# SUPPLEMENT A-1

Determination of CI and CI thresholds

**Define sORF in intergenic regions**

ATG.....................................TAG
90bp

F$_A$ — 75 bp
F$_B$ — 75 bp
F$_C$ — 75bp
F$_D$ — 75 bp
F$_E$ — 75bp
F$_F$ — 75bp

**Calculate Coding Index (CI):**

$$CI = \frac{P(coding|F_A) + P(coding|F_B) \dots + P(coding|F_F)}{7}$$
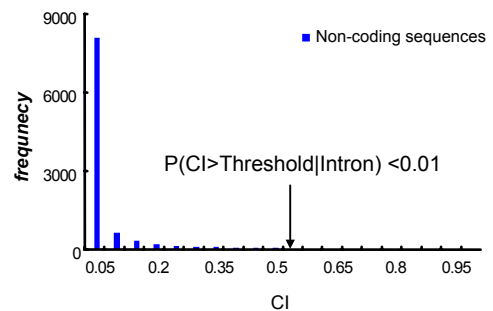
**Define coding sequences**

If (CI value of an sORF > the threshold),
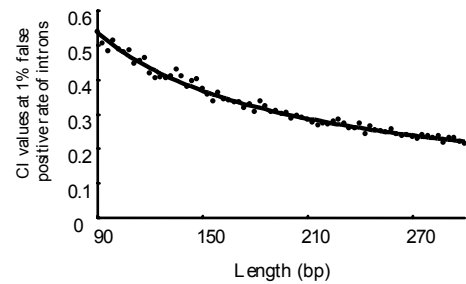the the sORF is a coding sequence

**Generate random sequences for coding and non-coding sequences with different length**

Estimate Coding Index (CI)

**Determine CI threshold based on intron random seq at a particular length**

■ Non-coding sequences

P(CI>Threshold|Intron) <0.01

*frequnecy*

CI

**Define CI thresholds**

CI values at 1% false positive rate of introns

Length (bp)

## SUPPLEMENT A-2

Determination of CI thresholds based on simulation studies of CI values of NCDS-like random sequences



A. Arabidopsis data.

$y = 17.644x^{-0.7212}$

$R^2 = 0.989$
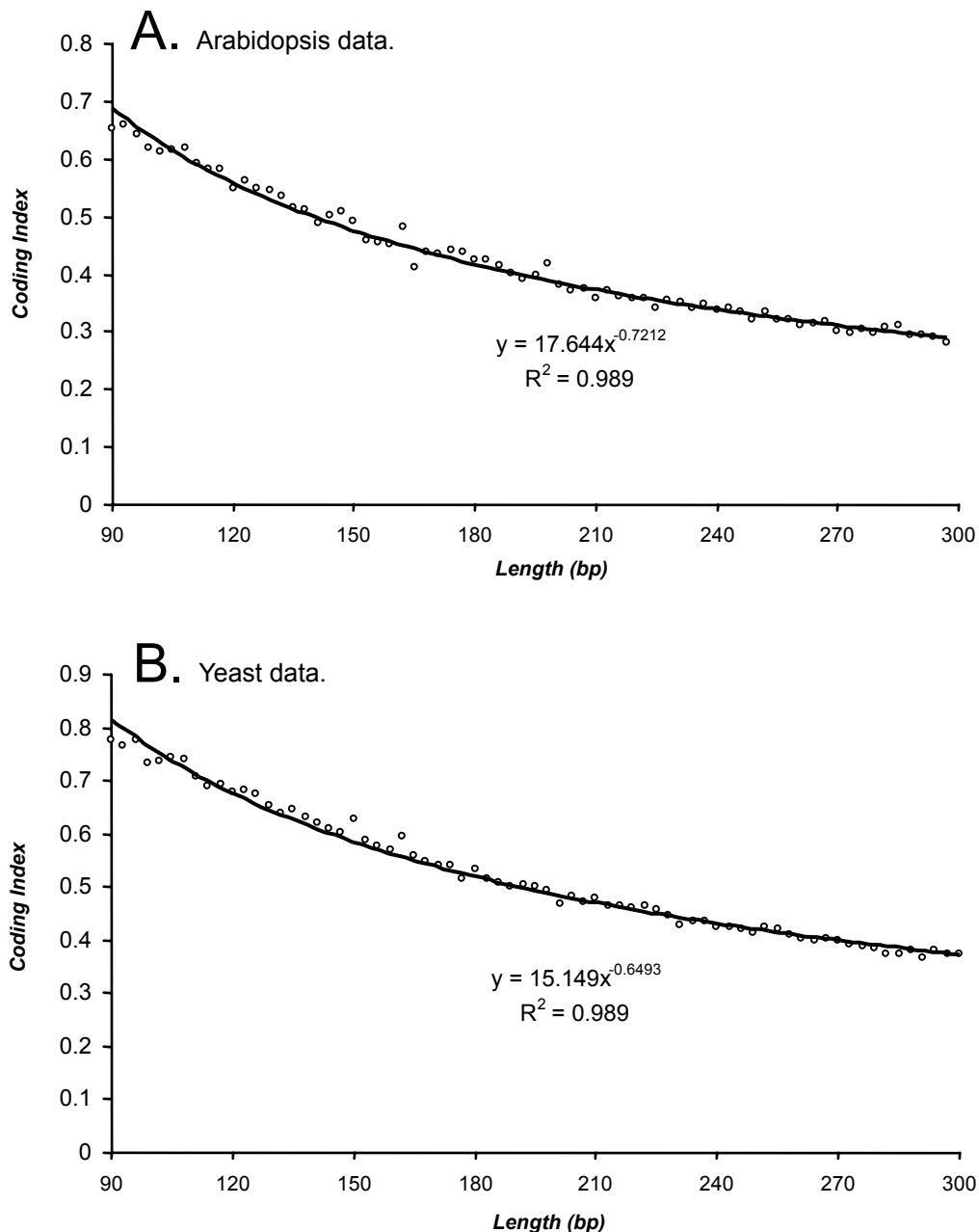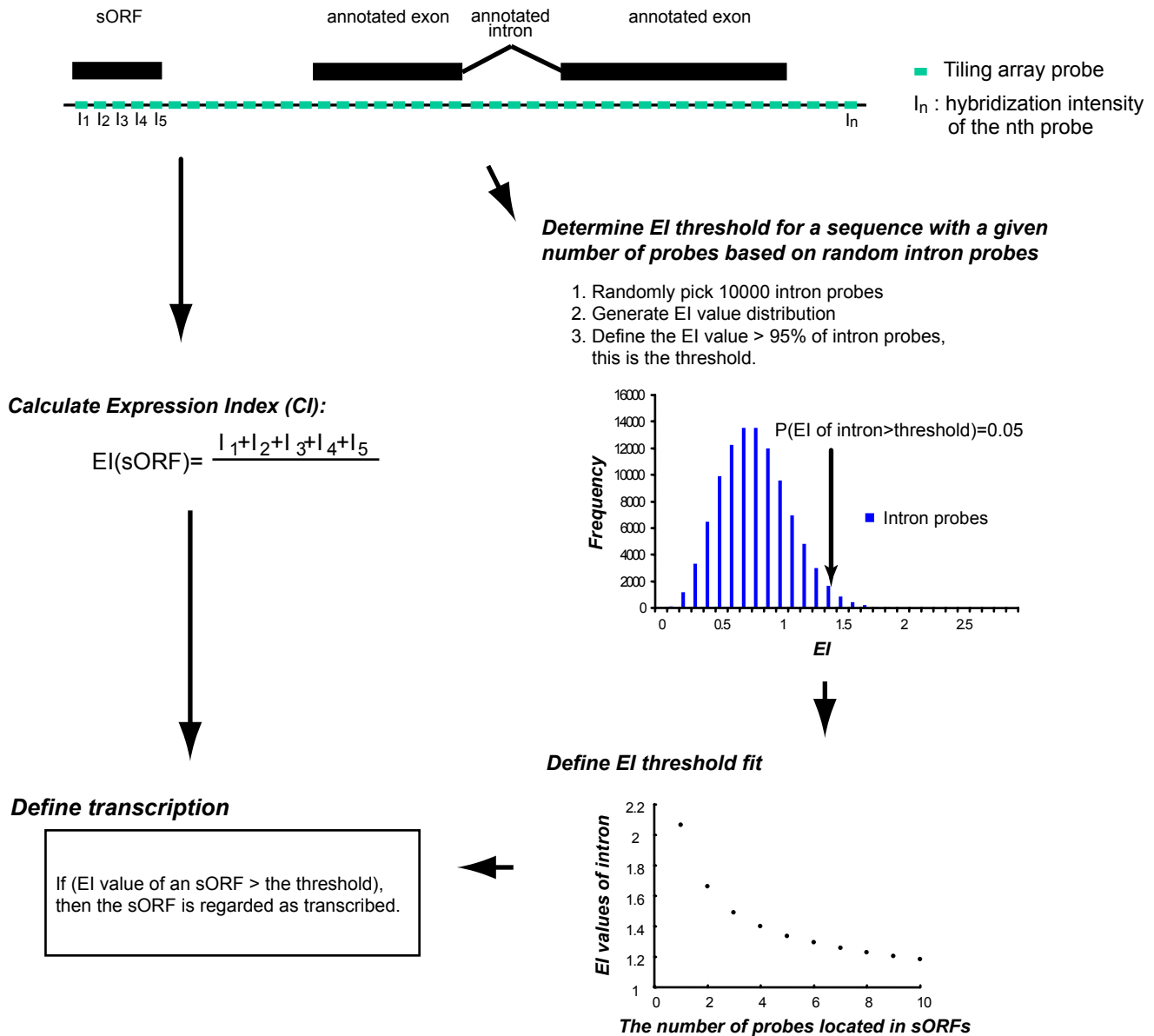
B. Yeast data.

$y = 15.149x^{-0.6493}$

$R^2 = 0.989$

Figure legend: The CI thresholds for (A) Arabidopsis thaliana and for (B) yeast. Each circle represents the CI threshold value of a particular size class at 1% false positive rate. The false positive rates were determined based on the distributions of the CI values of 100,000 random NCDS-like sequences for each size class (ranging from 90-300 nucleotides with 3 nt increment). The CI threshold values are fitted with the power. The equation for the fit and correlation coefficients are shown.

Determination of EI and EI thresholds

sORF          annotated exon   annotated   annotated exon
                               intron

$I_1$ $I_2$ $I_3$ $I_4$ $I_5$                                    $I_n$

■ Tiling array probe

$I_n$ : hybridization intensity
of the nth probe

*Determine EI threshold for a sequence with a given number of probes based on random intron probes*

1. Randomly pick 10000 intron probes
2. Generate EI value distribution
3. Define the EI value > 95% of intron probes, this is the threshold.

*Calculate Expression Index (CI):*

$$EI(sORF) = \frac{I_1 + I_2 + I_3 + I_4 + I_5}{}$$

P(EI of intron>threshold)=0.05

Frequency

■ Intron probes

EI

*Define EI threshold fit*

*Define transcription*

If (EI value of an sORF > the threshold), then the sORF is regarded as transcribed.

EI values of intron

The number of probes located in sORFs

**SUPPLEMENT A-4**

Determination of EI thresholds based on a given number of probes based on random intron probes
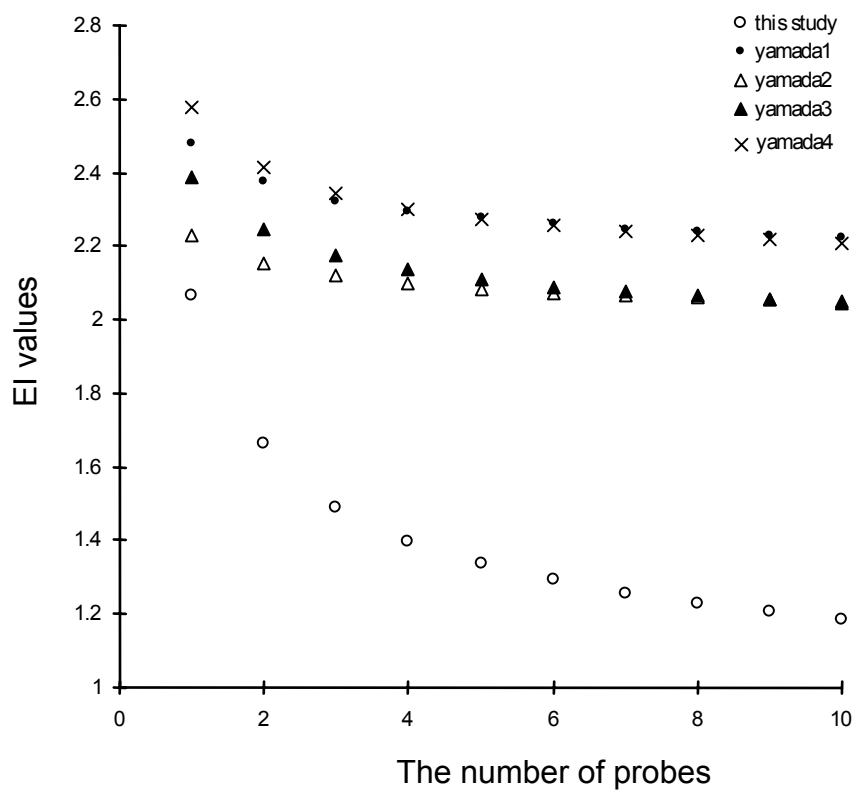


Figure legend: Each point represents the EI threshold value of a given number of probe(s) at 5% false positive rate. The false positive rates were determined based on the EI value distributions of 100,000 randomly sampled intron probes.