

Probe design

for microarrays by H. Bjørn Nielsen



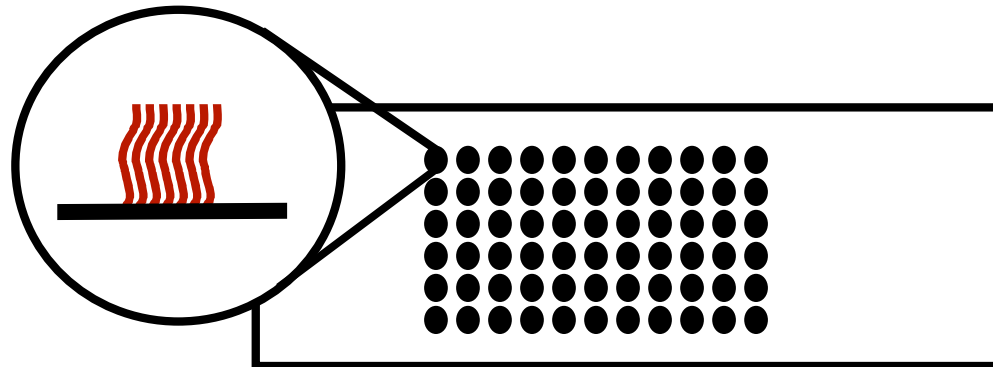
- What is a Probe
- Different Probe Types
- OligoWiz
- Probe Design
 - Cross Hybridization and Complexity
 - Affinity
 - Position

An Ideal Probe

must

CENTER FOR
RIBONUCLEIC ACID
SEQUENCING
ANALYSIS
CBS

- Discriminate well between its intended target and all other targets in the target pool
- Detect concentration differences under the hybridization conditions



Probe Type

comparisons

CENTER FOR
RIBOLOGY
CAL SEQU
ENCE ANA
LYSIS CBS

PCR products

Advantages

Inexpensive
Linkers can be applied

Disadvantages

Handling problems
Hard to design to avoid cross-
hybridization
Unequal amplification

Oligos

Can be designed for many
criteria
Easy to handle
Normalized concentrations
Linkers can be applied

Expensive
(Dkk. 100-150 per oligo)

Affymetrix GeneChip

High quality data
Standardized arrays
Fast to set up
Multiple probes per gene

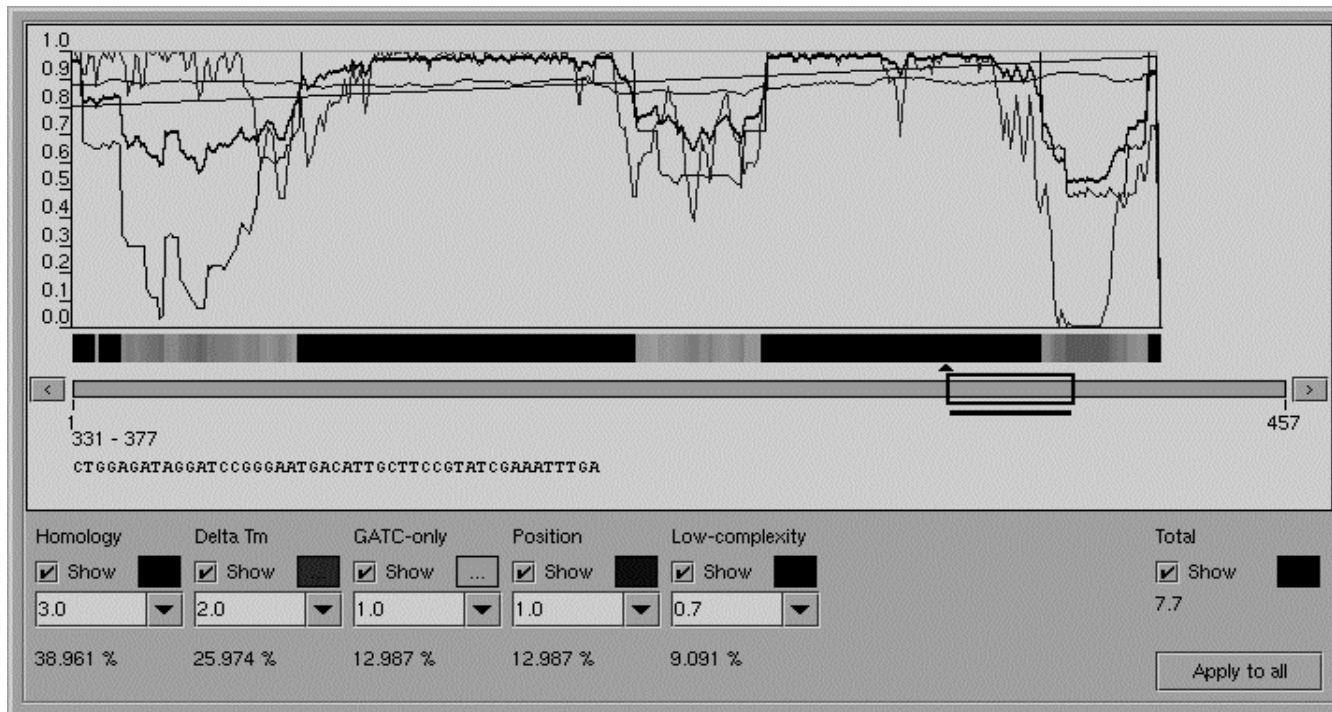
Expensive
Arrays available for limited number
of species

OligoWiz a Tool

for flexible probe design

CENTER FOR
RIBOLOGY
CALSEQU
ENCEANA
LYSIS CBS

A client-server application for designing oligonucleotides for microarrays



Download OligoWiz

and try the sample data set



Go to

www.cbs.dtu.dk/services/OligoWiz/OW2

Linked from course program

Download the

Program: [OligoWiz2.jar](#)

Sample data set: [yeast.fsa](#)

Annotation file: [yeast.ann](#)

Execute OligoWiz2 with the sample data set:

[Double click the icon \(PC/Mac\)](#)

[java -jar OligoWiz2.jar \(UNIX\)](#)

How to Avoid

cross-hybridization



From Kane et al. (2000) we learn that a 50' mer probe can detect significant false signal from a target that has

>75-80% homology to a 50' mer oligo

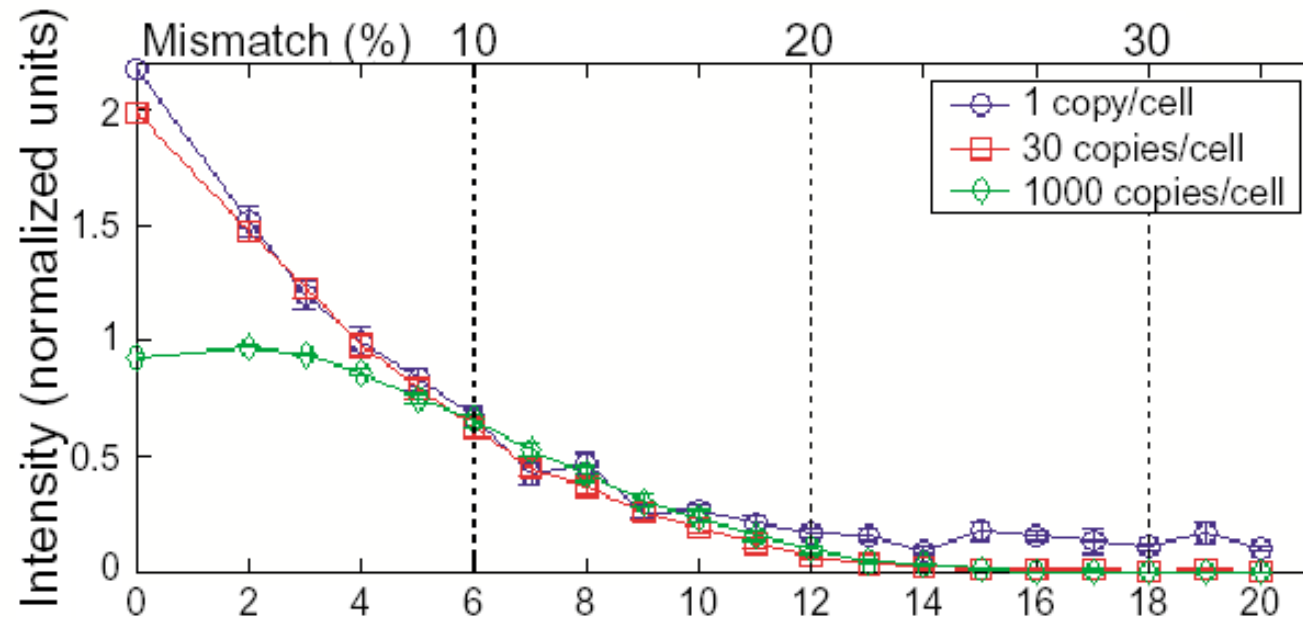
or a continuous stretch of >15 complementary bases

If we have substantial sequence information on the given organism, we can try to avoid this by choosing oligos that are not similar to any other expressed sequences.

Probe Specificity

Hughes *et al.* 2001

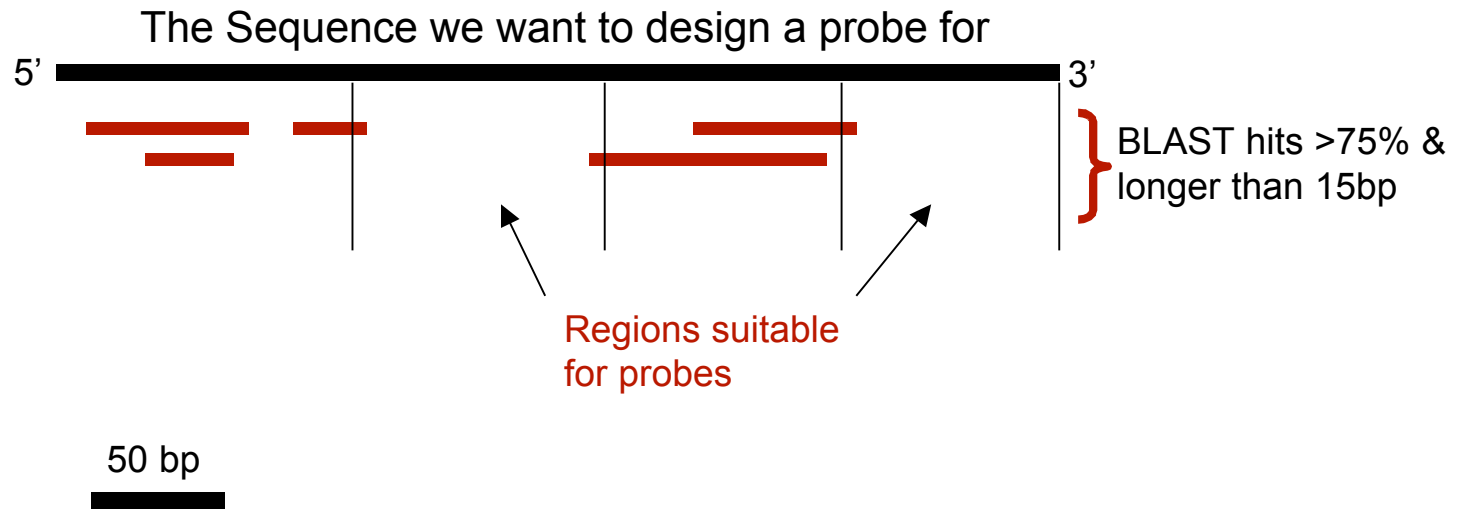
CENTER FOR
MOL
RBIOL
LOGI
CAL
SEQU
ENCE
ANAL
YSIS
CBS



Mapping Regions

without similarity to other transcripts

CENTER FOR
RIBOBIOLOGICAL
SEQUENCE
ANALYSIS
CBS



Filtering Self Detecting

BLAST hits out

CENTERFO
R BIOLOGI
CAL SEQU
ENCE ANA
LYSIS CBS

The Sequence we want to design a oligo for
5' _____ 3'



} BLAST hits >75%
& longer than 15bp

Sequence identical or
very similar to the query
sequence

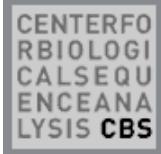
Therefore no BLAST hits with homology > 97% and
with a 'hit length vs. query length' ratio > 0.8,
are considered.

50 bp



Cross-hybridization

expressed as a 'homology score'

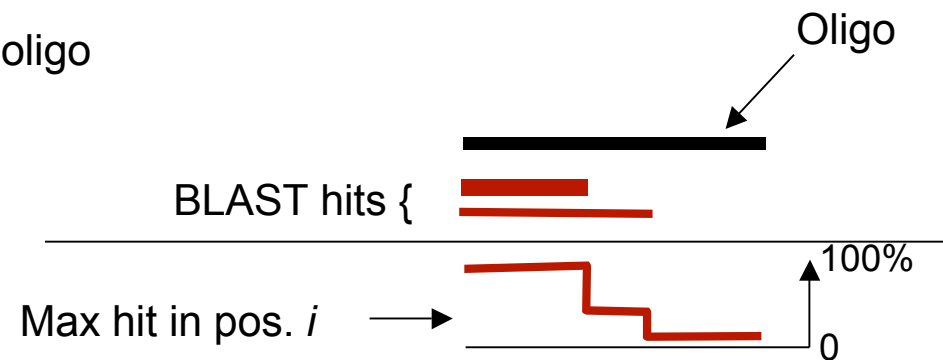


Only BLAST hits that passed filtering are considered

If m is the number of BLAST hits considered in position i .
Let $h=(h_{1_i}, \dots, h_{m_i})$ be the BLAST hits in position i in the oligo

$$\text{BLAST max score} = \frac{100 \times n - \sum_{i=1}^n \max(h_{1_i}, \dots, h_{m_i})}{100 \times n}$$

Where n is the length of the oligo



Similar Affinity

for all oligos



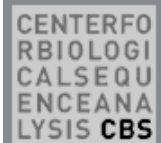
Another way of ensuring an optimal discrimination between target and non-target under hybridization is to design all the oligos on an array with similar affinity for their targets.

This will allow the experimentalist to optimize the hybridization conditions for all oligos by choosing the right hybridization temperature and salt concentration.

Commonly Melting Temperature (T_m) is used as a measure for DNA:DNA or RNA:DNA hybrid affinity.

Melting Temperature

difference



$$T_m(i) = \frac{1000\Delta H}{A + \Delta S + R \ln\left(\frac{C_t}{4}\right)} - 273.15 + 16.6 \log[Na^+]$$

Where ΔH (Kcal/mol) is the sum of the nearest neighbor enthalpy, A is a constant for helix initiation corrections, ΔS is the sum of the nearest neighbor entropy changes, R is the Gas Constant (1.987 cal deg⁻¹ mol⁻¹) and C_t is the total molar concentration of strands.

$$\Delta T_m \text{ score} = \left| T_m(i) - \frac{\sum_N^1 T_m(i)}{N} \right|$$

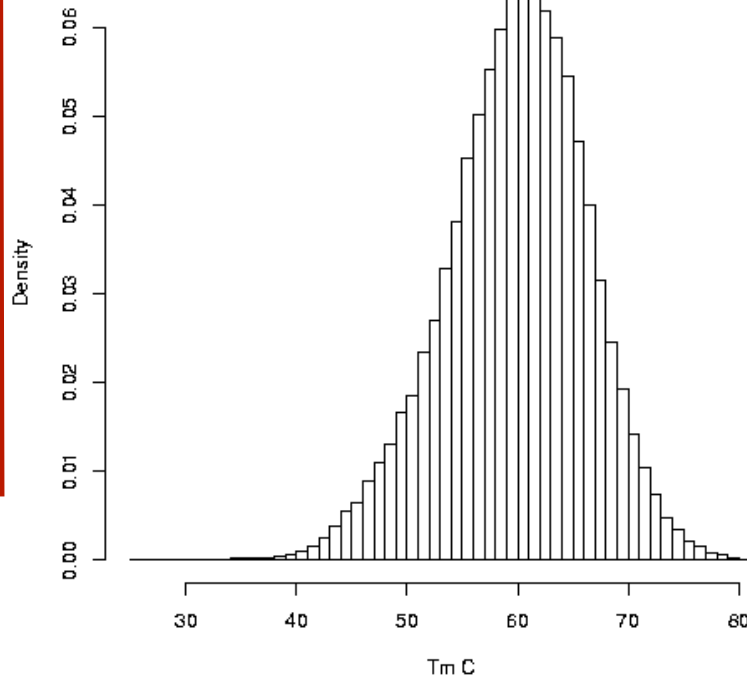
Where N is all oligos in all sequences.

Tm distributions

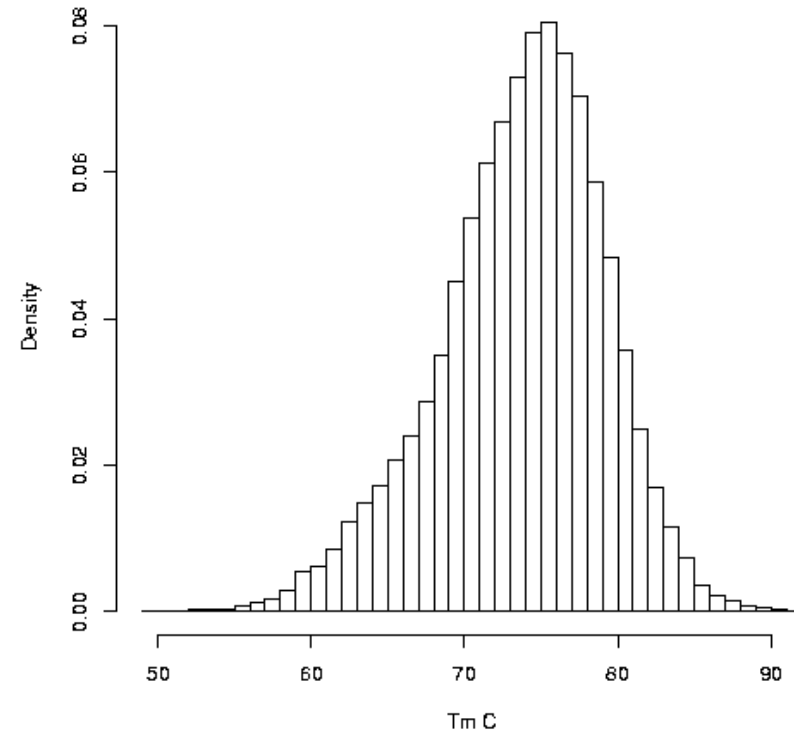
for 30'mers and 50'mers

CENTER FOR
RIBOLOGY
CALSEQU
ENCEANA
LYSIS CBS

Tm distribution for 30'mers (yeast)



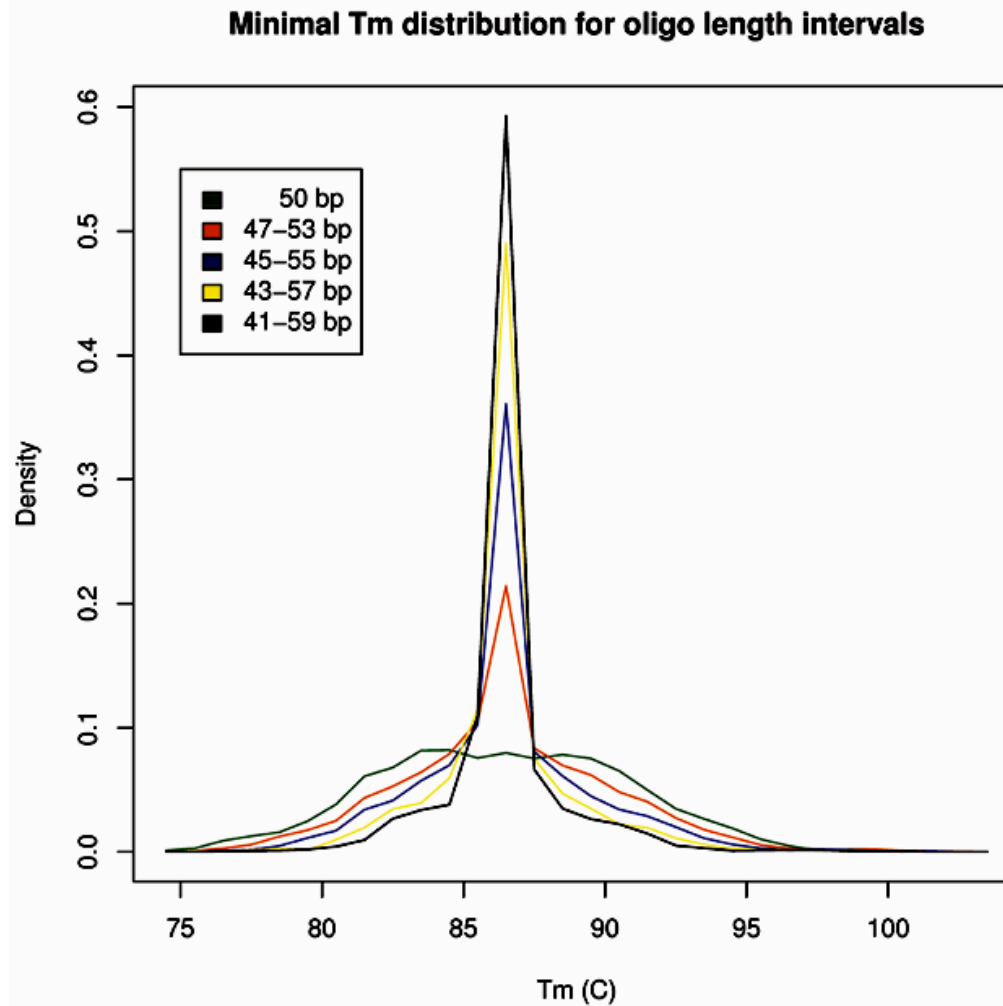
Tm distribution for 50'mers (yeast)



ΔT_m Distribution

for oligo length intervals

CENTER FOR
RIBIOLOGICAL
CALSEQU
ENCEANA
LYSIS CBS



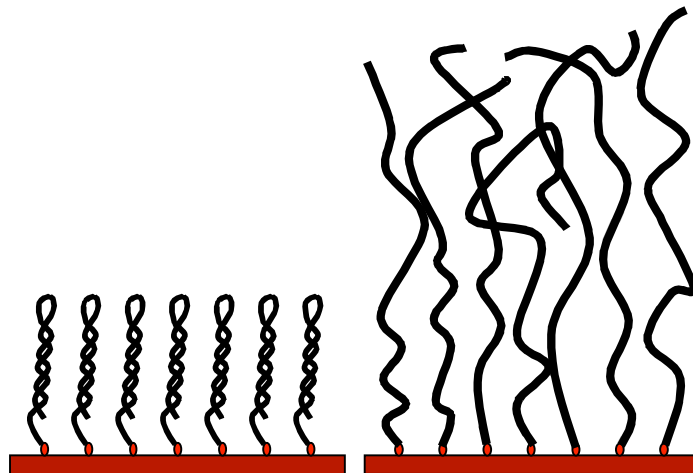
Avoid self annealing oligos

Sensitivity may be influenced

CENTER FOR
RIBIOMOLECULAR
CALCULATIONAL
ANALYSIS
LYSIS CBS

Probes that form strong hybrids with it self *i.e.* probes that fold should be avoided.

But, accurate folding algorithms like the one employed by mFOLD or RNAfold, is too time consuming, for large scale folding of oligos.



Time consumption:
mFOLD ~2 sec / 30' mer
Pr. gene (500bp) ~16 min.

Folding an oligonucleotide

an approximation

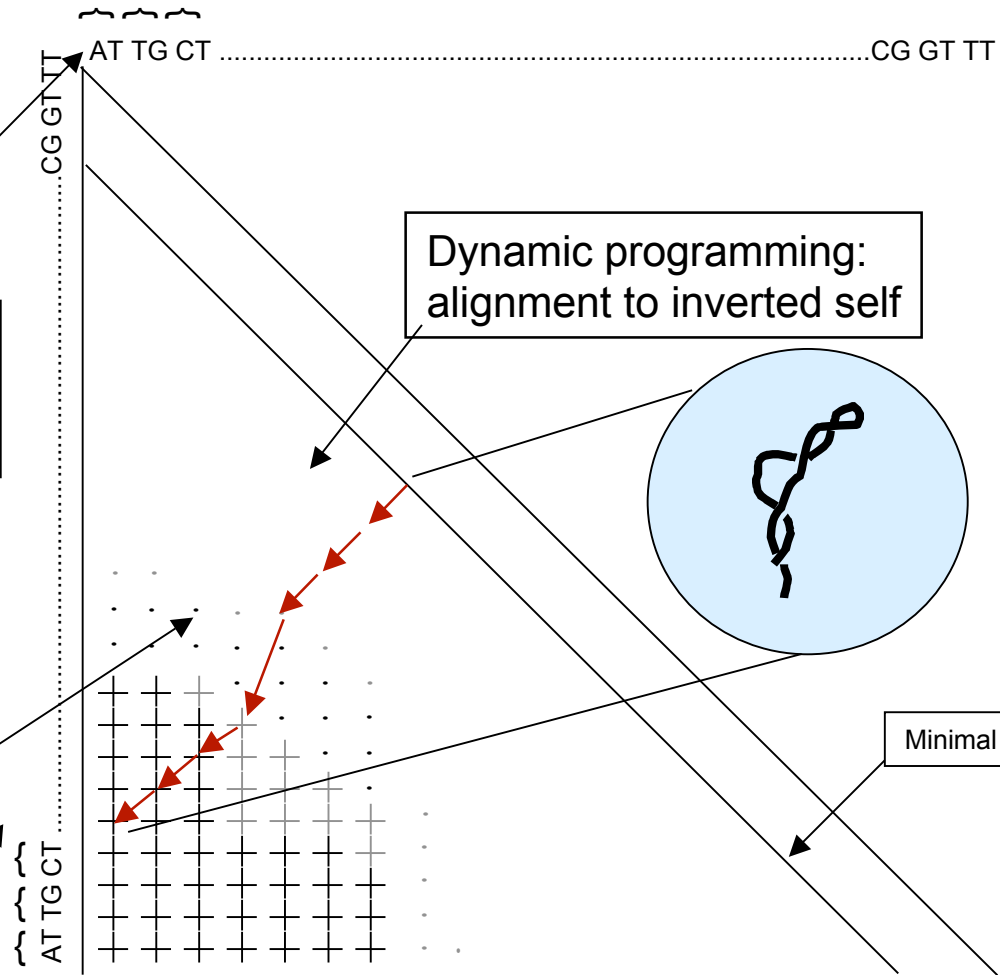
CENTER FOR
RIBONUCLEIC ACID
SEQUENCE ANALYSIS
CBS

The alignment is based on dinucleotides

Substitution matrix is based on binding energies

Dynamic programming:
alignment to inverted self

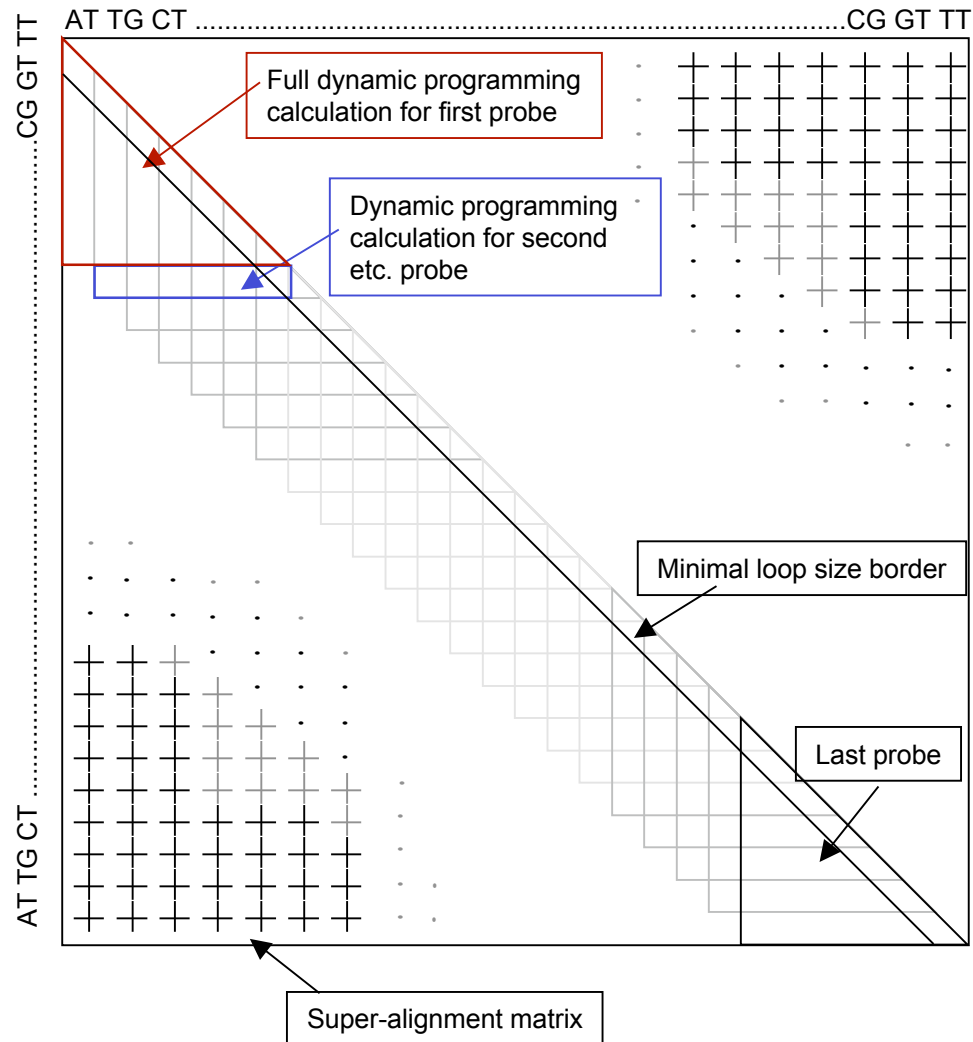
Minimal loop size border



Folding a lot of oligos

a fast heuristic implementation

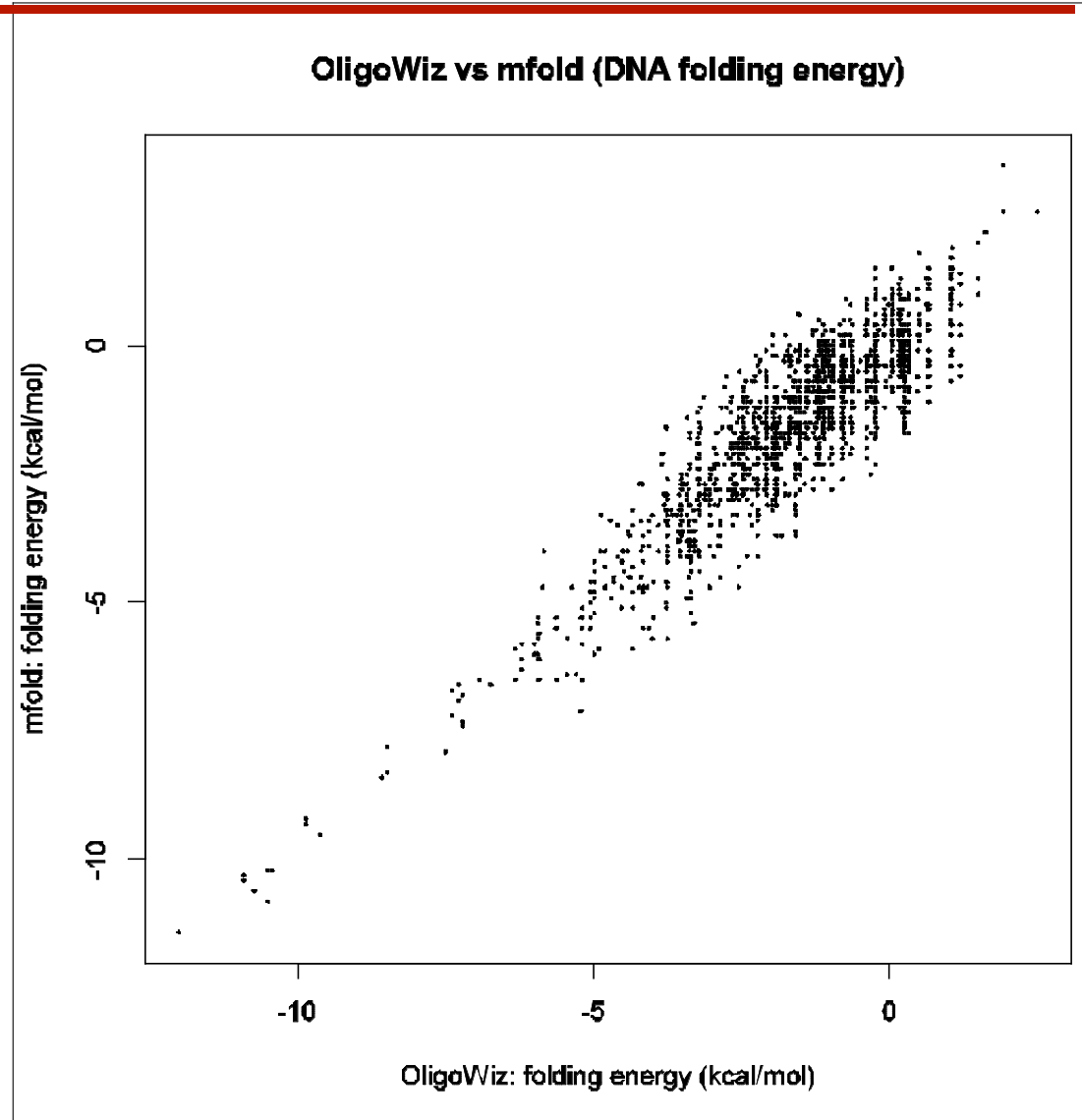
CENTER FOR
RIBBIOLOGICAL
SEQUENCE
ANALYSIS
CBS



Reasonably folding prediction

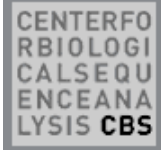
compared to mFOLD

CENTER FOR
RIBONUCLEIC ACID
SEQUENCE ANALYSIS
CBS



Probes With Very Common

sub sequences may result in unspecific signal



If the sub-fractions of an oligo are very common we define it as 'low-complex'

Oligo with low-complexity:

AAAAAAGGAGTTTTTTTTCAAAAACTTTTTAAAAAGCTTTAGGTTTTTA
(Human)

Oligo without low-complexity:

CGTGACTGACAGCTGACTGCTAGCCATGCAACGTCATAGTACGATGACT
(Human)

Low-complexity

expressed as a score

CENTER FOR
RIBIOLOGICAL
CALSEQUEN
CEANALYSIS
CBS

For a given transcriptome a list of information content from all 'words' with length wl (8bp) is calculated:

$$I(w) = \frac{f(w)}{tf(w)} \log_2 \frac{f(w)}{tf(w)} 4^{wl}$$

Where $f(w)$ is the number of occurrences of a pattern and $tf(w)$ is the total number of patterns of length wl .

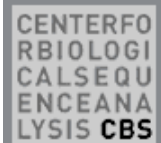
A low-complexity score for a given oligo is defined as:

$$\text{Low-complexity} = 1 - \text{norm} \left(\sum_{i=1}^{L-wl+1} I(w_i) \right)$$

Where **norm** is a function that normalizes to between 1 and 0, L is the length of the oligo and w_i is the pattern in position i .

Location of Oligo

within transcript



Labeling include reverse transcription of the mRNA and is sensitive to:

- RNA degradation
- Premature termination of cDNA synthesis
- Premature termination of cRNA transcription (IVT)

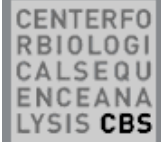
A 'Position Score' reflecting this (eukaryotes):

$$\text{Position score} = (1 - drp)^{\Delta 3'end}$$

Where *drp* is the chance of labeling termination pr. base

Placement of multiple oligos

in respect to annotation



Load annotation file:

Project/project properties : annotation file

Place 10 oligos in each gene:

Tools/oligo placement

Export the probes

to a file



Export the oligos to a FASTA file:

File/Export oligo

Have a look at it

Place probes

relative to annotation



Use the **Tools/oligo placement** tool, to place Oligos:

- 1) Exclusively in exons
- 2) Exclusively in introns
- 3) Exclusively on the exon-intron boundary

Using regular expressions

The “Introduction to regular expressions” Is linked from the course program.

Note: the yeast.ann file consists of lists of: {E,I,D,A,),(,.,. }
E=exon, I=intron, D= donor site, A=acceptor site,)(= end and beginning of exon and “.” = UTR