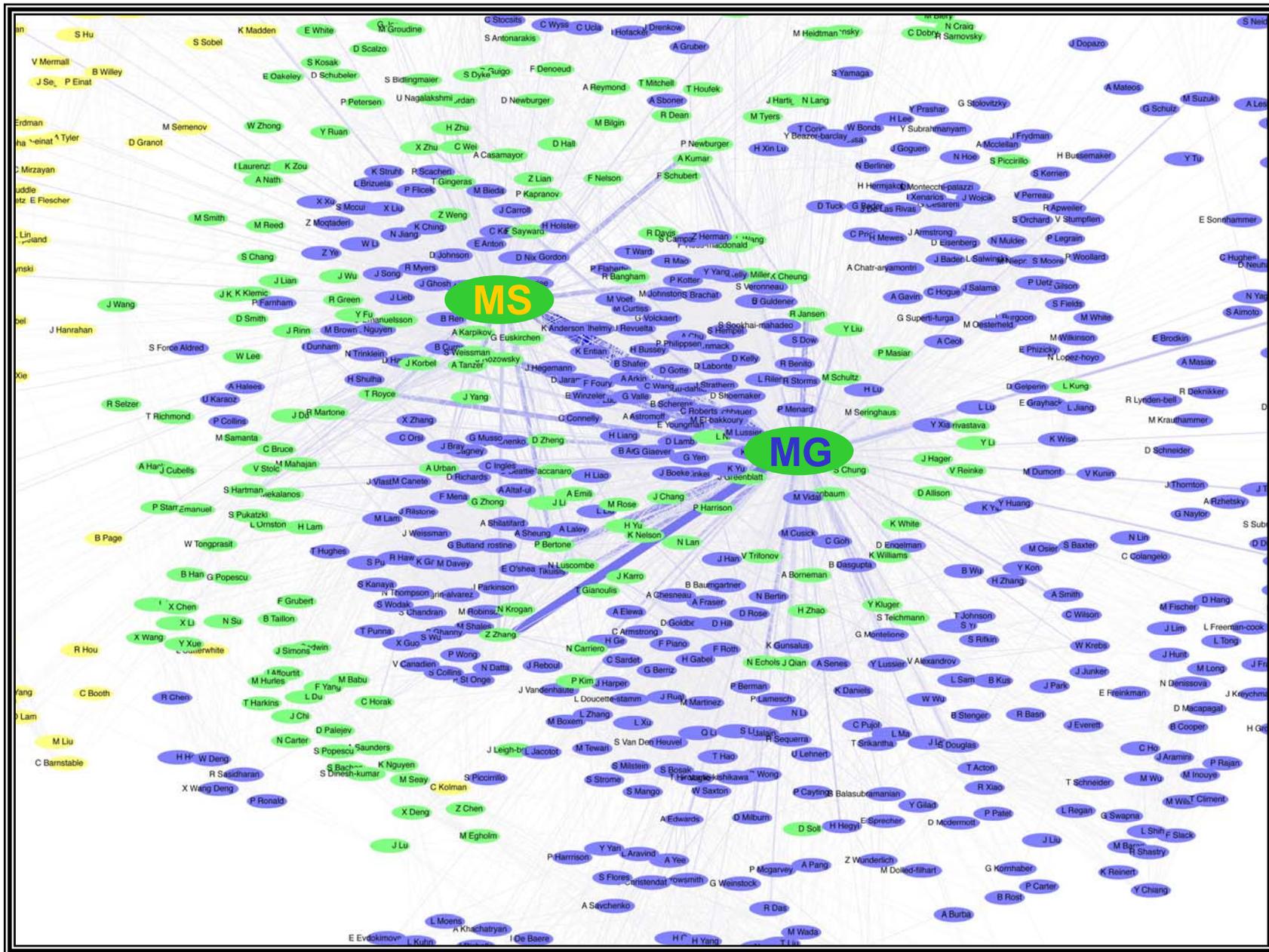
Acknowledgements

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



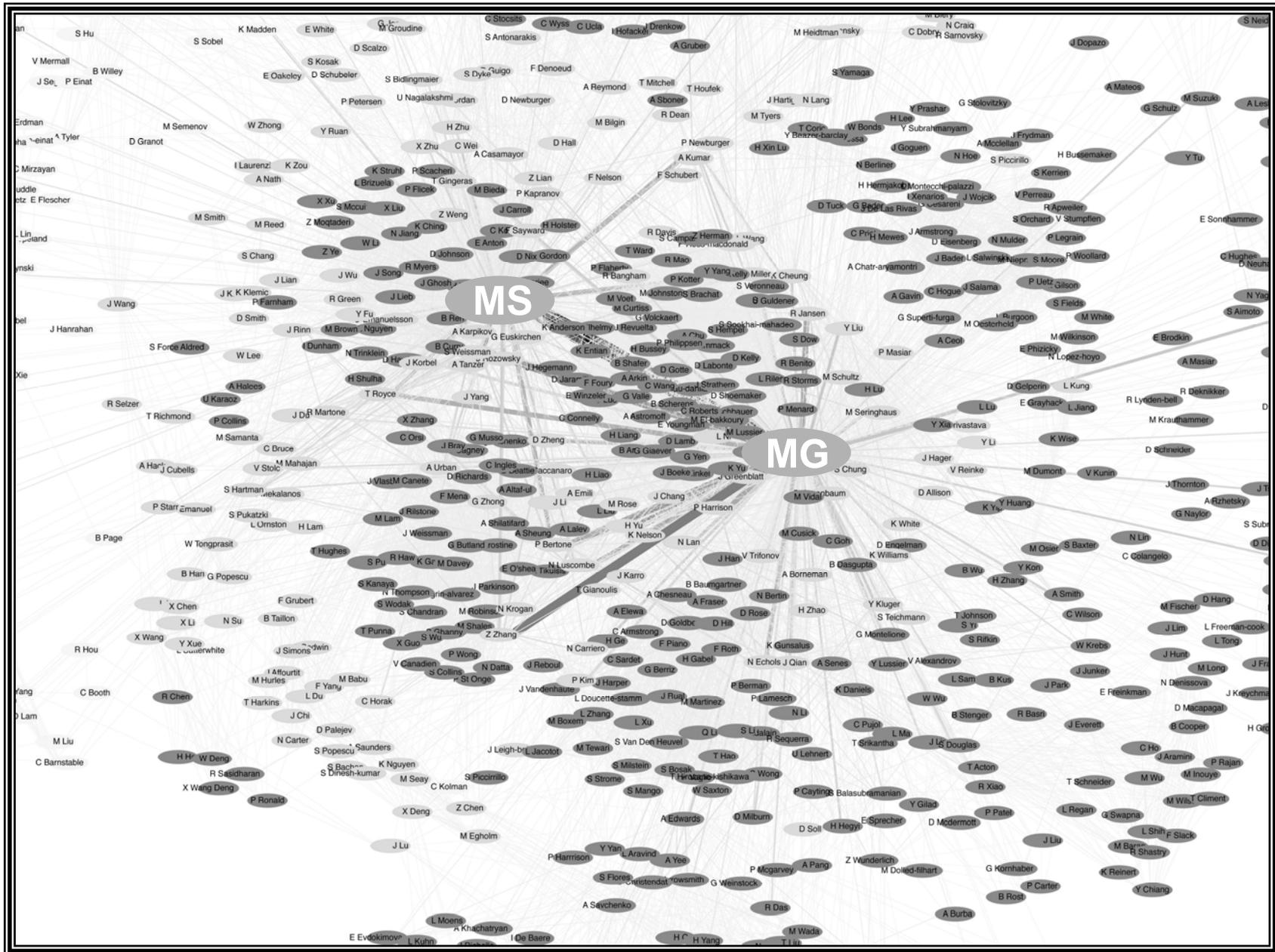
Acknowledgements

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



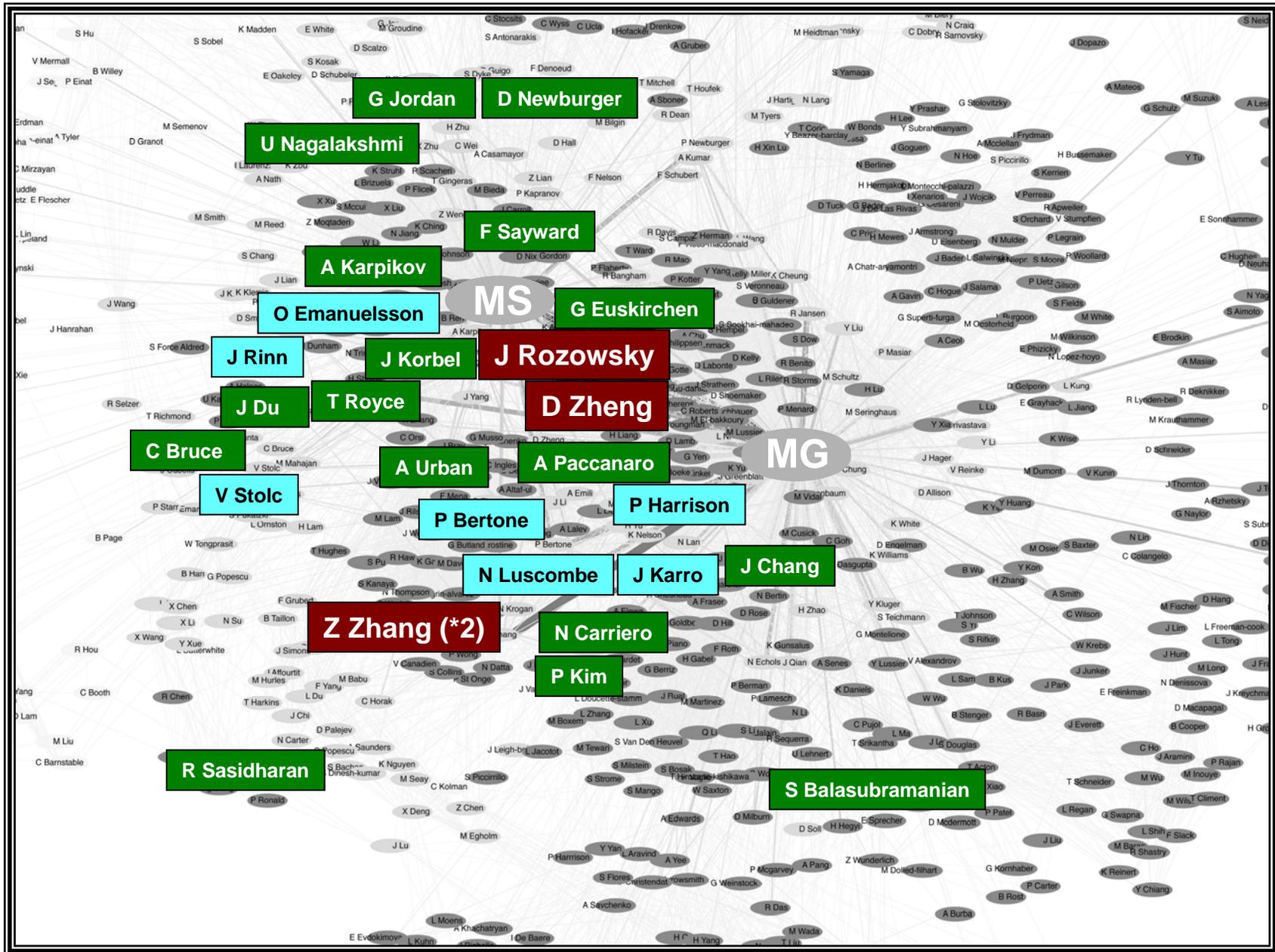
Acknowledgements

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



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).