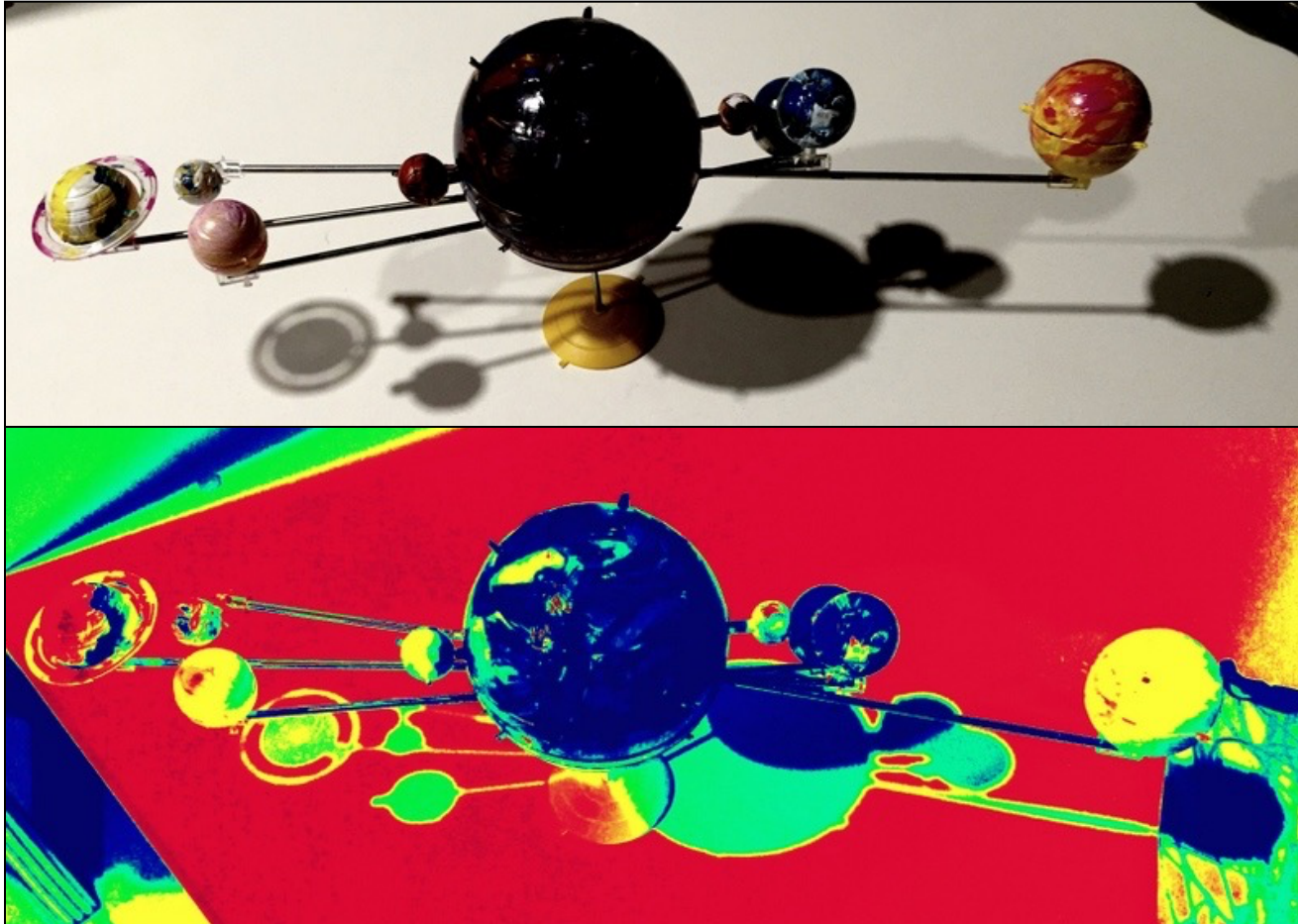# Transcriptome Analysis:
# Large-scale data, high-throughput pipelines & privacy considerations



Mark Gerstein, Yale

Slides freely downloadable from Lectures.GersteinLab.org
& "tweetable" (via @markgerstein). See last slide for more info.

# Large-scale RNA

- Recent advent of many consortia & group producing large scale RNA-seq following on DNA sequencing

- Often this is of human subjects (eg TCGA, PCAWG, GTEx) and needs to be protected

- Useful to build tools & approaches that interact with these data
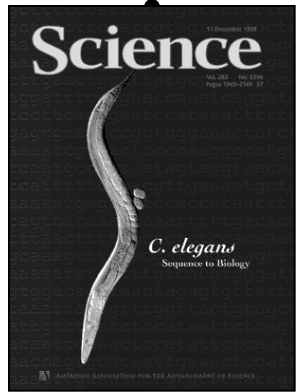
The Human Genome Project

ENCODE Pilot

ENCODE Production

Comparative ENCODE

Epigenome Roadmap
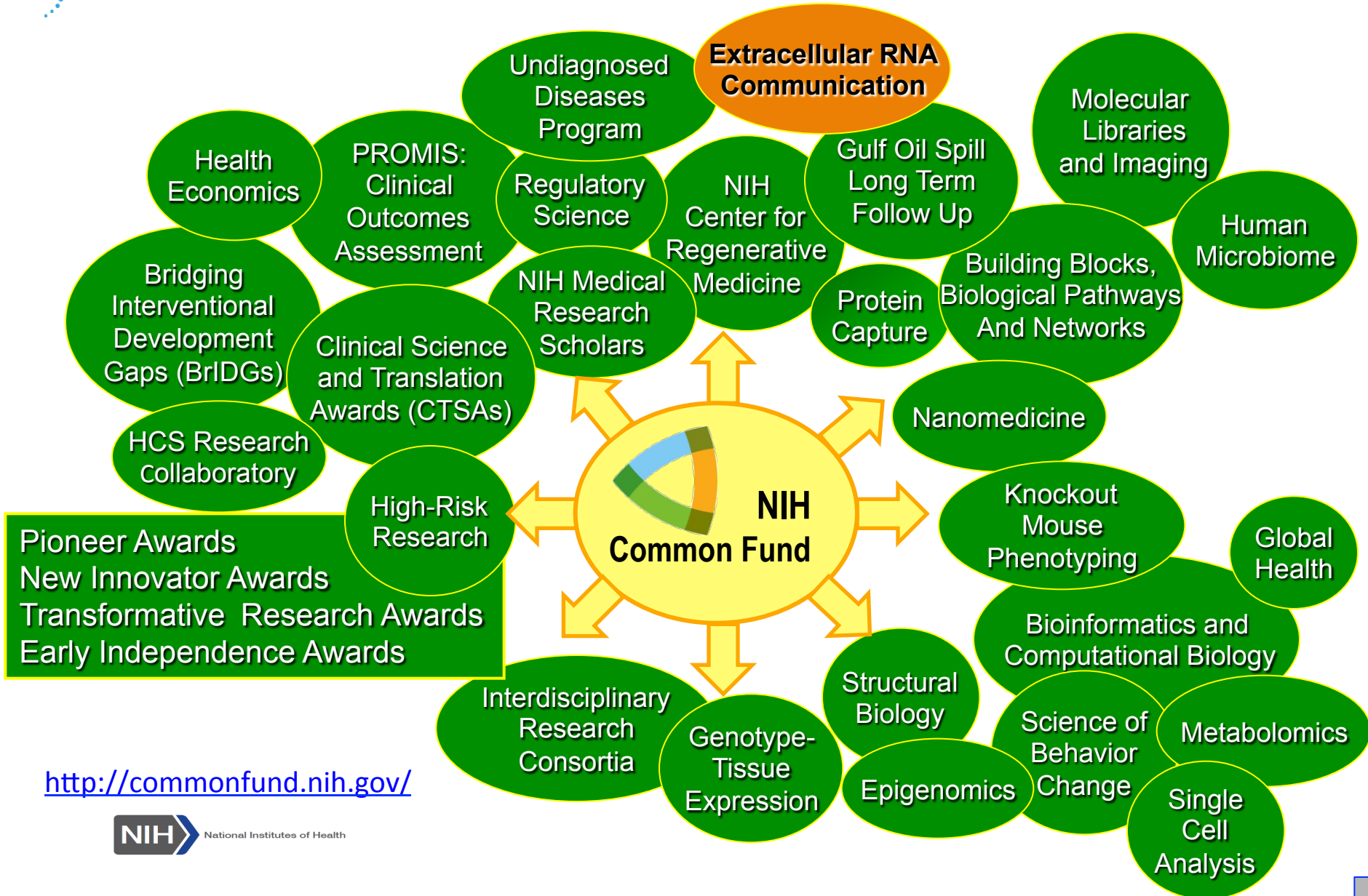
2000  2005  2010  2015

Worm Genome

modENCODE

1000 Genomes Pilot

1000 Genomes Phase 3

GTEx

# Common Fund Programs



http://commonfund.nih.gov/

# ERCC Organization

**Data Management & Resource Repository (DMRR)**

- **Curated Gene Context, Network Modules, Pathways**

**Data Coordination Center (DCC)**

**Aleks Milosavljevic, PI**

**Scientific Outreach Component (SOC) David Galas, PI**

**Data Integration & Analysis Component (DIAC)**

**Mark Gerstein, PI**

**DCC-Admin Core**

**RFA RM-12-011** Reference Profiles

**RFA RM-12-012** Biogenesis, Biodistr, Uptake, and Effector Function

**RFA RM-12-013** Clinical Utility for Biomarkers

**RFA RM-12-014** Clinical Utility for Therapy Development

- **Develop Metadata Stds**
- **Create & Host exRNA Atlas**
- **Data Analysis & Visualizaton**

- **Develop Analysis Pipelines/Tools**
- **Integrative Analysis**
- **Profile/Gene/Network Modules**

# eXRNA Portal



# www.exrna.org

# Regional exRNA Mapping Centers and Data Coordination Center



Pacific Northwest Res. Inst. (Seattle)

U. Mass Med School (Worcester)

UC Michigan

UCSF

Beth Israel (Boston)

UC San Diego

NCBI
Washington, D.C.

Baylor College
of Medicine (Houston):
Data Coordination Center

Bioinformatics Research Laboratory

# RNA-Seq Profiling of Human exRNAs – Multiple Biological Fluids

- **Short & long non-coding, non-coding, circular coding**

- **Plasma**
- **Serum**
- **Urine**
- **Saliva**
- **Cerebrospinal fluid**
- **Cord blood**
- **Seminal fluid**
- **Bronchoaveolar fluid**
- **Placenta**

- **Endogenous vs exogenous (environment/diet)**

# ExRNA Atlas

| Analysis: EXR-DGPLAS00-AN<br>**Status: Protect**<br>Analysis Type: Reference Alignment | Total Mapped Reads | Other Genomic Loci | rRNAs | miRNAs | tRNAs | piRNAs | snoRNAs | Rfam RNAs | Plant and Virus miRNAs |
|---|---|---|---|---|---|---|---|---|---|
| EXR-DGPLAS01-BS<br>SM11_norm1, Plasma, Scientific Control | 6822465 | 46.232% | 51.062% | 2.296% | 0.216% | 0.011% | 0.002% | 0.000% | 0.181% |
| EXR-DGPLAS02-BS<br>SM12_norm2, Plasma, Scientific Control | 6318178 | 47.003% | 51.495% | 1.150% | 0.186% | 0.010% | 0.002% | 0.000% | 0.155% |
| EXR-DGPLAS03-BS<br>SM13_norm3, Plasma, Scientific Control | 5943384 | 48.815% | 49.283% | 1.542% | 0.210% | 0.012% | 0.002% | 0.000% | 0.135% |
| EXR-DGPLAS04-BS<br>SM1_crc1, Plasma, Colorectal Carcinoma | 1490041 | 51.510% | 47.307% | 0.924% | 0.119% | 0.016% | 0.002% | 0.000% | 0.122% |
| EXR-DGPLAS05-BS<br>SM2_crc2, Plasma, Colorectal Carcinoma | 2872815 | 45.518% | 52.513% | 1.643% | 0.191% | 0.018% | 0.003% | 0.000% | 0.114% |
| EXR-DGPLAS06-BS<br>SM3_crc3, Plasma, Colorectal Carcinoma | 3690661 | 46.638% | 51.518% | 1.498% | 0.215% | 0.018% | 0.002% | 0.000% | 0.110% |
| EXR-DGPLAS07-BS<br>SM6_uc1, Plasma, Ulcerative Colitis | 3304333 | 47.004% | 51.877% | 0.661% | 0.345% | 0.020% | 0.006% | 0.000% | 0.088% |
| EXR-DGPLAS08-BS<br>SM7_uc2, Plasma, Ulcerative Colitis | 3613655 | 45.448% | 53.622% | 0.590% | 0.196% | 0.018% | 0.002% | 0.000% | 0.125% |
| EXR-DGPLAS09-BS<br>SM8_uc3, Plasma, Ulcerative Colitis | 4719012 | 44.496% | 51.860% | 3.297% | 0.218% | 0.015% | 0.003% | 0.000% | 0.111% |

# Transcriptome Analysis:
## Large-scale data, high-throughput pipelines & privacy considerations
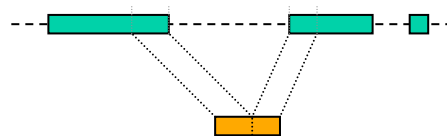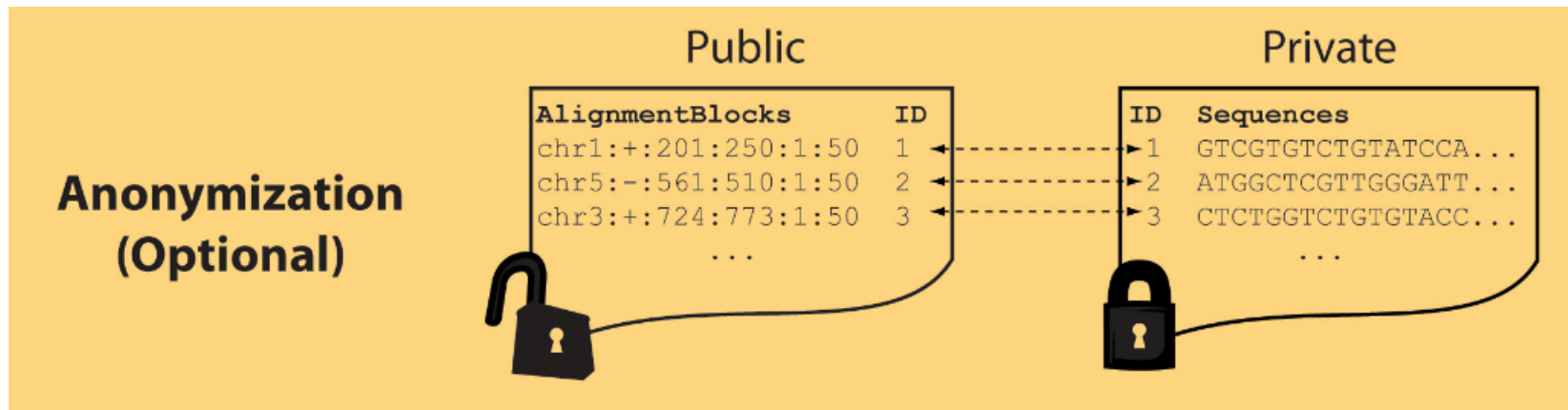
- Large-scale data from consortia #1
  - exRNA.org

- Long-RNA pipeline
  - RSeqTools & anonymized MRF format

- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous

- Large-scale data from consortia #2
  - Co-authorship networks show data flow through key broker individuals

- Large-scale data from consortia #1
  - exRNA.org

- Long-RNA pipeline
  - RSeqTools & anonymized MRF format

- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous

- Large-scale data from consortia #2
  - Co-authorship networks show data flow through key broker individuals

# Light-weight formats

- Some lightweight format clearly separate public & private info., aiding exchange

- Files become much smaller

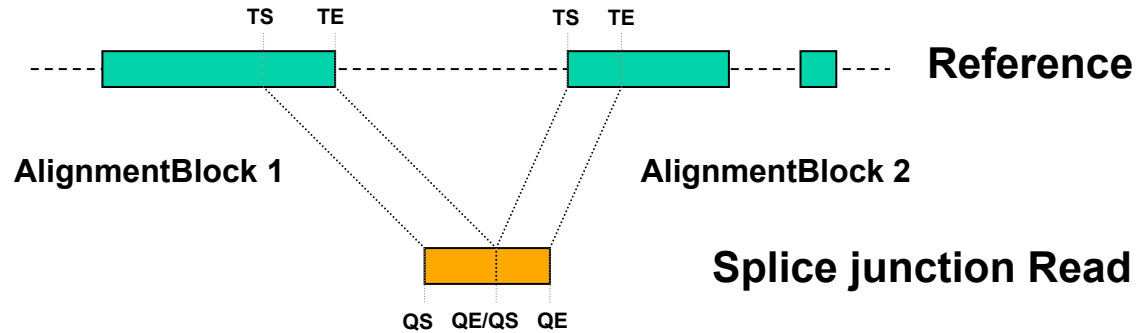- Distinction between formats to compute on and those to archive with – become sharper with big data



Mapping coordinates without variants (MRF)

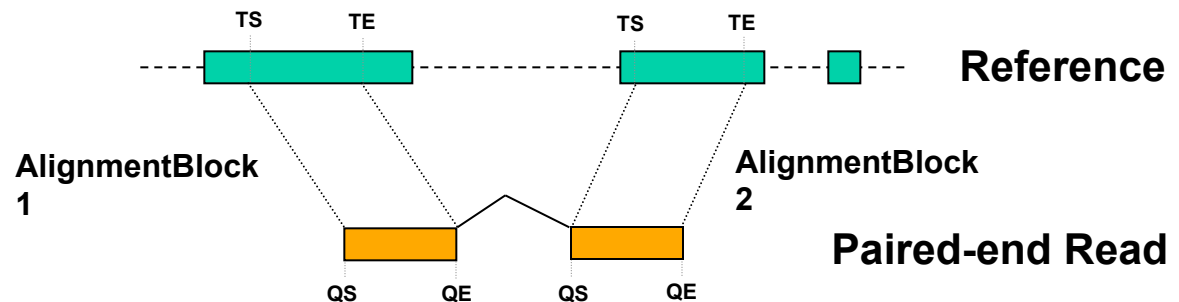Reads (linked via ID, 10X larger than mapping coord.)

# MRF Examples



`chr2:+:601:630:1:30,chr2:+:921:940:31:50`

**Reference**

TS  TE          TS  TE

**AlignmentBlock 1**        **AlignmentBlock 2**

**Splice junction Read**

QS  QE/QS  QE

Legend: TS = TargetStart, TE = TargetEnd, QS = QueryStart, QE = QueryEnd
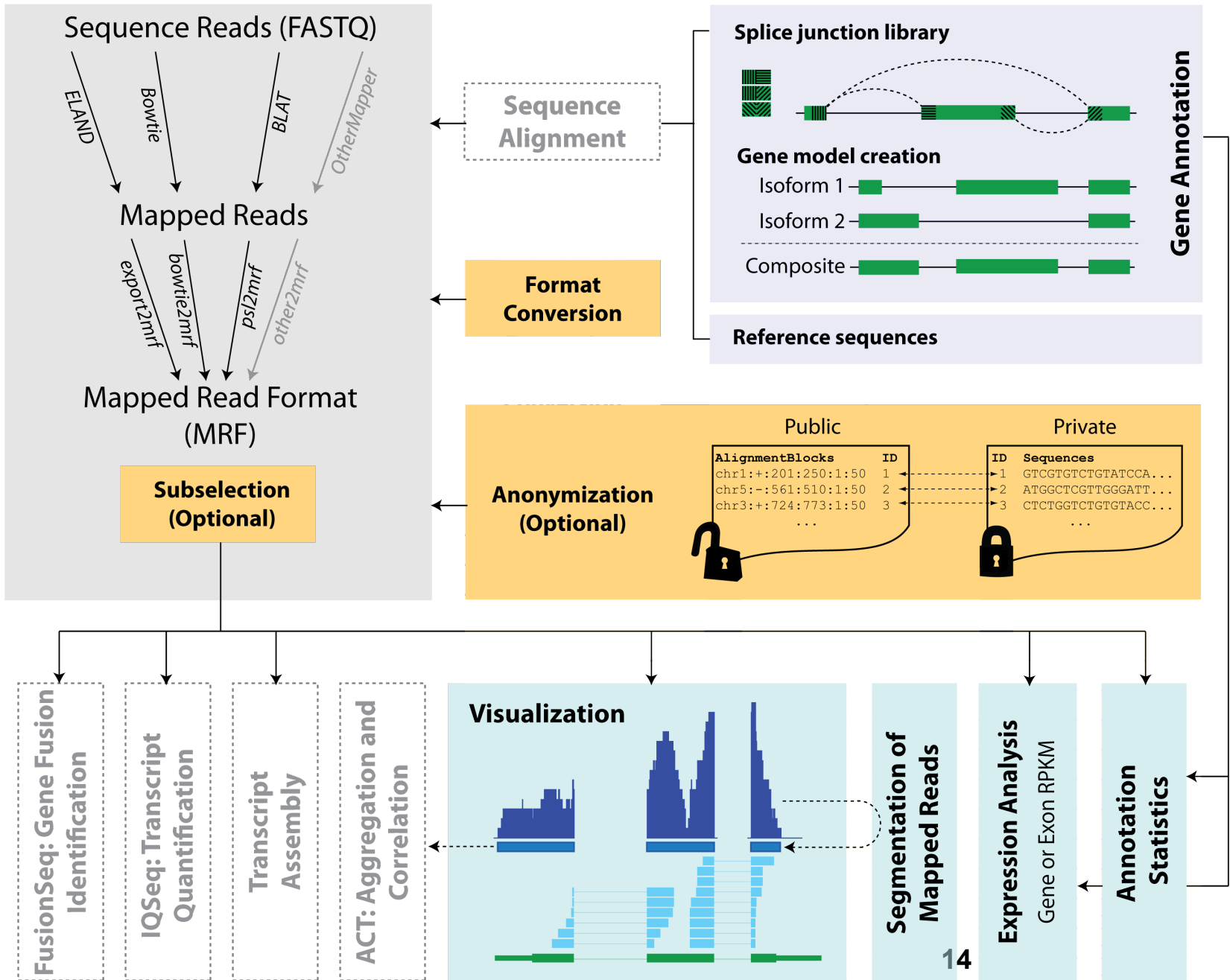
**10X Compression Ex.**

**Raw ELAND** export file has uncompressed file size: ~4 GB; total number of reads: ~20 million; number of mapped reads: ~12 million .

**MRF file** is significantly smaller (~400 MB uncompressed, ~130 MB compressed with gzip).

**BAM file**
has a size of ~1.2 GB.

Reference based compression (ie CRAM) is similar but it stores actual variant beyond just position of alignment block
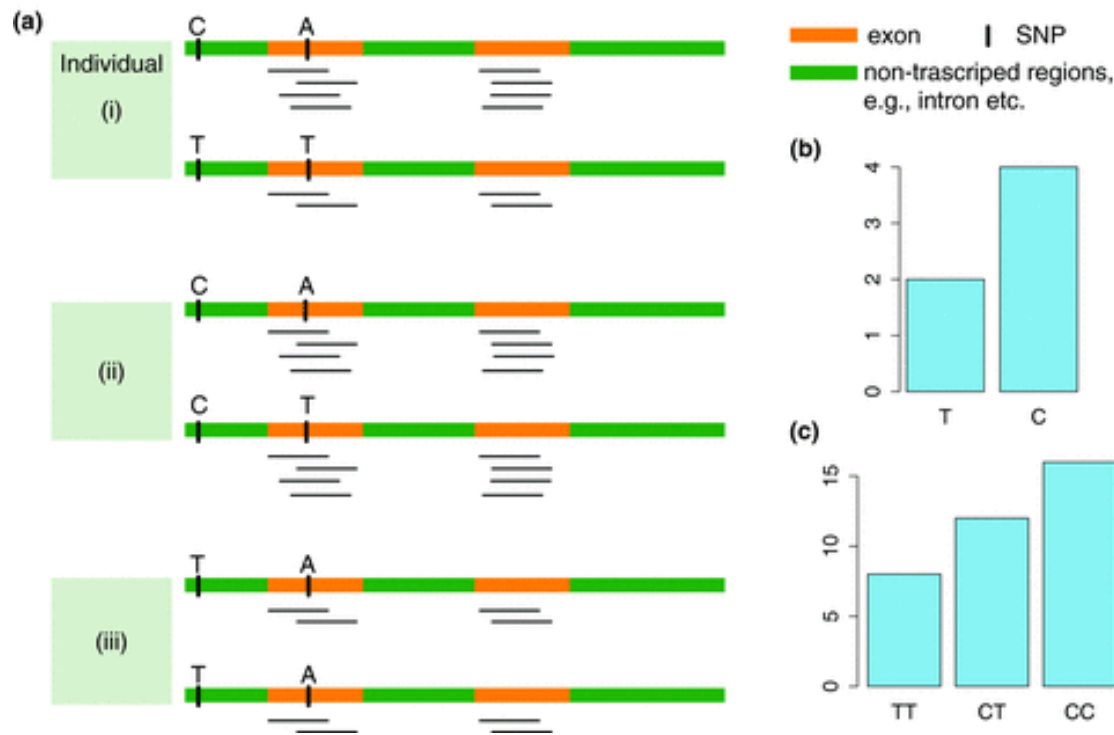
`chr9:+:431:480:1:50|chr9:+:945:994:1:50`

**Reference**

TS       TE          TS       TE

**AlignmentBlock 1**        **AlignmentBlock 2**

**Paired-end Read**

QS      QE      QS      QE

Legend: TS = TargetStart, TE = TargetEnd, QS = QueryStart, QE = QueryEnd

[Habegger et al., Bioinformatics ('11)]

Sequence Reads (FASTQ)

*ELAND*  *Bowtie*  *BLAT*  *OtherMapper*

Mapped Reads

*export2mrf*  *bowtie2mrf*  *psl2mrf*  *other2mrf*

Mapped Read Format (MRF)

**Subselection (Optional)**

**Format Conversion**

**Sequence Alignment**

**Splice junction library**

**Gene model creation**

Isoform 1
Isoform 2
Composite

**Gene Annotation**

**Reference sequences**

**Anonymization (Optional)**

Public

| AlignmentBlocks | ID |
|---|---|
| chr1:+:201:250:1:50 | 1 |
| chr5:-:561:510:1:50 | 2 |
| chr3:+:724:773:1:50 | 3 |
| ... | |

Private

| ID | Sequences |
|---|---|
| 1 | GTCGTGTCTGTATCCA... |
| 2 | ATGGCTCGTTGGGATT... |
| 3 | CTCTGGTCTGTGTACC... |
| | ... |

**FusionSeq: Gene Fusion Identification**

**IQSeq: Transcript Quantification**

**Transcript Assembly**

**ACT: Aggregation and Correlation**

**Visualization**

**Segmentation of Mapped Reads**

**Expression Analysis** Gene or Exon RPKM

**Annotation Statistics**

14

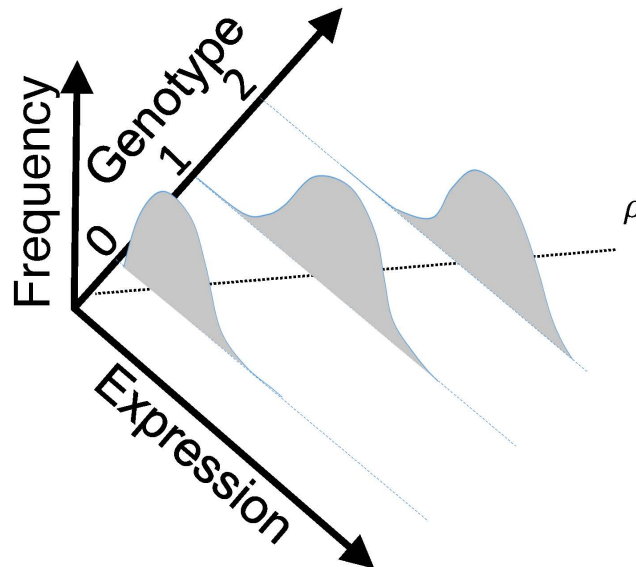# The Genboree Workbench: Web-based Data Management & Analysis

- Large-scale data from consortia #1
  - exRNA.org

- Long-RNA pipeline
  - RSeqTools & anonymized MRF format

- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous

- Large-scale data from consortia #2
  - Co-authorship networks show data flow through key broker individuals

# eQTL Mapping Using RNA-Seq Data



(a)

Individual (i), (ii), (iii)

exon | SNP
non-trascriped regions, e.g., intron etc.

(b)

(c)

[*Biometrics 68(1) 1–11*]

Genotype 0 1 2
Frequency
Expression
$\rho$

- eQTLs are genomic loci that contribute to variation in mRNA expression levels

- eQTLs provide insights on transcription regulation, and the molecular basis of phenotypic outcomes

- eQTL mapping can be done with RNA-Seq data

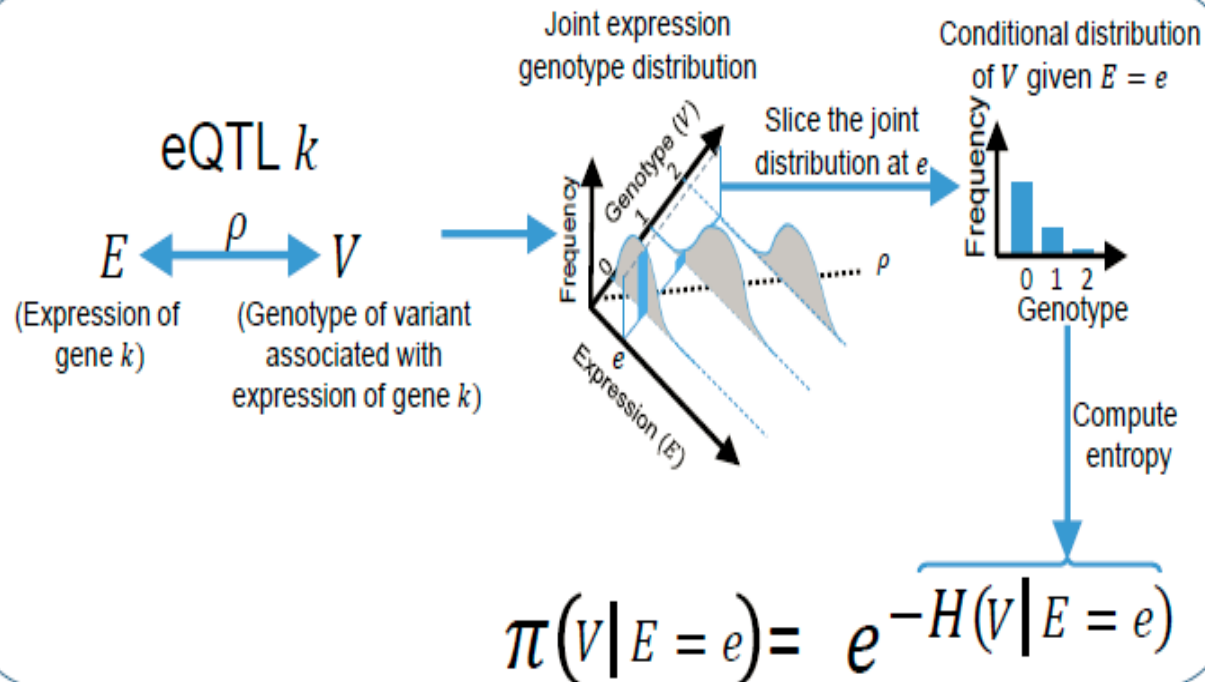# Information Content and Predictability



$$ICI \begin{pmatrix} \text{Individual has variant} \\ \text{genotypes } g_1, g_2, ..., g_n \\ \text{for variants } V_1, V_1, ..., V_n \end{pmatrix} = \log\left(\cfrac{1}{\text{Frequency of } V_1 \text{ genotype}}\right) + \log\left(\cfrac{1}{\text{Frequency of } V_2 \text{ genotype}}\right) + ... + \log\left(\cfrac{1}{\text{Frequency of } V_n \text{ genotype}}\right)$$

$g_1 = 2$  $g_2 = 1$  $g_n = 2$

$V_1$ genotype frequencies    $V_2$ genotype frequencies    $V_n$ genotype frequencies

eQTL $k$

$$E \xleftrightarrow{\rho} V$$

(Expression of gene $k$)    (Genotype of variant associated with expression of gene $k$)

Joint expression genotype distribution

Slice the joint distribution at $e$

Conditional distribution of $V$ given $E = e$
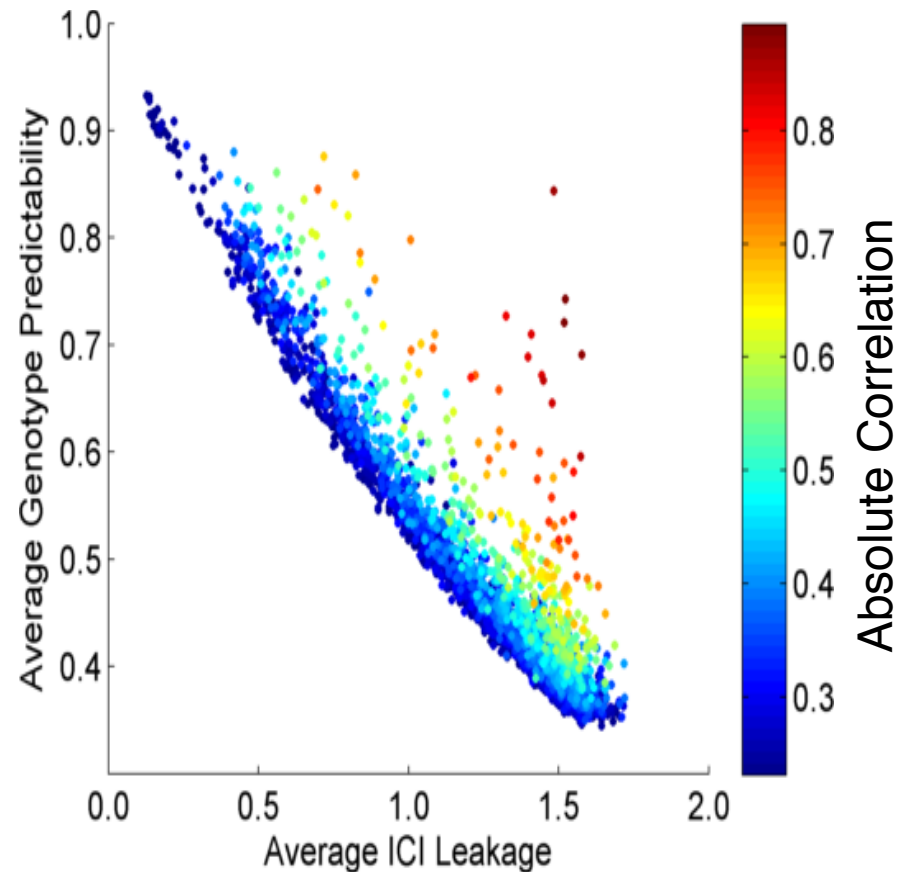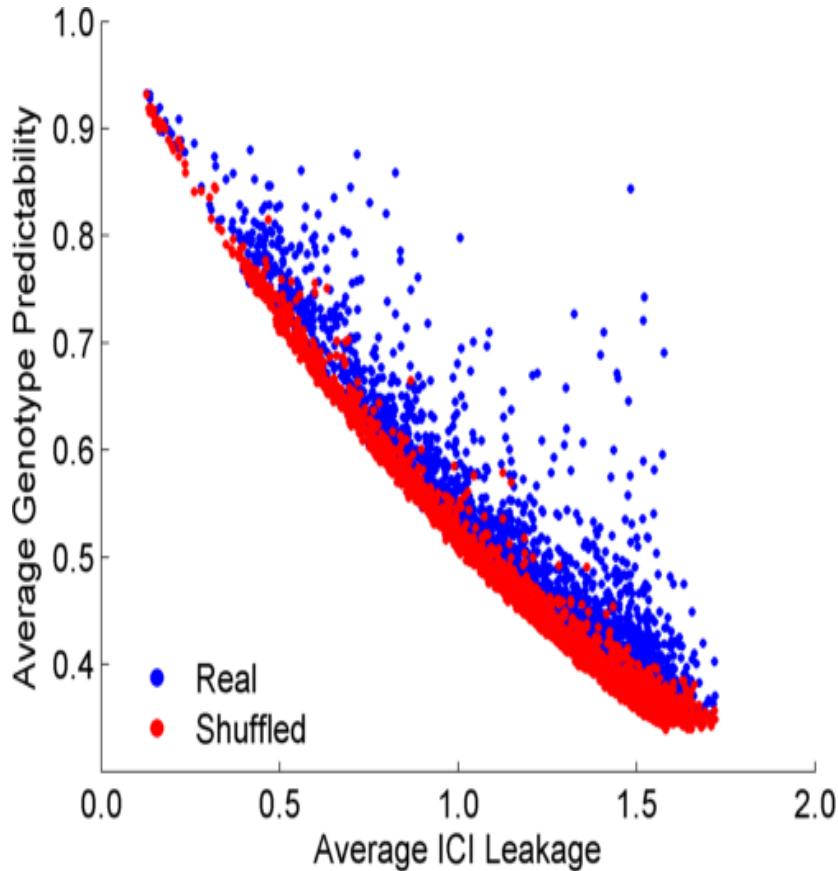
Compute entropy

$$\pi(V|E = e) = e^{-H(V|E = e)}$$

# Representative Expression, Genotype, eQTL Datasets

- mRNA sequencing for 462 individuals
  - Publicly availableQuantification for protein coding genes
- Approximately 3,000 cis-eQTL (FDR<0.05)
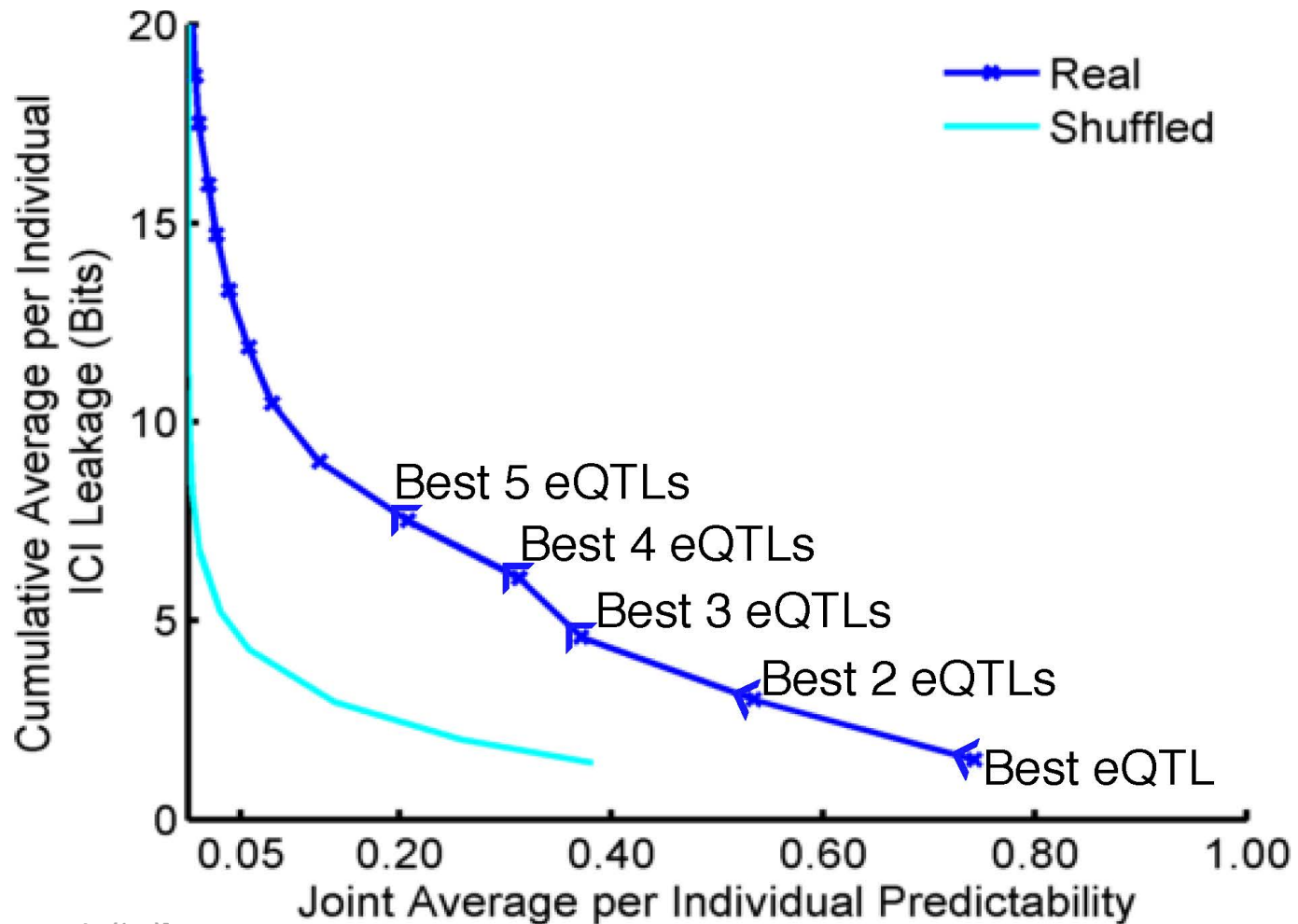- Genotypes are available from the 1000 Genomes Project

# Per eQTL and ICI Cumulative Leakage versus Genotype Predictability
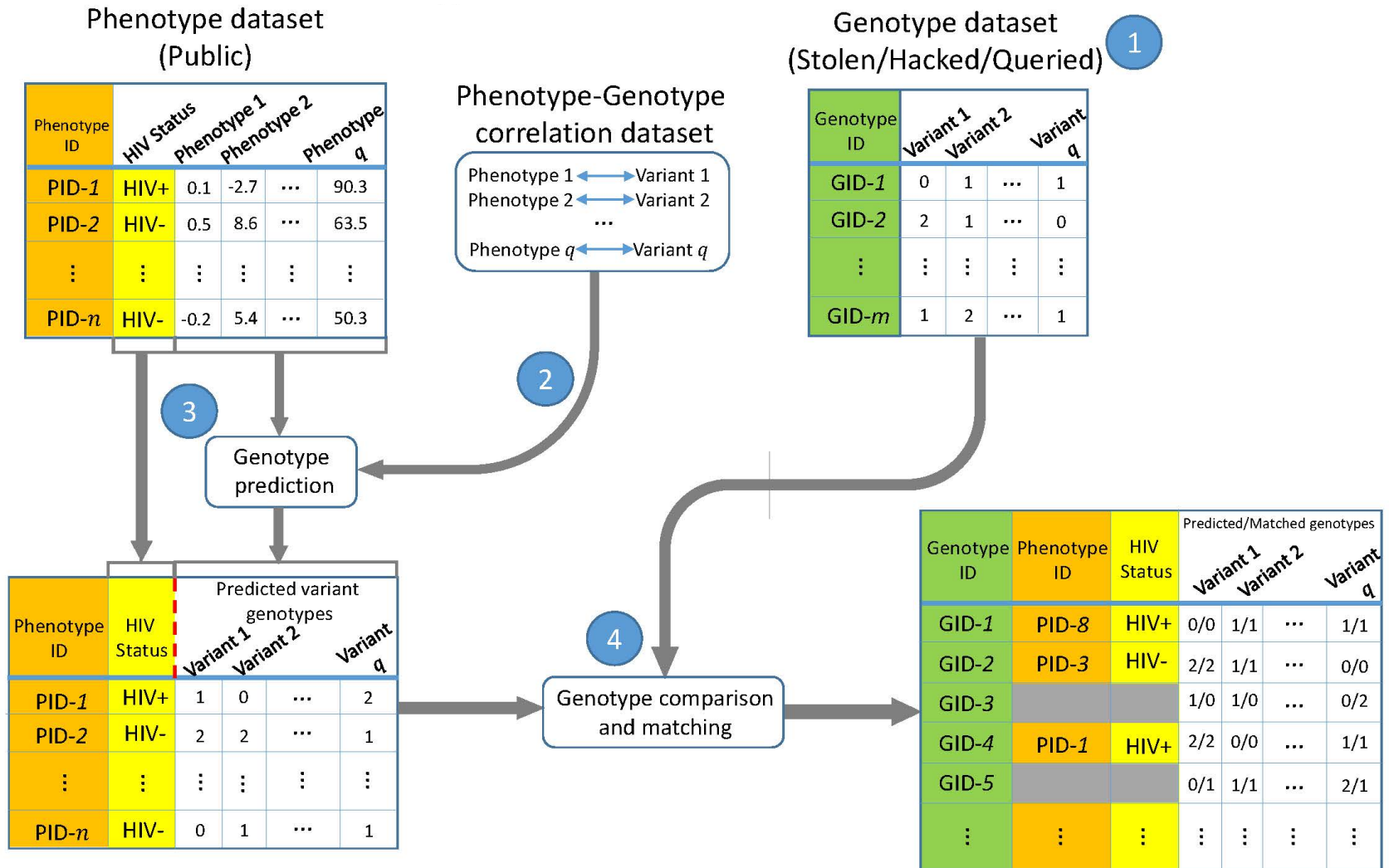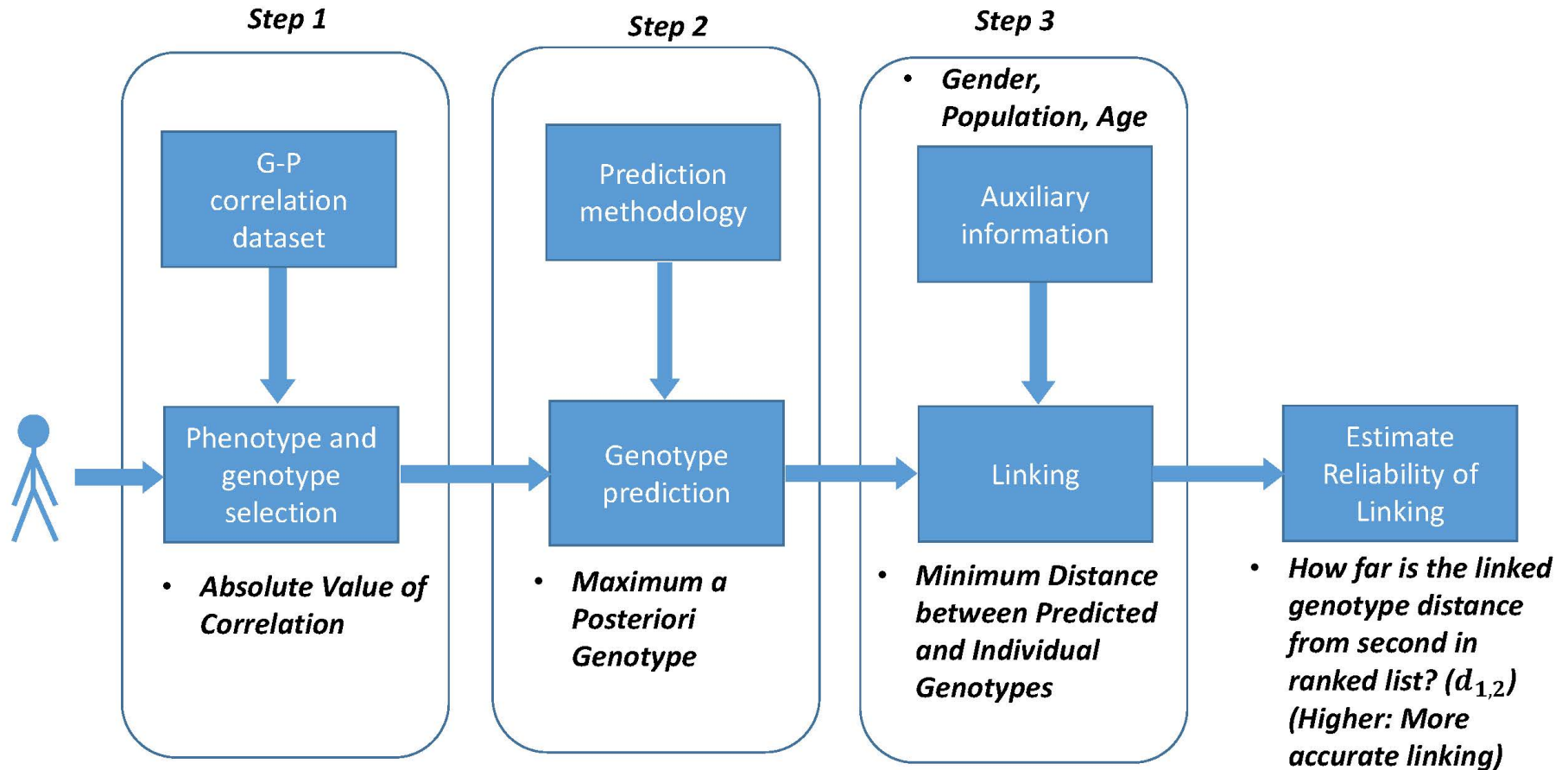
Colors by absolute correlation
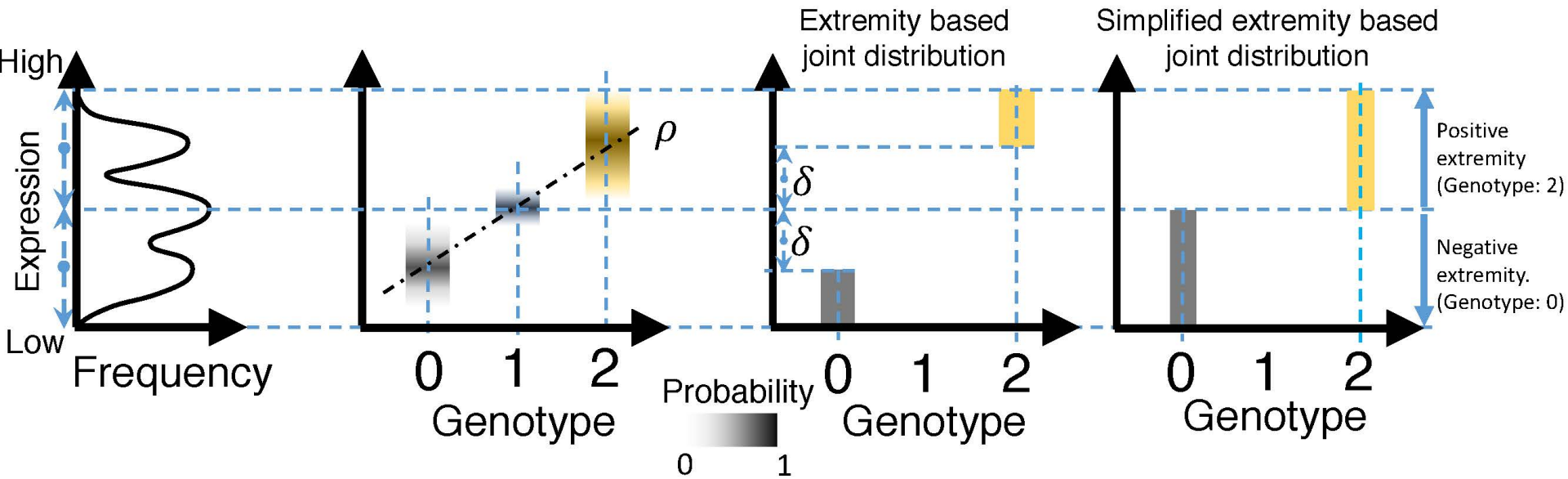
# Cumulative Leakage versus Joint Predictability

# Linking Attack Scenario

# Steps in Instantiation of a (Mock) Linking Attack
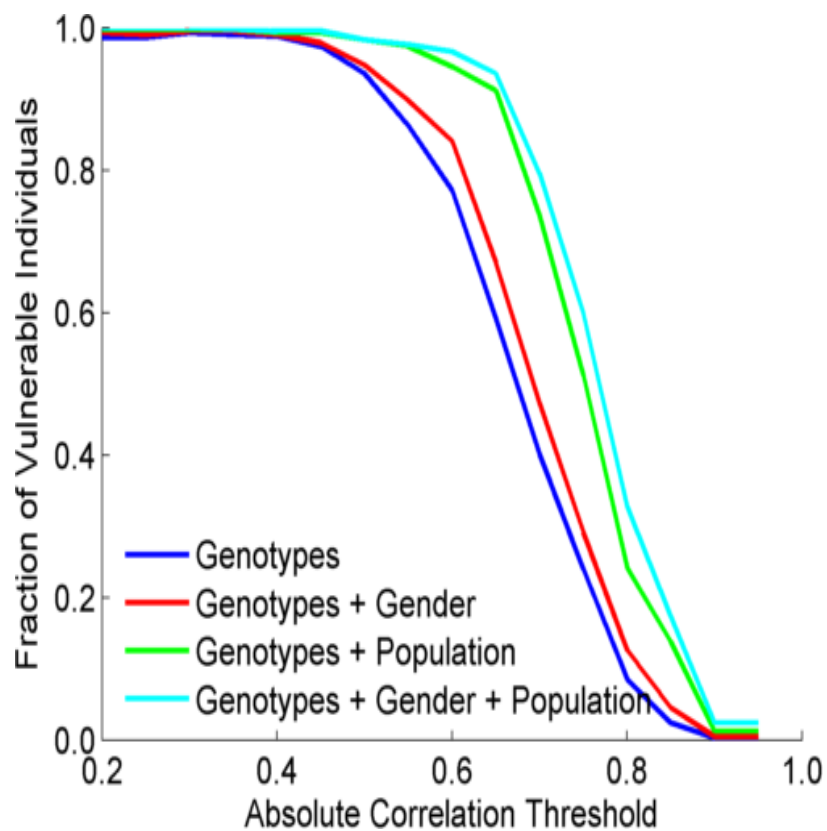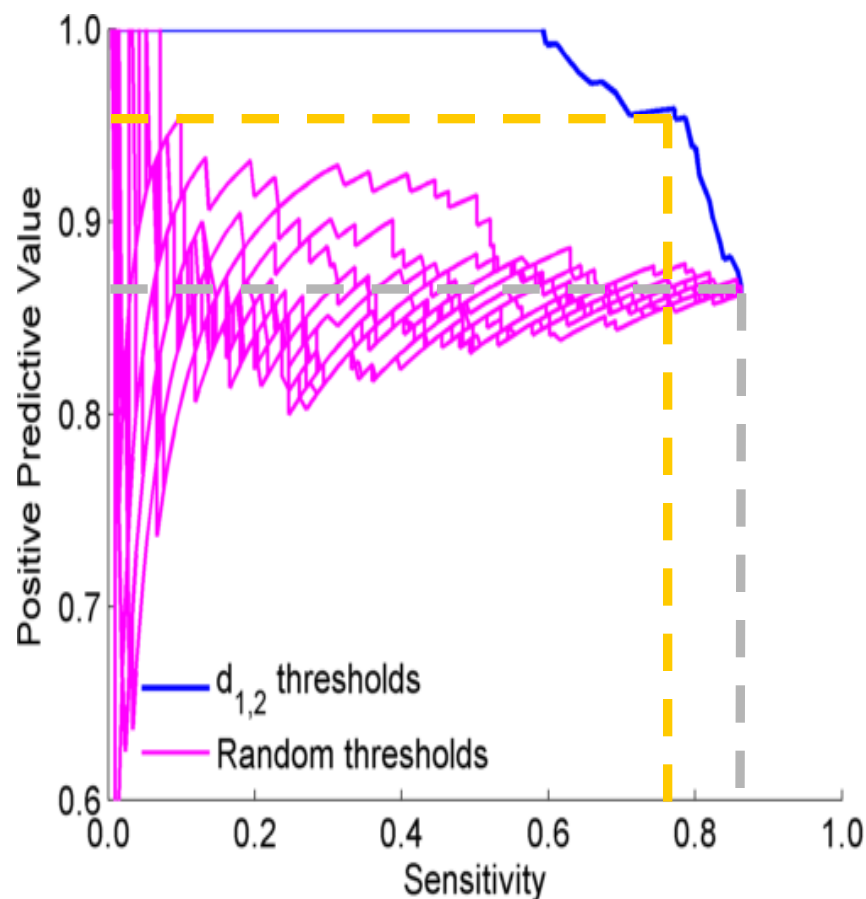


Step 1

Step 2

Step 3

G-P correlation dataset

Prediction methodology

- *Gender, Population, Age*

Auxiliary information

Phenotype and genotype selection

Genotype prediction

Linking

Estimate Reliability of Linking

- *Absolute Value of Correlation*

- *Maximum a Posteriori Genotype*

- *Minimum Distance between Predicted and Individual Genotypes*

- *How far is the linked genotype distance from second in ranked list? ($d_{1,2}$) (Higher: More accurate linking)*

[Harmanci et al. Nat. Meth. ('16)]

**Extremity based joint distribution**

**Simplified extremity based joint distribution**

High

Low

Expression

Frequency

$\rho$

Genotype

0 1 2

Probability

0 — 1

$\delta$

$\delta$

Genotype

0 1 2

Positive extremity (Genotype: 2)

Negative extremity. (Genotype: 0)

Genotype

0 1 2

Levels of Expression-Genotype Model Simplifications:

$E_k$  $p(E_k, V_k)$

$E_k$  $p(E_k, V_k)$

$E_k$  $p(E_k, V_k)$

$E_k$  $p(E_k, V_k)$

$V_k$

0 1 2

$\mu$  $\sigma_1$  $\sigma_2$  $\sigma_3$

$V_k$

0 1 2

$\mu$  $\sigma$

$V_k$

0 1 2

$e_{mid}$

$V_k$

0 1 2

[Harmanci et al. Nat. Meth. ('16)]

# Extremity based linking with homozygous genotypes

# Attacker can estimate the reliability of linkings



Sensitivity: Fraction of correctly linked Individuals among all individuals

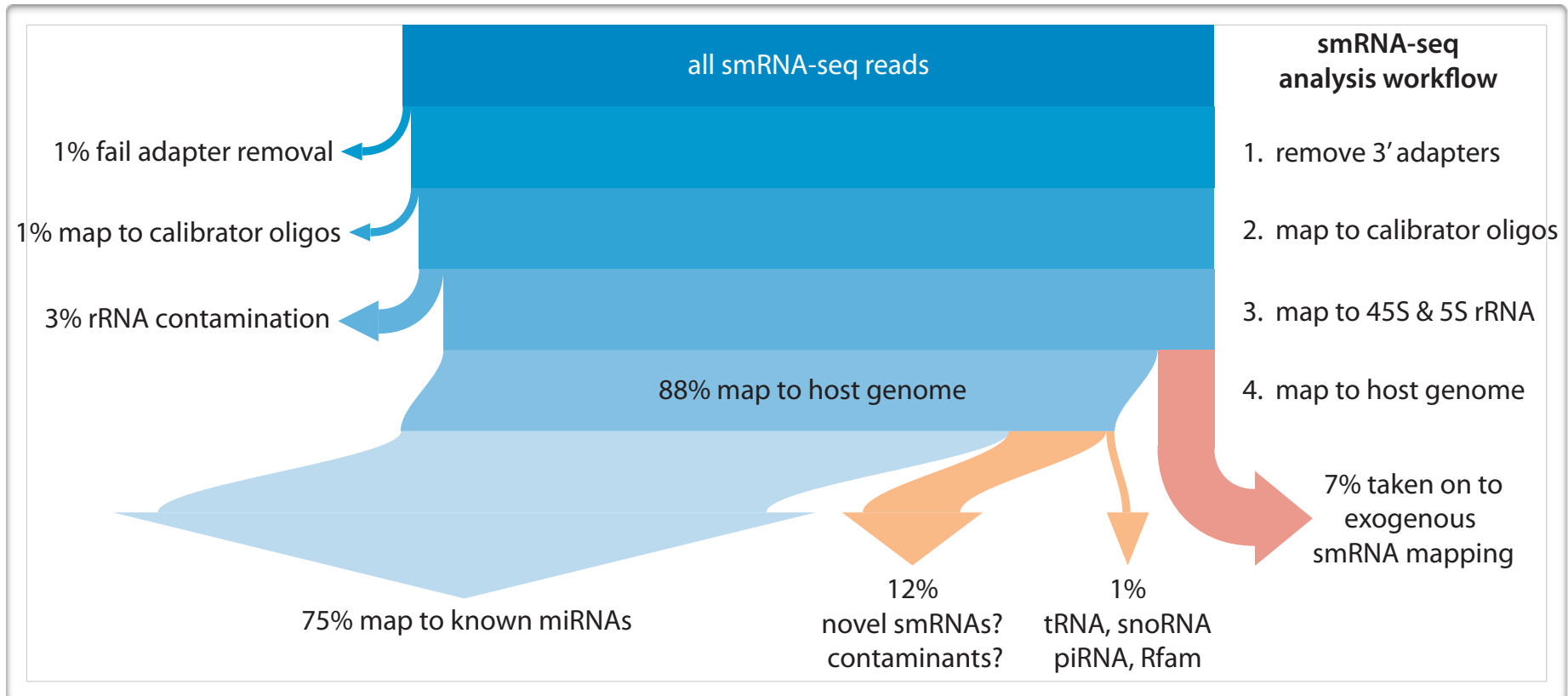PPV: Fraction of correctly linked individuals among selected individuals

- Large-scale data from consortia #1
  - exRNA.org

- Long-RNA pipeline
  - RSeqTools & anonymized MRF format

- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous

- Large-scale data from consortia #2
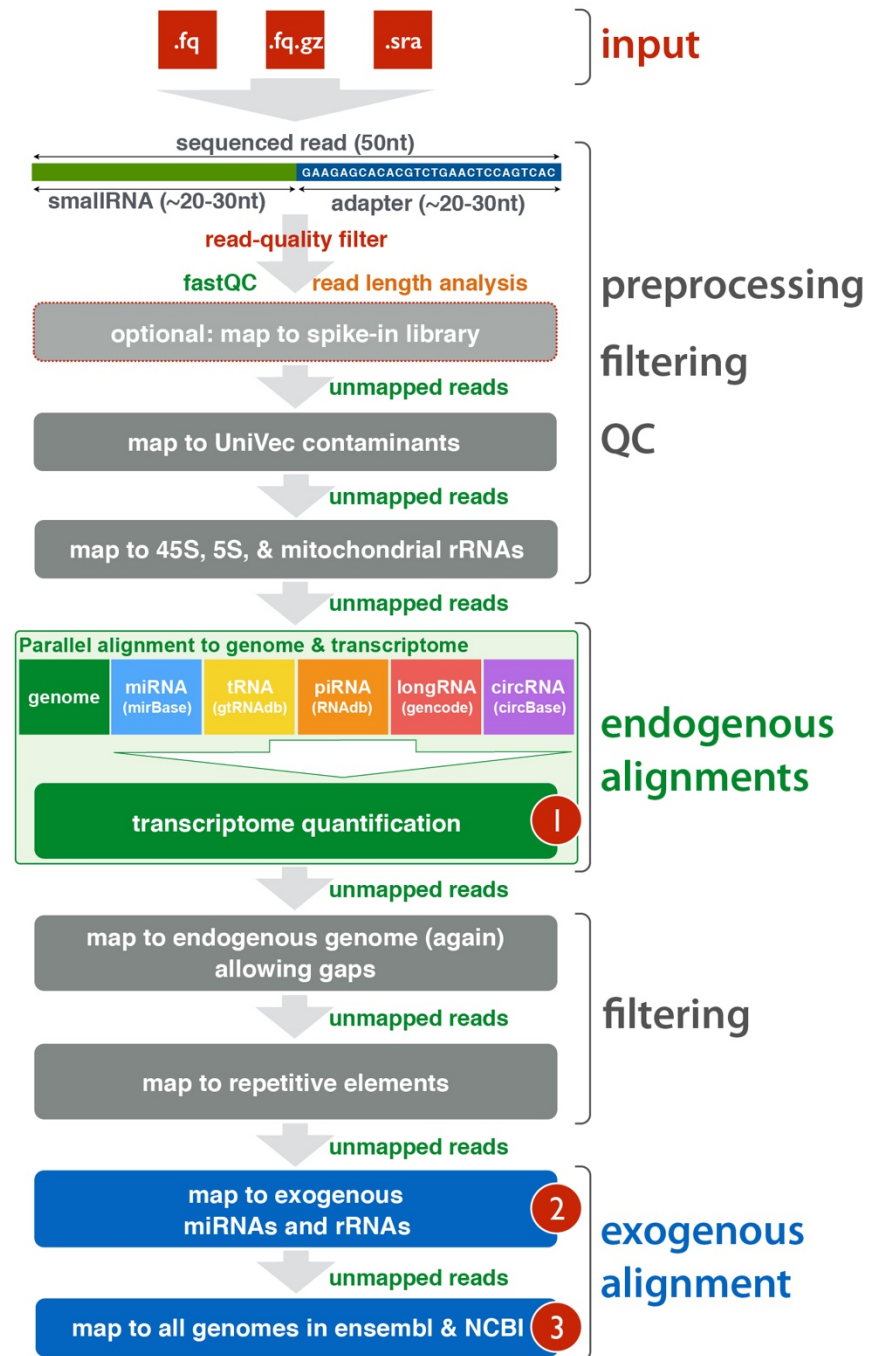  - Co-authorship networks show data flow through key broker individuals

- Large-scale data from consortia #1
  - exRNA.org

- Long-RNA pipeline
  - RSeqTools & anonymized MRF format

- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous

- Large-scale data from consortia #2
  - Co-authorship networks show data flow through key broker individuals

# for a typical cellular sample...



smRNA-seq analysis workflow

all smRNA-seq reads

1% fail adapter removal
1% map to calibrator oligos
3% rRNA contamination
88% map to host genome

1. remove 3' adapters
2. map to calibrator oligos
3. map to 45S & 5S rRNA
4. map to host genome

7% taken on to exogenous smRNA mapping

75% map to known miRNAs

12% novel smRNAs? contaminants?

1% tRNA, snoRNA piRNA, Rfam

- exRNA samples typically much noisier

- cascade of read-alignment steps mitigates contamination

**ex**tra-**ce**llular **R**NA
**p**rocessing **t**oolkit

- automatic pre-processing and QC of sequence reads

- explicit filtering of contaminants & rRNA

- quantification of spike-in sequences and many different smallRNA biotypes

- support for random barcodes (Bioo)

- choice of 3 end-points:

  (1) endogenous only

  (2) exogenous miRNA +rRNA

  (3) exogenous genomes

# total reads by biotype

- large contribution from miRNA and mRNA

- also some signal from exogenous sequences

# exceRpt @ Genboree.org



- extremely simple to use (1 input, 1 output)

- can process multiple samples in parallel

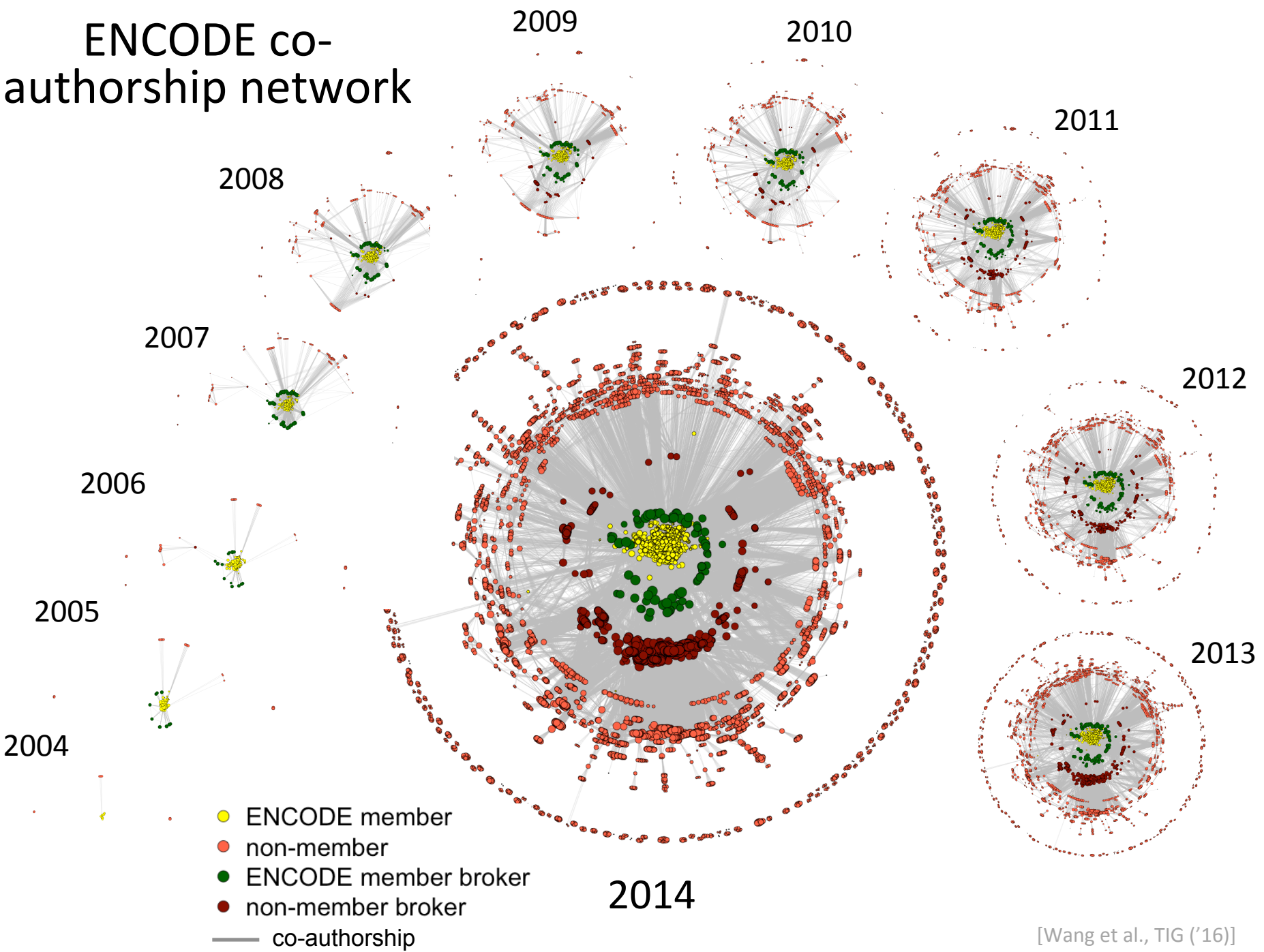- very customisable (choice of smRNA libs, calibrators, etc)

- Large-scale data from consortia #1
  - exRNA.org

- Long-RNA pipeline
  - RSeqTools & anonymized MRF format

- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous

- Large-scale data from consortia #2
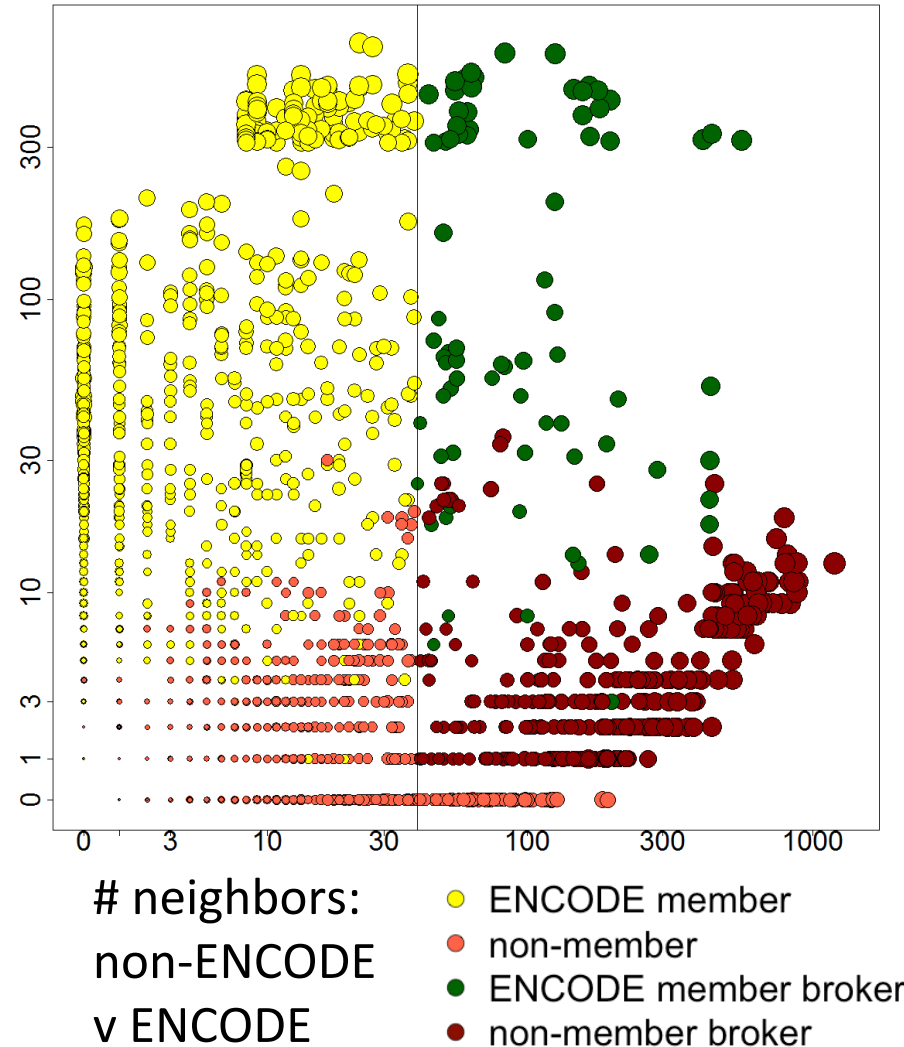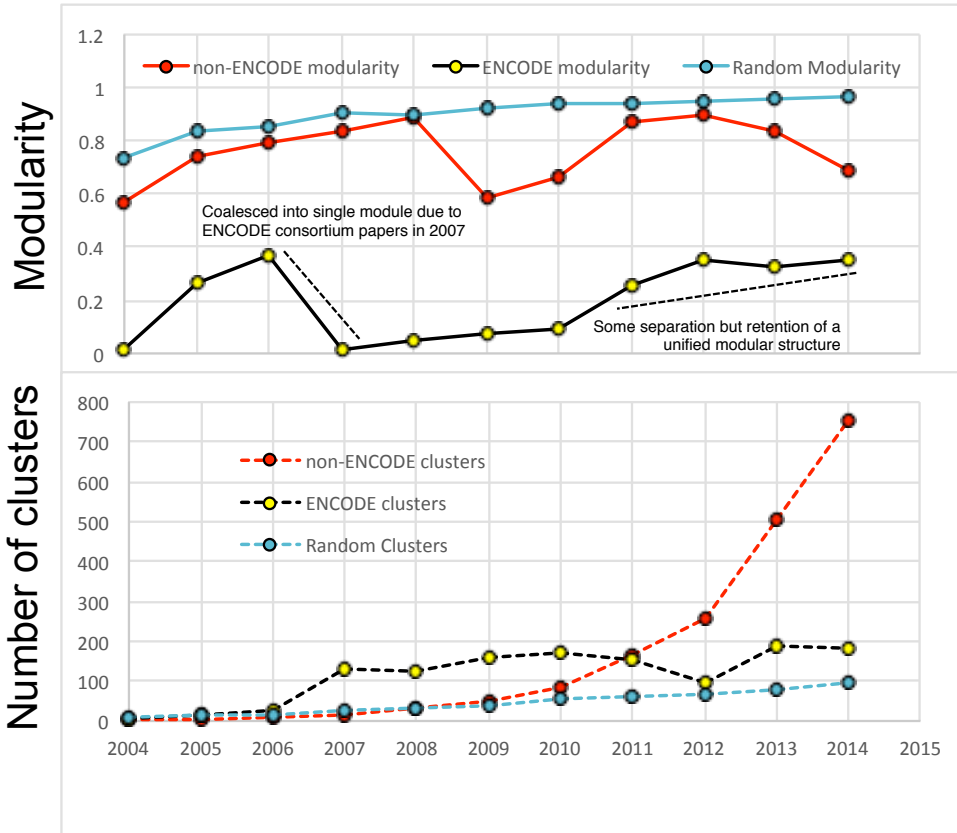  - Co-authorship networks show data flow through key broker individuals

**Papers authored by ENCODE consortium members vs. those that use ENCODE data but were not funded by ENCODE**

# ENCODE co-authorship network



2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014

Legend:
- ● ENCODE member (yellow)
- ● non-member (salmon)
- ● ENCODE member broker (green)
- ● non-member broker (dark red)
- ── co-authorship

[Wang et al., TIG ('16)]

# Network statistics highlight
# change in modularity with consortium rollouts (L)
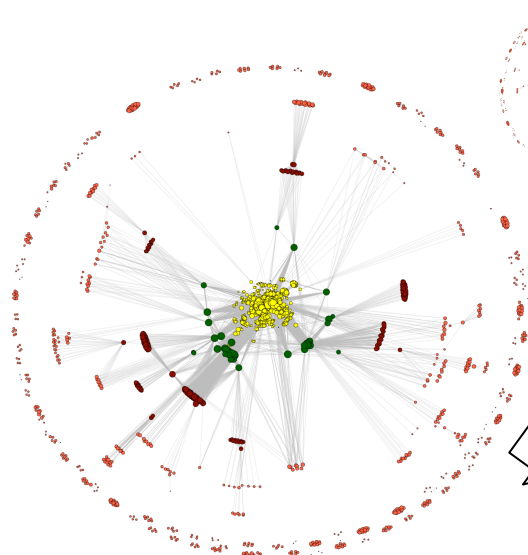# & importance of broker role (R)



# neighbors:
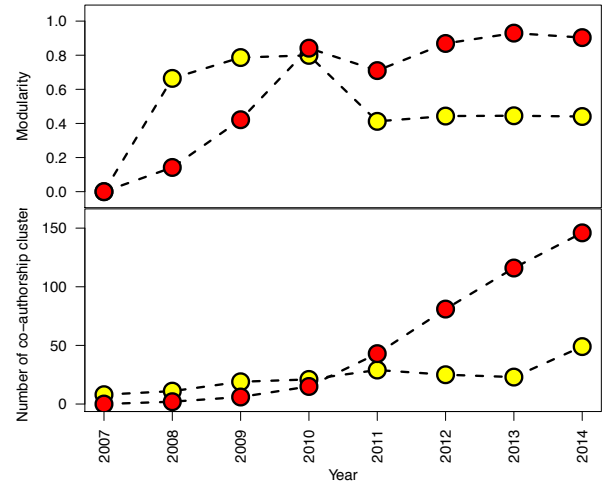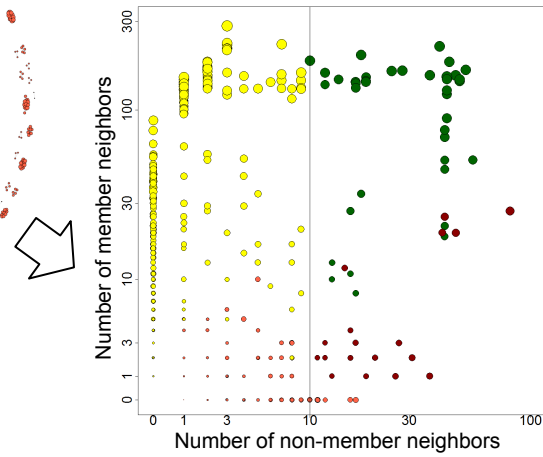non-ENCODE
v ENCODE

- ENCODE member
- non-member
- ENCODE member broker
- non-member broker

2014    2013    2012    2011    2010    2009    2008    2007

**modENCODE**

Number of member neighbors

Number of non-member neighbors

Modularity

Number of co-authorship cluster

Year

consortium
member

non–member

member
broker

non–member
broker

consortium
network

non–consortium
network

random
network

co–authorship

- Large-scale data from consortia #1
  - exRNA.org
- Long-RNA pipeline
  - RSeqTools & anonymized MRF format
- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous
- Large-scale data from consortia #2
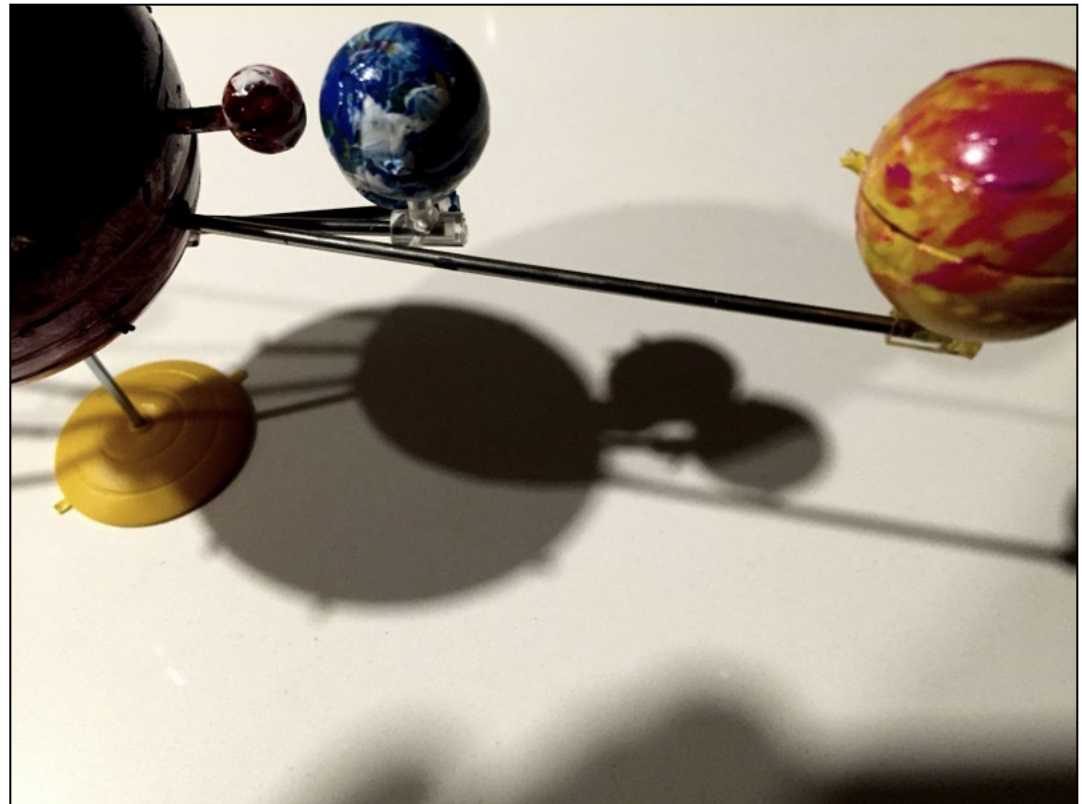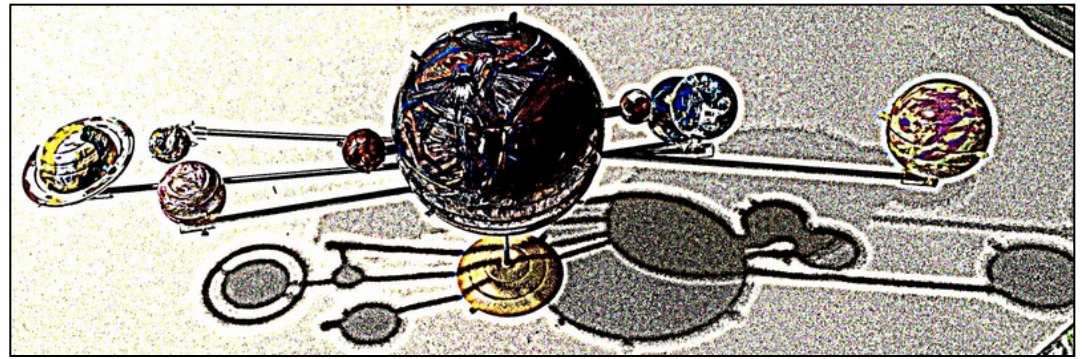  - Co-authorship networks show data flow through key broker individuals

- Large-scale data from consortia #1
  - exRNA.org
- Long-RNA pipeline
  - RSeqTools & anonymized MRF format
- Privacy risk in eQTLs
  - Quantifying & removing variant info from expression levels + eQTLs
  - Linking Attack using extreme expression levels

- Short-RNA pipeline
  - exceRpt
  - Appreciable mapping beyond miRNAs & mRNAs, including exogenous
- Large-scale data from consortia #2
  - Co-authorship networks show data flow through key broker individuals

**PrivaSeq**.gersteinlab.org
**A Harmanci**

**RSEQtools**.gersteinlab.org
**L Habegger**, A Sboner,
TA Gianoulis, J Rozowsky,
A Agarwal, M Snyder

"Encode authors"
**D Wang**, KK Yan,
J Rozowsky, E Pan

**exRNA.org** & **exceRpt**
**R Kitchen**
**J Rozowsky**
A Milosavljevic
M Roth
S Subramanian



# Acknowledgments

# Extra

# Info about content in this slide pack

- General PERMISSIONS
  - This Presentation is copyright Mark Gerstein, Yale University, 2015.
  - Please read permissions statement at www.**gersteinlab.org/misc/permissions.html** .
  - Feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).
  - Paper references in the talk were mostly from Papers.GersteinLab.org.

- PHOTOS & IMAGES. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see http://streams.gerstein.info .
  - In particular, many of the images have particular EXIF tags, such as kwpotppt , that can be easily queried from flickr, viz: http://www.flickr.com/photos/mbgmbg/tags/kwpotppt