# Prediction of Gene Expression in Yeast using Conserved Sequence Templates
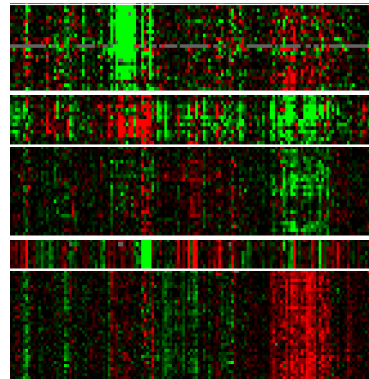
Derek Chiang & Alan Moses

Molecular & Cell Biology, UCB
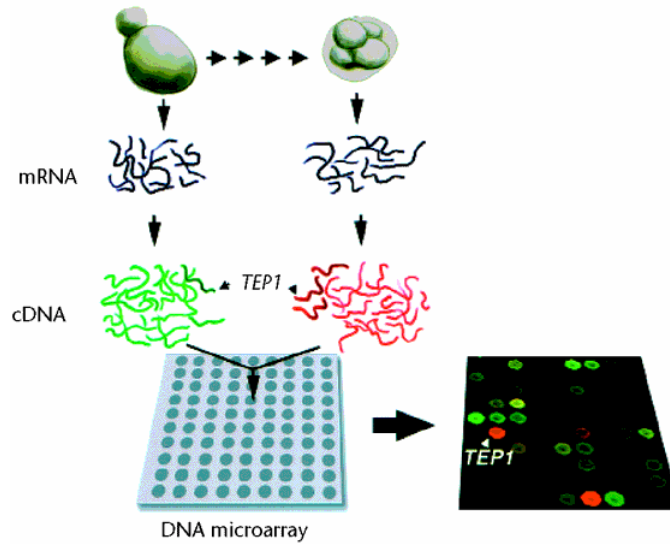
Eisen Lab

---

# PROBLEM

What information controlling gene expression
is encoded in genome sequences?

```
>Saccharomyces cerevisiae chr V
CGTCTCCTCCAAGCCCTGTTGTCTCTTACCC
GGATGTTCAACCAAAAGCTACTTACTACCTT
TATTTTATGTTTACTTTTTATAGATTGTCTT
TTTATCCTACTCTTTCCCACTTGTCTCTCGC
TACTGCCGTGCAACAAACACTAAATCAAAAC
AGTGAAATACTACTACATCAAAACGCATATT
CCCTAGAAAAAAAAATTTCTTACAATATACT
ATACTACACAATACATAATCACTGACTTTCG
TAACAACAATTTCCTTCACTCTCCAACTTCT
CTGCTCGAATCTCTACATAGTAATATTATAT
CAAATCTACCGTCTGGAACATCATCGCTATC
CAGCTCTTTGTGAACCGCTACCATCAGCATG
TACAGTGGTACCTTCGTGTTATCTGCAGCGA
GAACTTCAACGTTTGCCAAATCAAGCCAATG
TGGTAACAACCACACCTCCGAAATCTGCTCC
AAAAGATACTCCAGTTTCTGCCGAAATGTTT
```
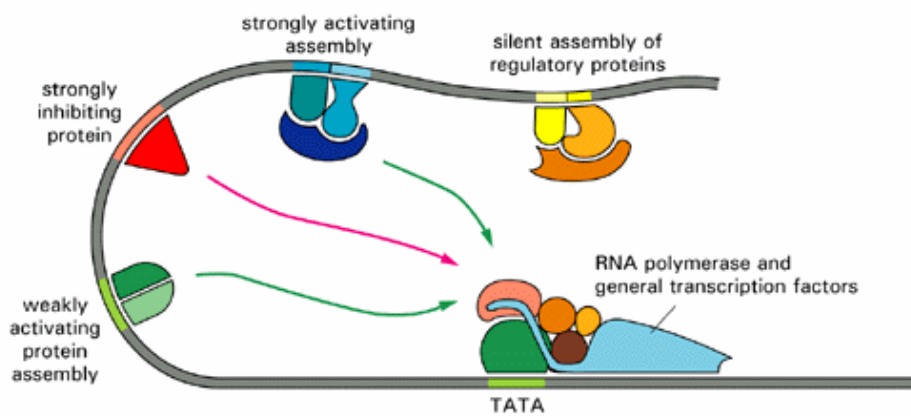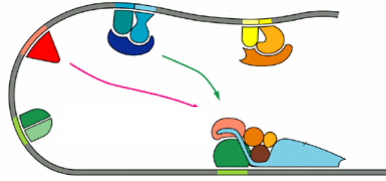
Features?

# Gene Expression: Experiment
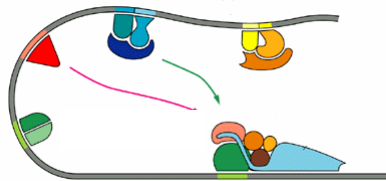


# Gene Expression: Mechanism

## Promoter Structure

# Conserved Sequence Rules

Hidden Variables

Sites $\qquad$ $S_1$, $S_2$, …
Positions $\qquad$ $P_{11}$, $P_{12}$, …
$\qquad\qquad\qquad$ $P_{21}$, $P_{22}$, …

# Conserved Sequence Rules

Hidden Variables

Sites $\qquad$ $S_1$, $S_2$, …
Positions $\qquad$ $P_{11}$, $P_{12}$, …
$\qquad\qquad\qquad$ $P_{21}$, $P_{22}$, …

PRIORS from sequence conservation

# Conserved Sequence Rules



Hidden Variables

Sites $S_1, S_2, \ldots$
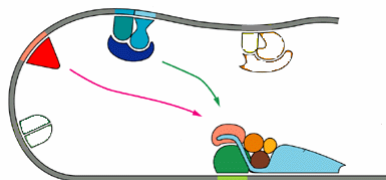Positions $P_{11}, P_{12}, \ldots$
$P_{21}, P_{22}, \ldots$

$S_2$

$S_1$

**Conserved Sequence Rules**

$\{\ S_1 > 0;\ \ S_2 > 0;$
$\min_{j}(\ |\ P_{1j} - P_{2j}\ | < 30\ )\ \}$

---

# Conserved Word Pairs

➢ **Finding Conserved Words**

➢ Evaluate Word Pairs

  1) Joint Conservation
  2) Close Spacing

➢ Validate with Gene Expression

## Finding Conserved Words

> 3860 CLUSTALW alignments from MIT

> Conserved word

$C_w$ ≡ Found in same position in 3+ genomes within 600 bp of gene start

```
> MET28_YAP5
Scer    AACCTAAAACCAAAAAAAA-A-AAATAAGTCACGTGCACT
Spar    AATAAAAAATAGACTAACA-A-ATTGCGGTCACGTGCACT
Smik    AATCCCAGGCCAAAAACCAGA-AATTGAGTCACGTGCAGT
Sbay    GTCACGTGCCCCGACGGCCCCACAACTGTGGCATCCATCTT
```
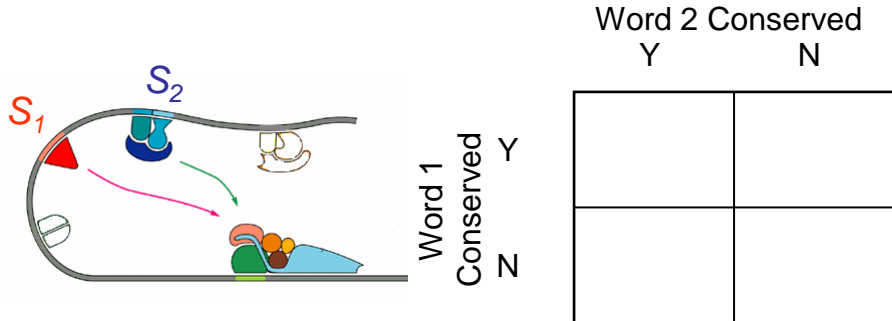
## Conserved Word Pairs

> Finding Conserved Words

> **Evaluate Word Pairs**

1) Joint Conservation
2) Close Spacing

> Validate with Gene Expression

# TEST: Joint Word Conservation



|  | Word 2 Conserved Y | N |
|---|---|---|
| Word 1 Conserved Y | | |
| N | | |

# TEST: Joint Word Conservation



|  | Word 2 Conserved Y | N |
|---|---|---|
| Word 1 Conserved Y | 32 | 162 |
| N | 134 | 3226 |

**Chi-square Test for Independence**
(Yates adjustment)

$$\chi^2 = \sum_k \frac{(|\,O_k - E_k\,| - \frac{1}{2})^2}{E_k}$$

## TEST: Joint Word Conservation

➢ 2090 words  (Length 6: Word-Rev complement )

➢ $2.06 \times 10^6$ Word PAIRS    (Exclude overlap)

**HISTOGRAM**

Empirical $\chi^2_{0.005} = 30$

Bonferroni $\chi^2_{0.05} = 31$



y-axis: $-\log_{10}($ Frequency $)$

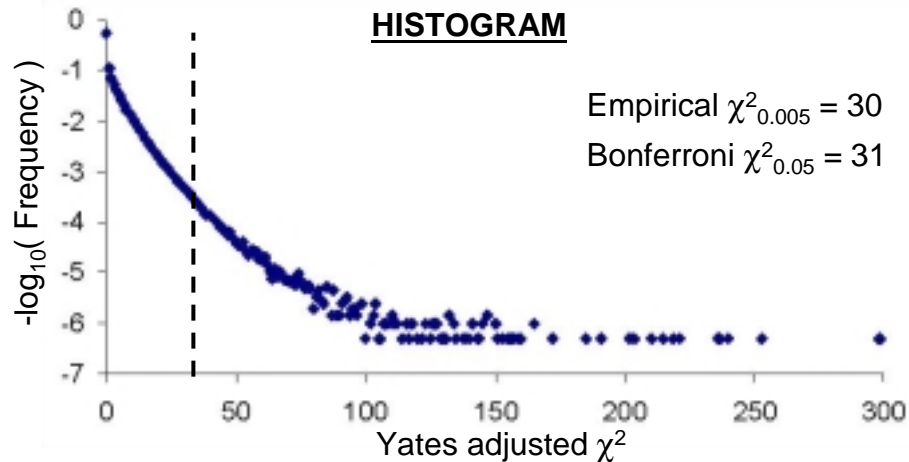x-axis: Yates adjusted $\chi^2$

## TEST: Close Spacing



$P_1 \longleftrightarrow P_2$

For conserved genes $g$ : 1 … $N$

$J_g = ( \overrightarrow{P_{1g}}, \overrightarrow{P_{2g}} )$ sampled jointly
from position distributions

Test Statistic

$$X = \frac{1}{N} \sum_g \min_{k \in g_i} |P_{1k} - P_{2k}|$$

NULL:      $X$ obtained from random sampling

ALT:      $X$ <u>smaller</u> than expected from sampling

# TEST: Close Spacing

## Nonparametric Bootstrap

Avg. Min. Distance

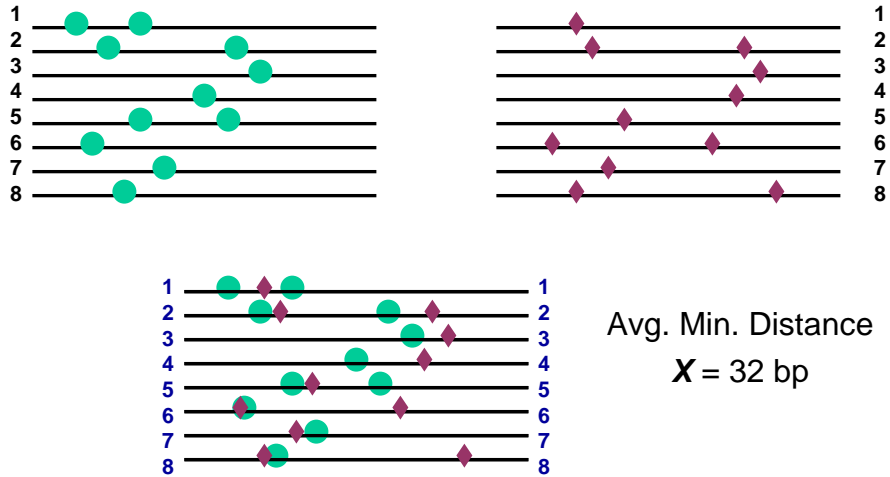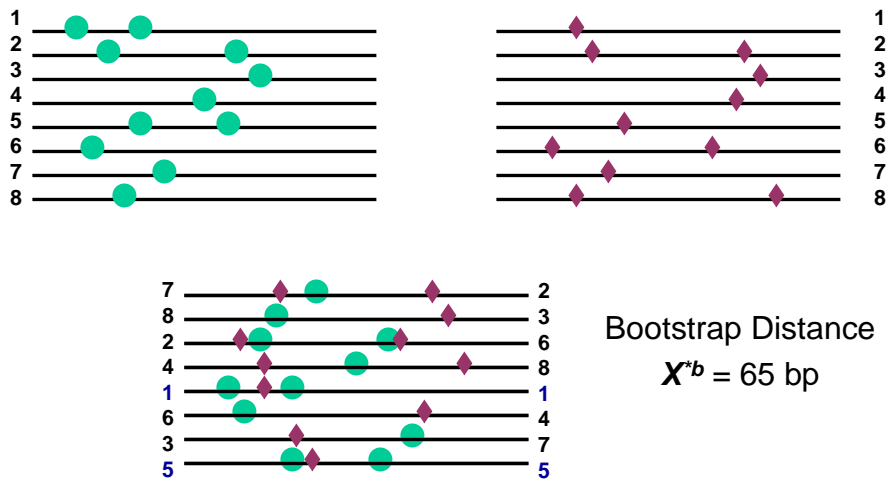$$X = 32 \text{ bp}$$

# TEST: Close Spacing

## Nonparametric Bootstrap

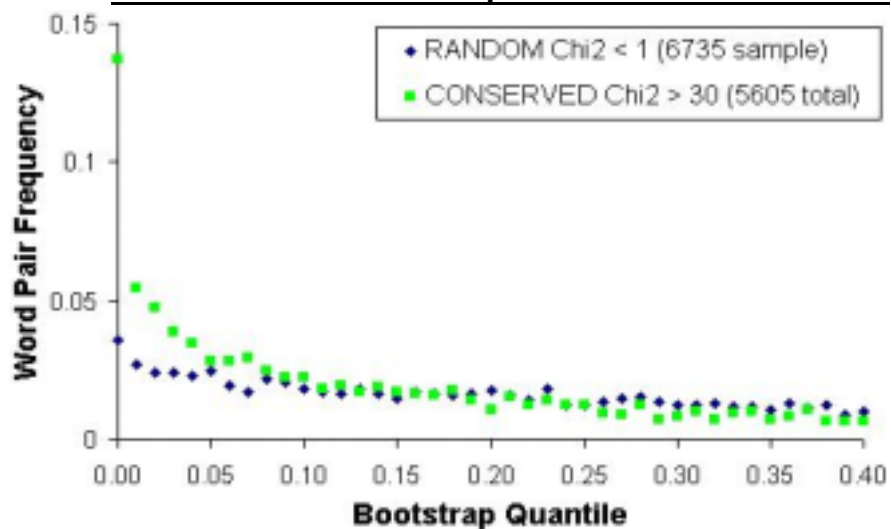Bootstrap Distance

$$X^{*b} = 65 \text{ bp}$$

# TEST: Close Spacing

## Nonparametric Bootstrap

- Position distributions  $\Pi_1 = \{\ \overrightarrow{P_{11}},\ \dots\ ,\ \overrightarrow{P_{1v}}\ \}$

  $\Pi_2 = \{\ \overrightarrow{P_{21}},\ \dots\ ,\ \overrightarrow{P_{2w}}\ \}$

- Data  $J_1 = (\ \overrightarrow{P_{11}},\ \overrightarrow{P_{21}}\ ),\ \dots\ ,\ J_n = (\ \overrightarrow{P_{1n}},\ \overrightarrow{P_{2n}}\ )$

- Bootstrap $J_1^{*b} = (\ \overrightarrow{P_{11}}^{*b},\ \overrightarrow{P_{21}}^{*b}\ ),\ \dots\ ,\ J_n = (\ \overrightarrow{P_{1n}}^{*b},\ \overrightarrow{P_{2n}}^{*b}\ )$

  resampled with replacement from $\Pi_1$, $\Pi_2$

- Record quantile of $X$ in 100000 samples of $X^{*b}$
  (empirical null distribution)

---

# TEST: Close Spacing

### HISTOGRAM of Bootstrap Quantiles for Word Pairs



Legend:
- ◆ RANDOM Chi2 < 1 (6735 sample)
- ■ CONSERVED Chi2 > 30 (5605 total)

X-axis: Bootstrap Quantile
Y-axis: Word Pair Frequency

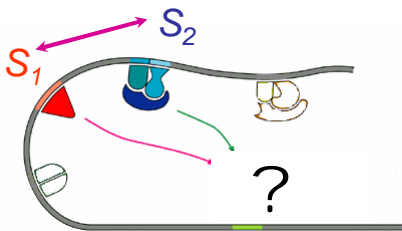## Conserved Word Pairs

➢ Finding Conserved Words

➢ Evaluate Word Pairs

   1) Joint Conservation
   2) Close Spacing

➢ **Validate with Gene Expression**

## Validating Expression Subsets

**Conserved Sequence (SUBSETTING) Rules**



Genome (6000 genes)

$$\{\ S_1 > 0;\quad S_2 > 0;$$
$$\min_{j}(\ |\ P_{1j} - P_{2j}\ | < d\ )\ \}$$

SUBSET ($N$ genes)

Assess gene expression

## Validating Expression Subsets

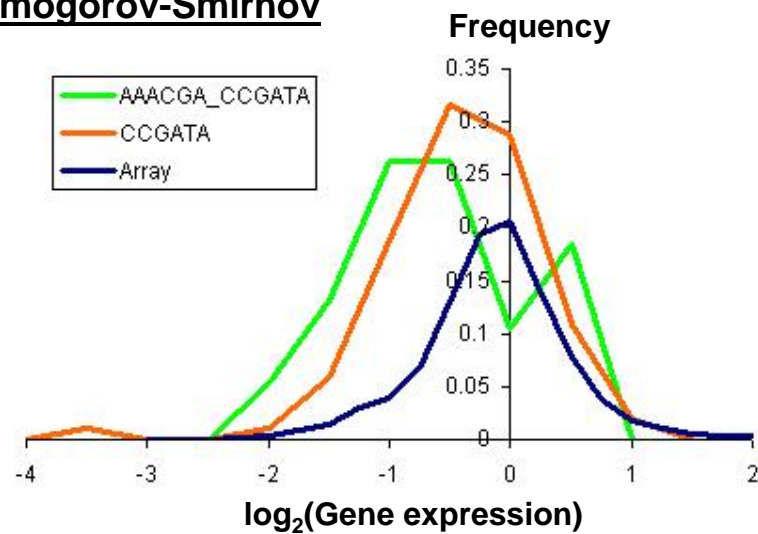**Some Nonparametric Tests**

1) Subset mean                  (Mann-Whitney)

2) Subset distribution        (Kolmogorov-Smirnov)

3) Subset weighted correlation   ( ? )

4) Kernel density classification   (Mixture of normals)

Find optimal word pair distance $d^*$ for each test …
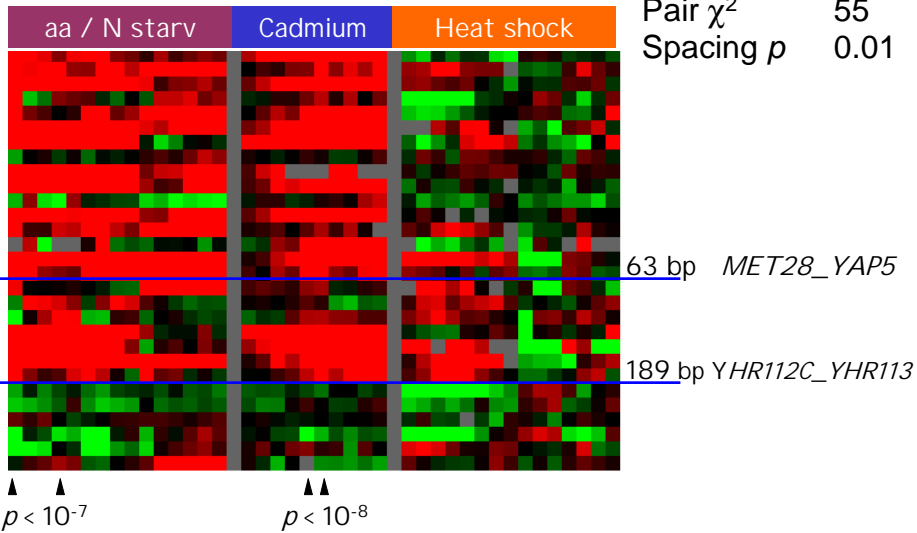
---

## Validating Expression Subsets

**Kolmogorov-Smirnov**
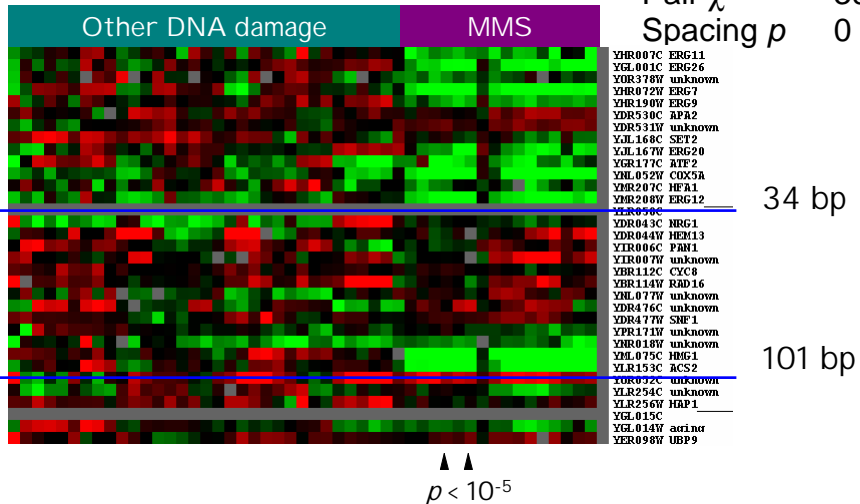
## Conserved Word Pairs – EXAMPLES

AACTGT-CACGTG (Cbf1-Met31)

| | |
|---|---|
| $N$ | 29 |
| Pair $\chi^2$ | 55 |
| Spacing $p$ | 0.01 |



aa / N starv  |  Cadmium  |  Heat shock

63 bp  *MET28_YAP5*

189 bp Y*HR112C_YHR113*

$p < 10^{-7}$      $p < 10^{-8}$

---

## Conserved Word Pairs – EXAMPLES

AAACGA-CCGATA  (Upc2-Hap1)

| | |
|---|---|
| $N$ | 32 |
| Pair $\chi^2$ | 39.5 |
| Spacing $p$ | 0 |



Other DNA damage  |  MMS

YHR007C ERG11
YGL001C ERG26
YOR378W unknown
YHR072W ERG7
YHR190W ERG9
YDR530C APA2
YDR531W unknown
YJL168C SET2
YJL167W ERG20
YGR177C ATF2
YNL052W COX5A
YMR207C HFA1
YMR208W ERG12
YLR050C

34 bp

YDR043C NRG1
YDR044W HEM13
YIR006C PAN1
YIR007W unknown
YBR112C CYC8
YBR114W RAD16
YNL077W unknown
YDR476C unknown
YDR477W SNF1
YPR171W unknown
YMR018W unknown
YML075C HMG1
YLR153C ACS2

101 bp

YOR052C unknown
YLR254C unknown
YLR256W HAP1
YGL015C
YGL014W asino
YER098W UBP9

$p < 10^{-5}$

## Pipeline Summary



| | Label | TEMPLATE | RANDOM |
|---|---|---|---|
| # Genes > 10 | | 5608 | 8768 |
| Spacing < 0.05 | | 1774 | 915 |
| K-S < $10^{-5}$ & Improve > 10 | | 334 | 29 |
| Known TF | | 110 | 6 |

TEMPLATE (Chi2 > 30)
RANDOM (1 of 10 Chi2 < 1)

# Word Pairs: 0, 2000, 4000, 6000, 8000, 10000

## Future Directions

➢ Better gene expression subset tests    (Timecourse)

➢ More flexible sequence models   (IUPAC, Self-dimer)

➢ Automate distance cutoff          (Distance *d*)

➢ Parameter optimization:  8 threshold values!
  ( Conserved: # aligned genomes & # upstream bp,
    Joint conservation $\chi^2$, Bootstrap quantile, K-S probs,
    Min gene #, Distance *d* )

## Acknowledgements



Michael Eisen

Eisen Lab
Justin Fay
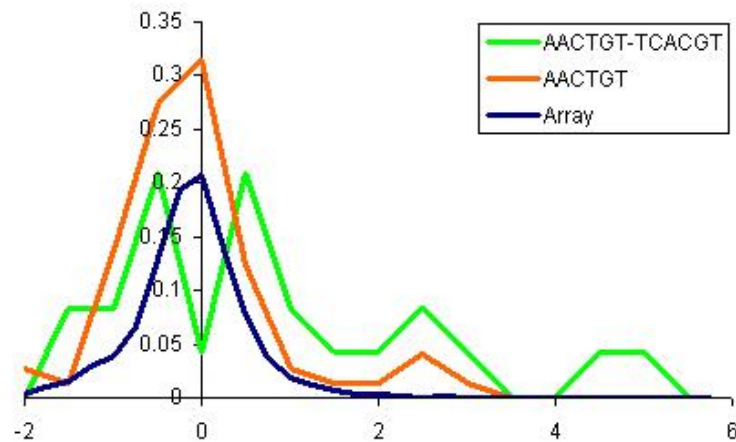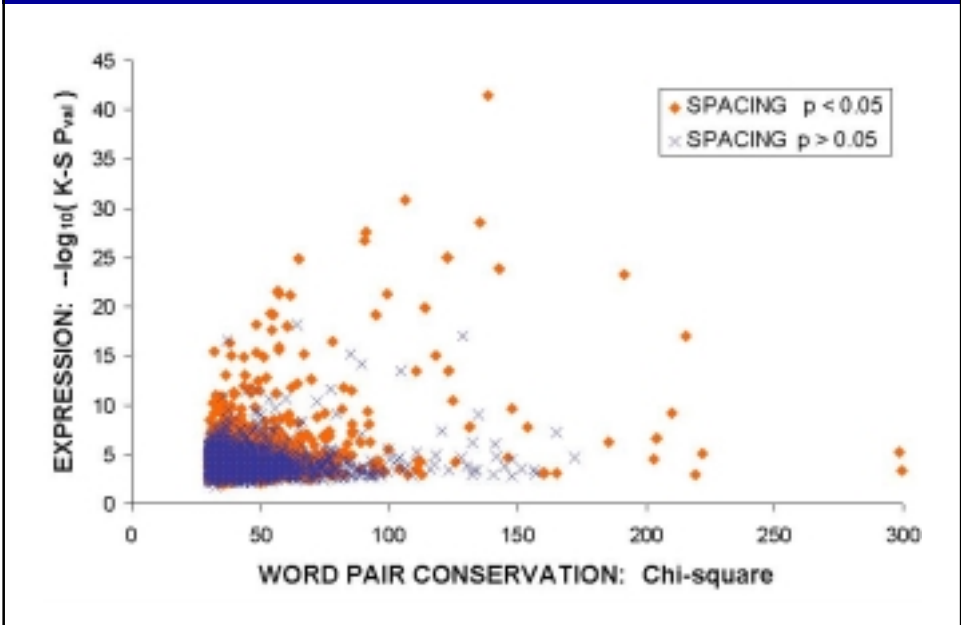Audrey Gasch
Hunter Fraser
Venky Nandagopal
Dan Pollard
Ben Lewis

## Validating Expression Subsets

**Kolmogorov-Smirnov**

Comparing Expression Samples



Comparing Expression Samples