# KIRMES

## Kernel-based Identification of Regulatory Modules in Euchromatic Sequences

GCB 2008 Dresden                 September 10, 2008
Sebastian J. Schultheiss  <sebi@tuebingen.mpg.de>
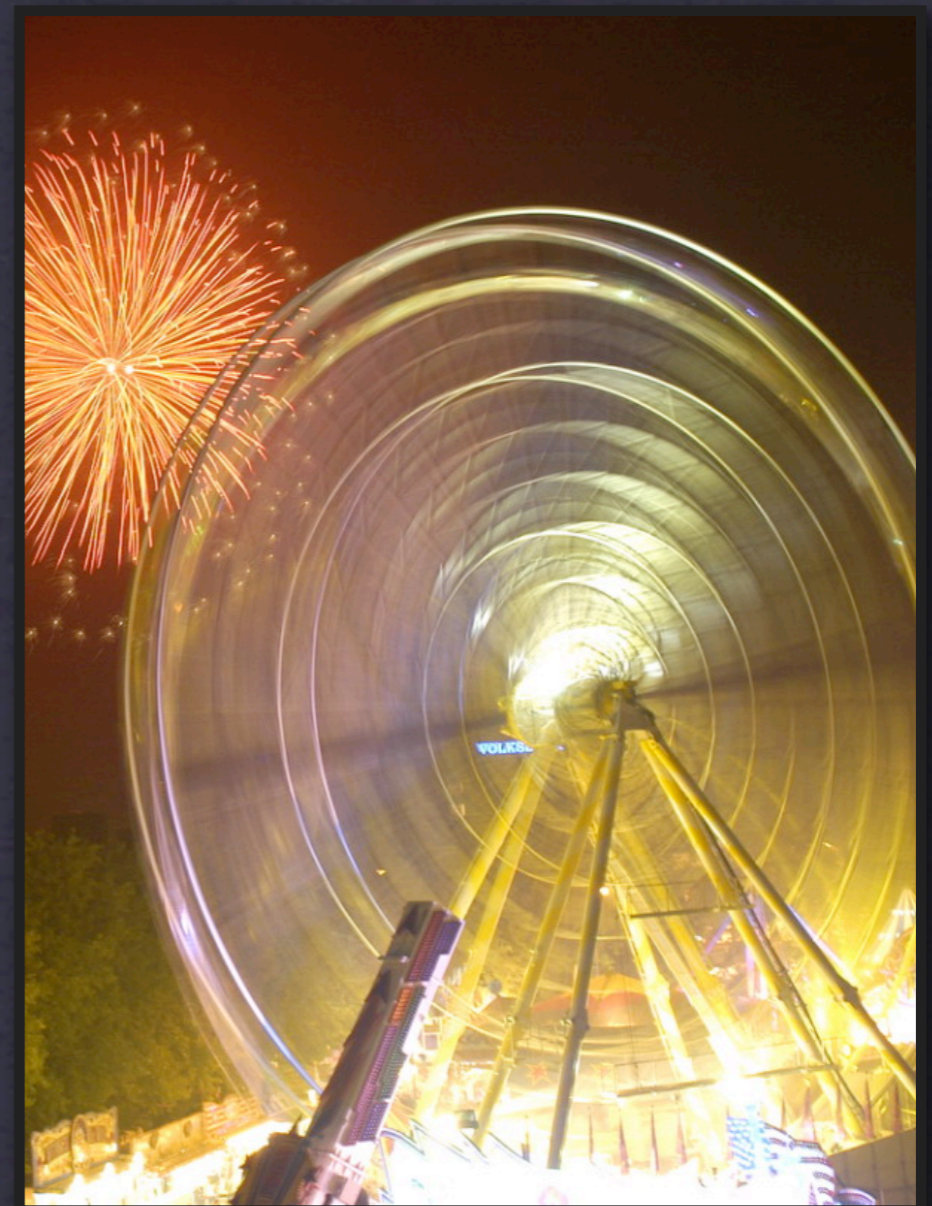Wolfgang Busch, Jan U. Lohmann, Oliver Kohlbacher, and Gunnar Rätsch

Friedrich Miescher Laboratory
of the Max Planck Society

MAX-PLANCK-GESELLSCHAFT

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

ZBIT

# KIRMES Overview

▸ Combine motif finding with SVMs

▸ Model degenerate motifs and regulatory modules of promoters

▸ **Input:** sets of co-expressed genes

▸ **Output:** a function that classifies genes as co-regulated or not



First hit in Google Image Search for "KIRMES"
bocholt.de

# Transcriptional Regulation

▸ **Transcription factors** (TFs) are proteins regulating transcription (activate, repress, ...)

▸ TFs bind to euchromatic sequences, *e.g.* promoter regions, **regulatory regions** in introns, ...

▸ TFs recognize conserved **binding motifs** (TFBMs)

▸ TFs form complexes, **bind** to **modules** of TFBMs

▸ **Relative distance** of TFBMs to transcription start and among TFBMs is relevant

TFBM A          TFBM B          TFBM C   TFBM B          **Transcription Start**

# Transcriptional Regulation

- **Transcription factors** (TFs) are proteins regulating transcription (activate, repress, ...)

- TFs bind to euchromatic sequences, *e.g.* promoter regions, **regulatory regions** in introns, ...

- TFs recognize conserved **binding motifs** (TFBMs)

- TFs form complexes, **bind** to **modules** of TFBMs

- **Relative distance** of TFBMs to transcription start and among TFBMs is relevant

TF A

TFBM A          TFBM B                TFBM C   TFBM B          Transcription Start

# Transcriptional Regulation

▸ **Transcription factors** (TFs) are proteins regulating transcription (activate, repress, ...)

▸ TFs bind to euchromatic sequences, *e.g.* promoter regions, **regulatory regions** in introns, ...

▸ TFs recognize conserved **binding motifs** (TFBMs)

▸ TFs form complexes, **bind** to **modules** of TFBMs

▸ **Relative distance** of TFBMs to transcription start and among TFBMs is relevant



TF A

TF C

**TFBM A**    **TFBM B**    **TFBM C  TFBM B**    **Transcription Start**

# Transcriptional Regulation

▶ **Transcription factors** (TFs) are proteins regulating transcription (activate, repress, ...)

▶ TFs bind to euchromatic sequences, *e.g.* promoter regions, **regulatory regions** in introns, ...

▶ TFs recognize conserved **binding motifs** (TFBMs)

▶ TFs form complexes, **bind** to **modules** of TFBMs

▶ **Relative distance** of TFBMs to transcription start and among TFBMs is relevant



TF A

TF C TF B

TFBM A            TFBM B            TFBM C    TFBM B

Transcription Start

# Machine Learning Method

▶ SVM for a **2-class problem**: genes are **co-regulated** by the same mechanism or not

▶ Compare region similarity

  ▶ Straightforward: string kernel looks at whole regulatory region

  ▶ Our method: windows around motif positions, relative distances, conservation

**Transcription Start**

# Machine Learning Method

▶ SVM for a **2-class problem**: genes are **co-regulated** by the same mechanism or not

▶ Compare region similarity

   ▶ Straightforward: string kernel looks at whole regulatory region

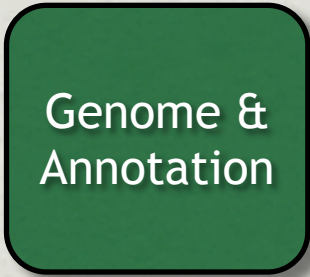   ▶ Our method: windows around motif positions, relative distances, conservation



**Transcription Start**

# Machine Learning Method

▶ SVM for a **2-class problem**: genes are **co-regulated** by the same mechanism or not

▶ Compare region similarity

  ▶ Straightforward: string kernel looks at whole regulatory region

  ▶ Our method: windows around motif positions, relative distances, conservation

**Transcription Start**

# Experimental Data

▶ Microarrays of *Arabidopsis thaliana*

▶ Statistical methods identify **co-expressed** genes in (several) experimental conditions

▶ Which are co-regulated?
Unclear for many genes

    ▶ not on chip

    ▶ change below detection threshold

▶ **Classifier** needed, train on co-expressed genes

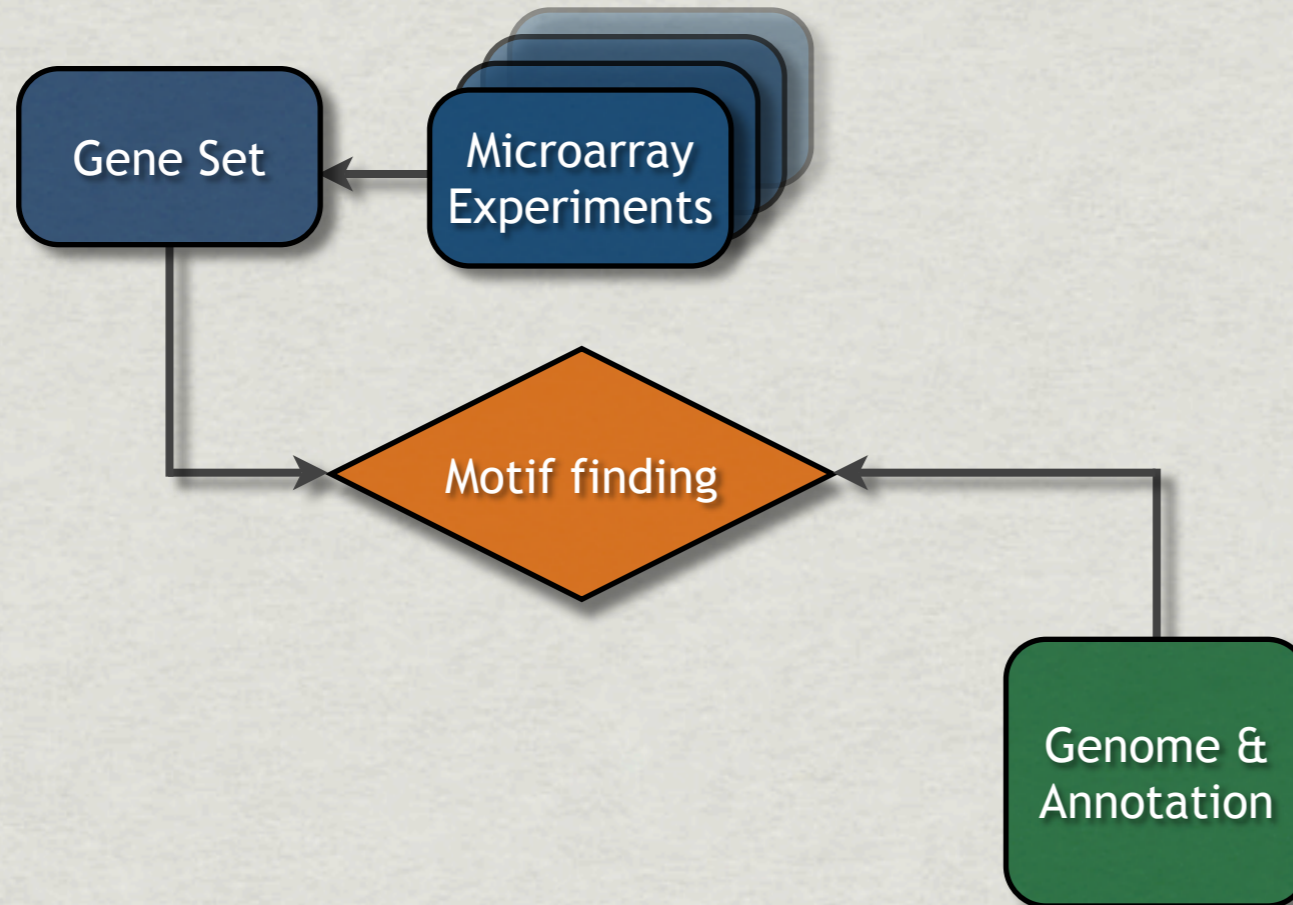▶ **Predict co-regulation** for whole genome

# KIRMES Workflow

Microarray
Experiments
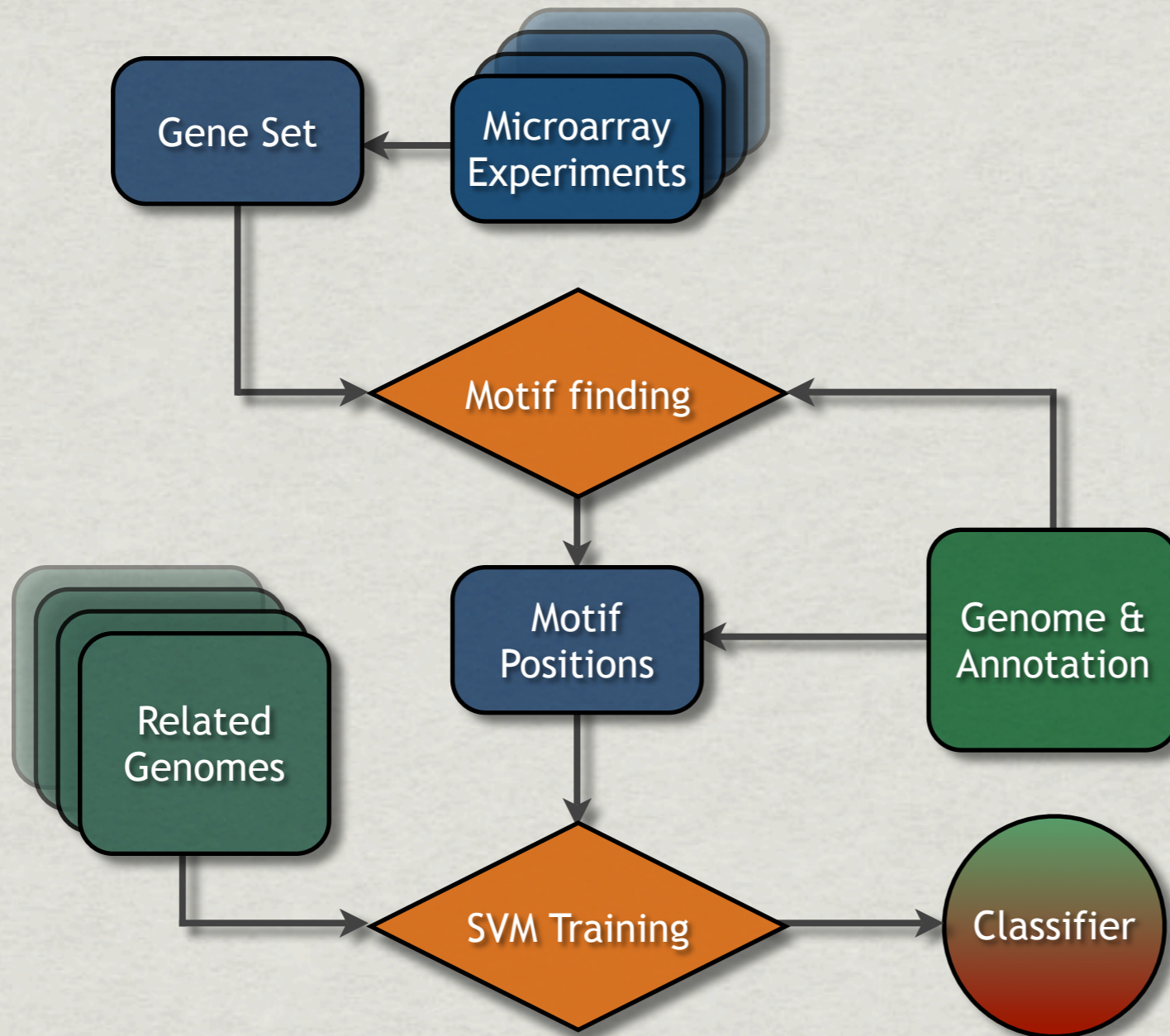
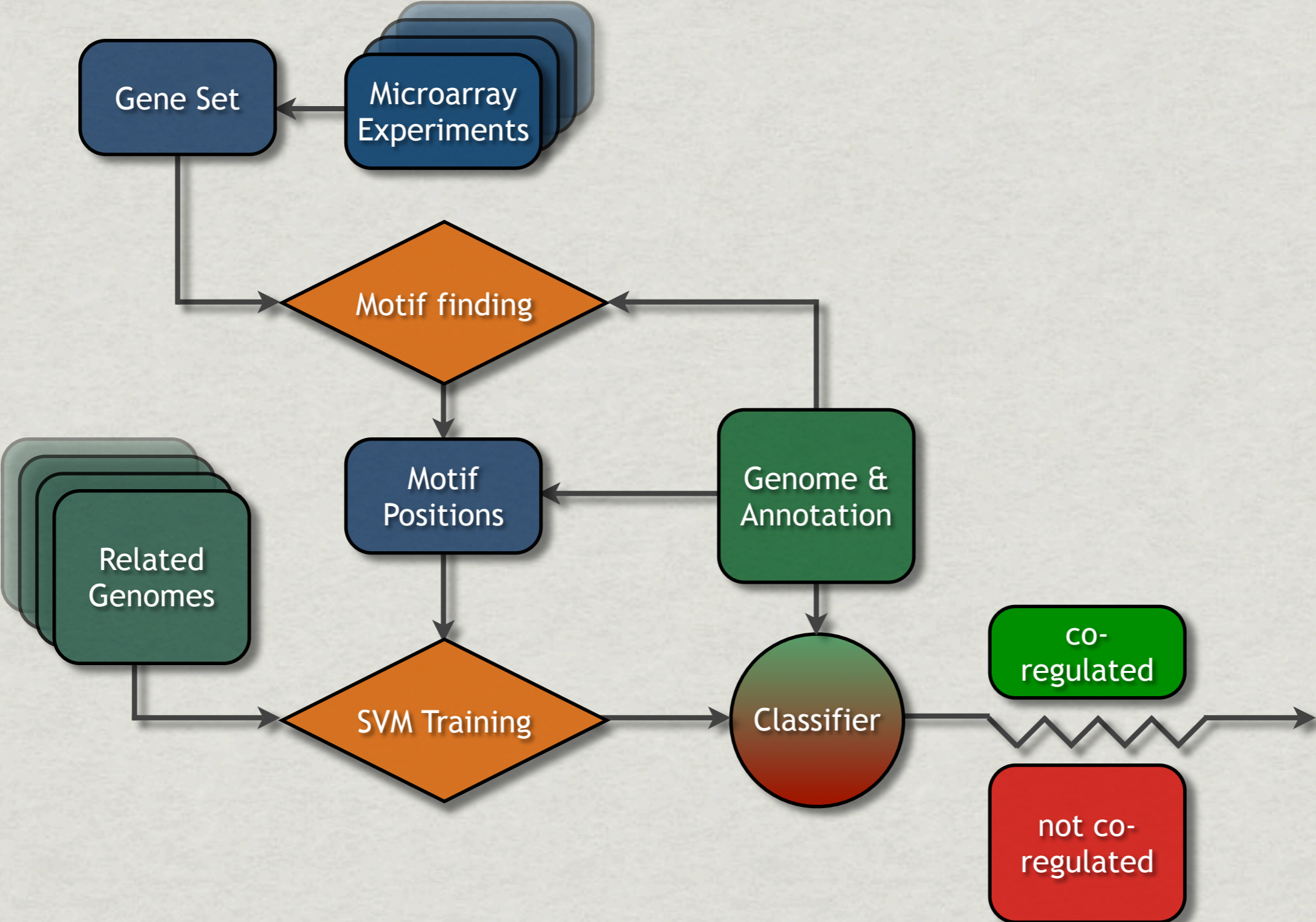Genome &
Annotation

# KIRMES Workflow

# KIRMES Workflow

# KIRMES Workflow

# KIRMES Workflow

# KIRMES Step by Step

1. Obtain sets of co-expressed genes

2. Select regulatory regions

3. Identify over-represented motifs

4. Learn to classify genes as co-regulated

5. Predict co-regulation for all genes

# Co-Expressed Genes

**1.**

▶ Obtain sets from *e.g.* microarray experiments

▶ Genes can be up- or down-regulated

▶ Obtain **negative** control set: genes invariantly expressed at high levels

# Co-Expressed Genes

**1.**

▶ Obtain sets from *e.g.* microarray experiments

▶ Genes can be up- or down-regulated

▶ Obtain **negative** control set: genes invariantly expressed at high levels

**set of gene sequences**

# Regulatory Regions

- ▶ Core **promoter**: 400 bp upstream from transcription start site, up to 3000 bp

- ▶ **Untranslated** regions, up- and downstream

- ▶ 1st **intron**, all introns, 1st exon, all exons

- ▶ **Downstream** region, *e.g.* 500 bp

**2.**

# Regulatory Regions

**2.**

▸ Core **promoter**: 400 bp upstream from transcription start site, up to 3000 bp

▸ **Untranslated** regions, up- and downstream

▸ 1$^{st}$ **intron**, all introns, 1$^{st}$ exon, all exons

▸ **Downstream** region, *e.g.* 500 bp

gene     selected gene     gene

# Regulatory Regions

**2.**

- ▶ Core **promoter**: 400 bp upstream from transcription start site, up to 3000 bp

- ▶ **Untranslated** regions, up- and downstream

- ▶ 1st **intron**, all introns, 1st exon, all exons

- ▶ **Downstream** region, *e.g.* 500 bp

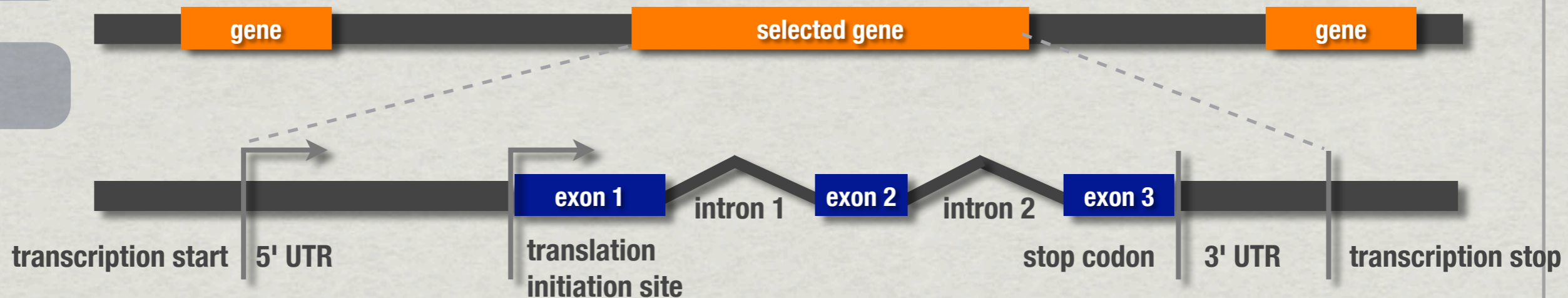gene    selected gene    gene

# Regulatory Regions

▸ Core **promoter**: 400 bp upstream from transcription start site, up to 3000 bp

▸ **Untranslated** regions, up- and downstream

▸ 1st **intron**, all introns, 1st exon, all exons

▸ **Downstream** region, *e.g.* 500 bp

**2.**

gene | selected gene | gene

exon 1 | intron 1 | exon 2 | intron 2 | exon 3

transcription start | 5' UTR | translation initiation site | stop codon | 3' UTR | transcription stop

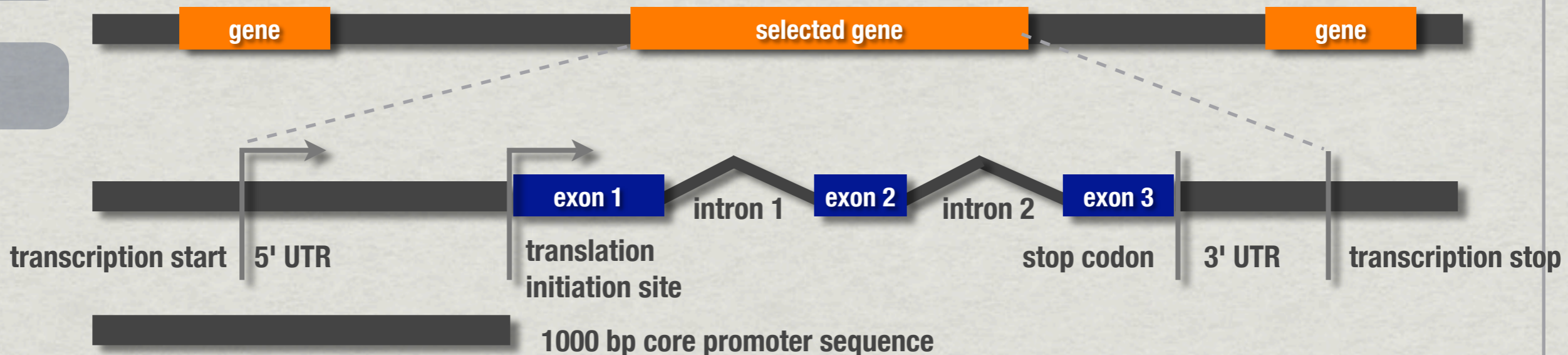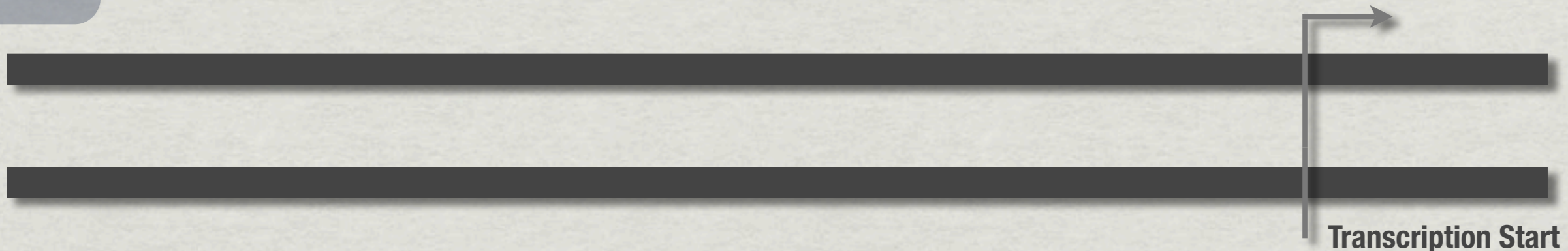# Regulatory Regions

**2.**

- ▶ Core **promoter**: 400 bp upstream from transcription start site, up to 3000 bp

- ▶ **Untranslated** regions, up- and downstream

- ▶ 1st **intron**, all introns, 1st exon, all exons

- ▶ **Downstream** region, *e.g.* 500 bp



gene | selected gene | gene

exon 1 | intron 1 | exon 2 | intron 2 | exon 3

transcription start | 5' UTR | translation initiation site | stop codon | 3' UTR | transcription stop

1000 bp core promoter sequence

# Identify over-represented motifs

- **Gibbs sampling** (Lawrence *et al.*, 1993)
  - identifies over-represented **weight matrices** that characterize TFBMs
- **Oligo counting** (faster)
  - count occurrences of all nucleotide sequences of length six
- Relative position, window around motifs constitute **SVM input vector**

**Transcription Start**

# Identify over-represented motifs

▶ **Gibbs sampling** (Lawrence *et al.*, 1993)
  ▶ identifies over-represented **weight matrices** that characterize TFBMs
▶ **Oligo counting** (faster)
  ▶ count occurrences of all nucleotide sequences of length six
▶ Relative position, window around motifs constitute **SVM input vector**

3.

**Transcription Start**

# Identify over-represented motifs

- **Gibbs sampling** (Lawrence *et al.*, 1993)
  - identifies over-represented **weight matrices** that characterize TFBMs
- **Oligo counting** (faster)
  - count occurrences of all nucleotide sequences of length six
- Relative position, window around motifs constitute **SVM input vector**

**3.**

**Transcription Start**
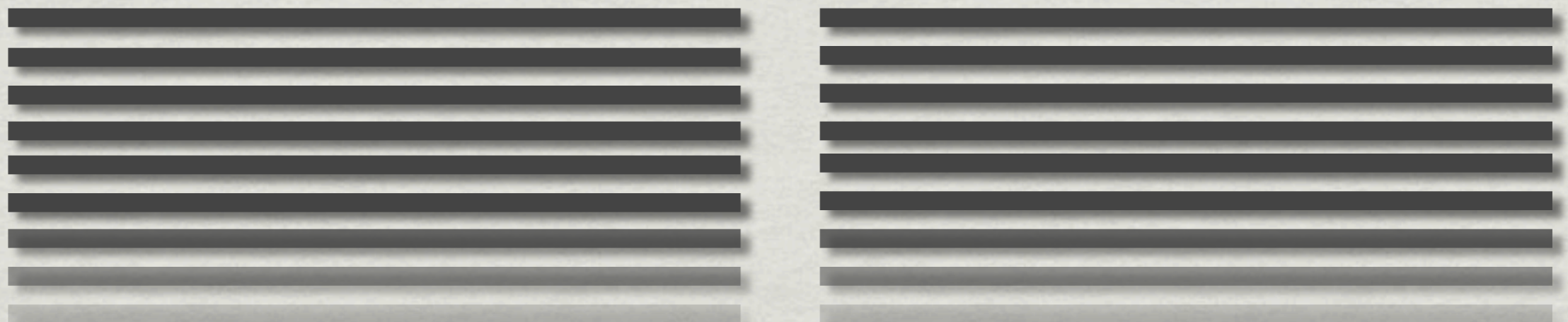
# Learn to Classify Genes

▶ Labeled training data sets with genes **coexpressed** on microarrays

▶ **Classification** output:
function that can classify all genes

▶ Optional output: top-ranking motifs

**4.**

# Learn to Classify Genes

▶ Labeled training data sets with genes **coexpressed** on microarrays

▶ **Classification** output:
function that can classify all genes

**4.**

▶ Optional output: top-ranking motifs

Classifier

# Predict Co-Regulation

**5.**

▶ Learning step creates classifier

▶ **Prediction of co-regulation**

  ▶ All other genes, not on chip

  ▶ Different regions than training data

Classifier

# Predict Co-Regulation

▶ Learning step creates classifier

▶ **Prediction of co-regulation**

   ▶ All other genes, not on chip

   ▶ Different regions than training data

**5.**

Classifier

# Machine Learning Methods

▶ String kernel: **Weighted Degree kernel**

▶ Uses whole regulatory region: can't identify modules, no positional interdependence

```
...AGTCAGATAGAGGACATCAGTAGACAGATTAAA...
       ||||||||      ||  ||      |||
...TTATAGATAGACAAAGACATCAGTAGACTTATT...
```

# Machine Learning Methods

▸ String kernel: **Weighted Degree kernel**

▸ Uses whole regulatory region: can't identify modules, no positional interdependence

```
...AGTCAGATAGAGGACATCAGTAGACAGATTAAA...
       ||||||||      ||   ||       |||
...TTATAGATAGACAAAGACATCAGTAGACTTATT...
```

▸ **With shifts**: still limited, no modules

```
...CGAACGCTACGTATTTTAGTCGGATTCGC...
    \\\\\\            ///////   ///
...TCGAACGAAAGGTTTTAGCCTGATGACGG...
```
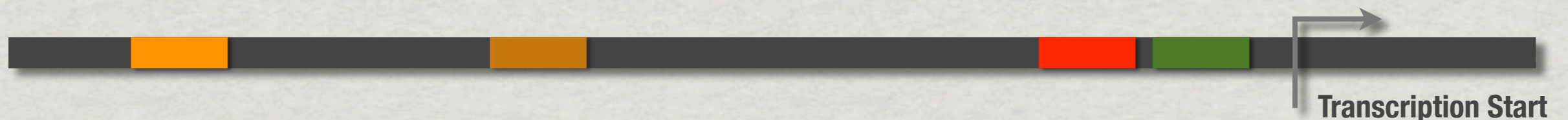
# Novel SVM Kernels

▶ Idea: compare **characteristic** sequence

▶ **Motif finder** identifies TFBMs over background

▶ Novel kernels:

  ▶ Relative **positional** information

  ▶ **Window** around motif positions

  ▶ Sequence **conservation**

▶ Sum up kernels (concatenates feature space)

**Transcription Start**

# Novel SVM Kernels

▸ Idea: compare **characteristic** sequence

▸ **Motif finder** identifies TFBMs over background

▸ Novel kernels:

  ▸ Relative **positional** information

  ▸ **Window** around motif positions

  ▸ Sequence **conservation**

▸ Sum up kernels (concatenates feature space)
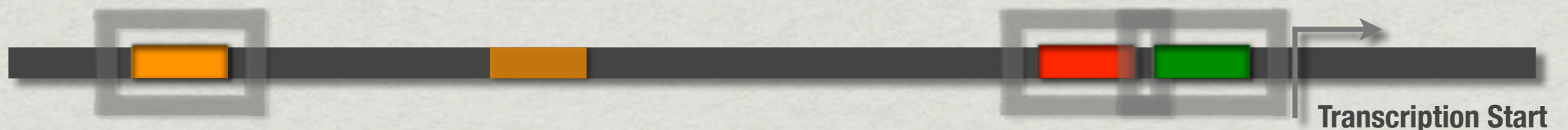
**Transcription Start**

# Novel SVM Kernels

▸ Idea: compare **characteristic** sequence

▸ **Motif finder** identifies TFBMs over background

▸ Novel kernels:

  ▸ Relative **positional** information

  ▸ **Window** around motif positions

  ▸ Sequence **conservation**

▸ Sum up kernels (concatenates feature space)

**Transcription Start**

# Novel SVM Kernels

▸ Idea: compare **characteristic** sequence

▸ **Motif finder** identifies TFBMs over background

▸ Novel kernels:

  ▸ Relative **positional** information

  ▸ **Window** around motif positions

  ▸ Sequence **conservation**

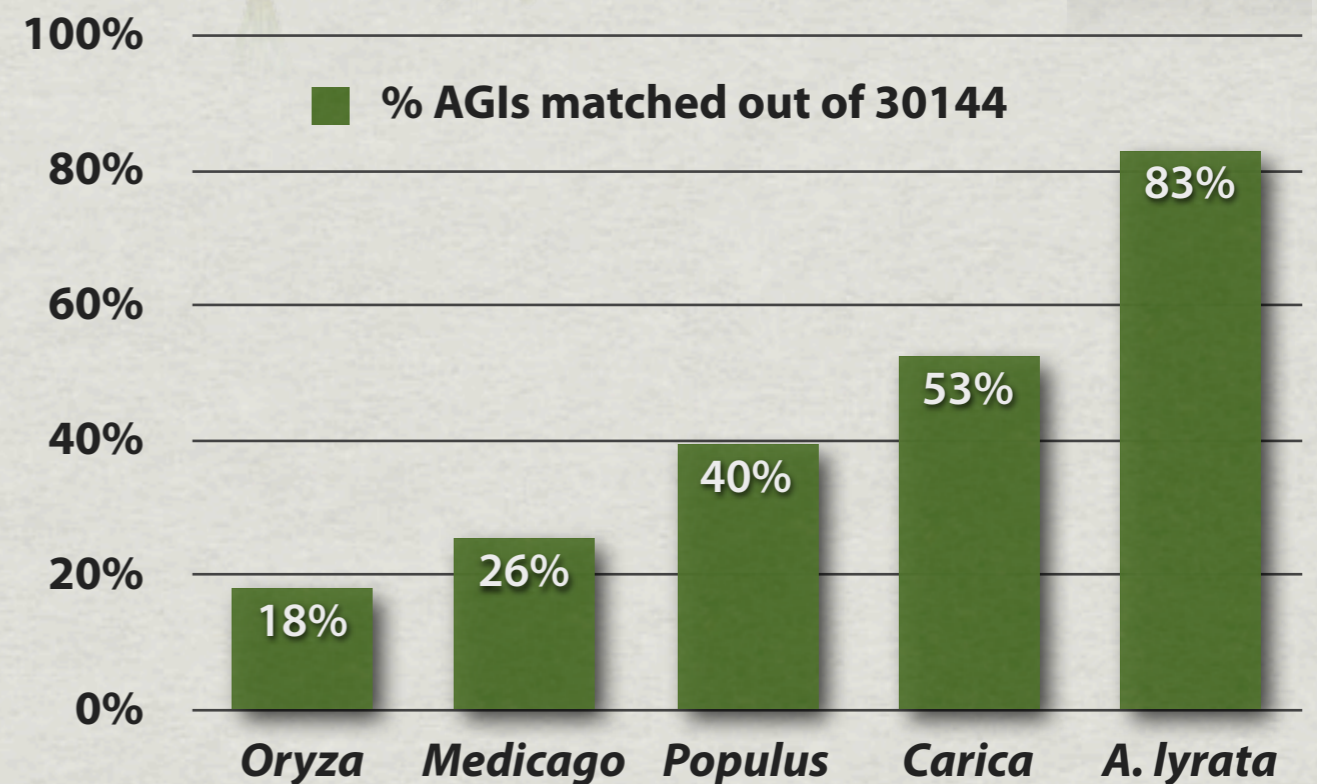▸ Sum up kernels (concatenates feature space)

**Transcription Start**

# Kernel for Regulatory Modules

▸ **RBF kernel**

  ▸ Relative position of motif matches to start and among each other

  ▸ Pairwise distances between match positions of motifs

▸ **Weighted degree kernel with shifts**

  ▸ Set of windows around best-matching motifs

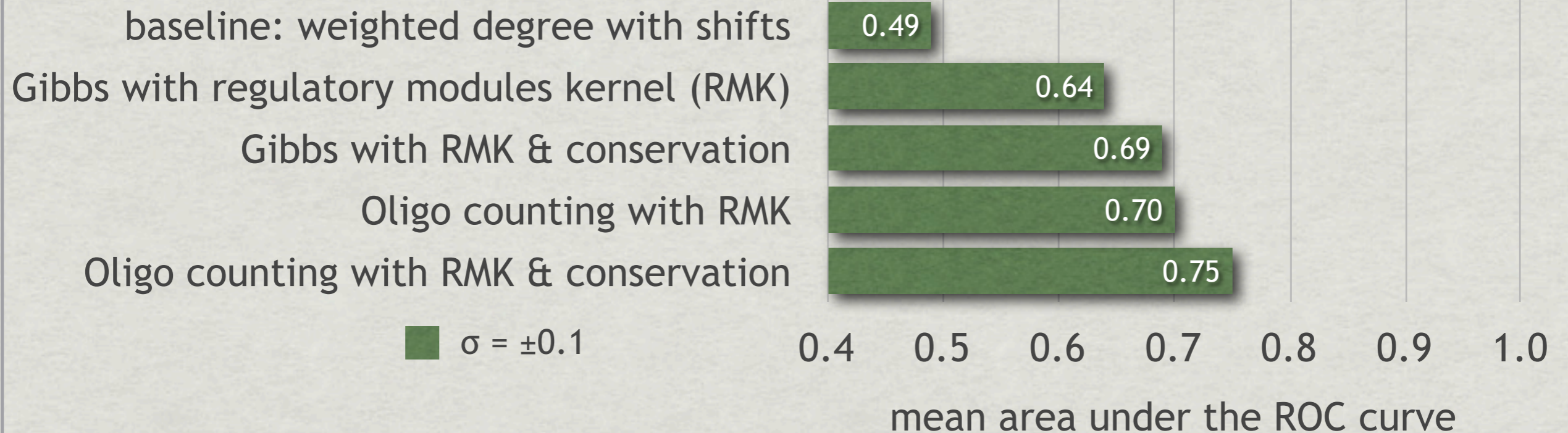  ▸ Highest similarity: all motifs appear in both sequences

# WDS Kernel with Conservation

- ▸ WDS kernel extended with conservation info

- ▸ Helpful, not essential

- ▸ *A. thaliana* AGI regions aligned with other plants

- ▸ How many of *A. thaliana*'s genes were aligned?



Bar chart — % AGIs matched out of 30144:

| Species | % AGIs matched |
|---------|----------------|
| Oryza | 18% |
| Medicago | 26% |
| Populus | 40% |
| Carica | 53% |
| A. lyrata | 83% |

# Results (Comparison)

▸ Sets of co-expressed genes from *A. thaliana* under different conditions, **auROC curve** evaluation

▸ Baseline method performs close to random guessing

▸ Oligo counting outperforms Gibbs



baseline: weighted degree with shifts — 0.49
Gibbs with regulatory modules kernel (RMK) — 0.64
Gibbs with RMK & conservation — 0.69
Oligo counting with RMK — 0.70
Oligo counting with RMK & conservation — 0.75

σ = ±0.1

0.4  0.5  0.6  0.7  0.8  0.9  1.0

mean area under the ROC curve

# Results (Biology)

▸ **WUSCHEL**, a key player in *Arabidopsis* stem cell maintenance

▸ Binding site part of KIRMES output

▸ TFBM verified* (SELEX, gel shift)

▸ Regulatory network: ChIP on chip binding

    ▸ **PWM**: 17% of genes

    ▸ **KIRMES**: 64% of genes



*A. thaliana* **wild type**     WUSCHEL overexpressing mutant
*W. Busch *et al.*, in preparation

# Discussion

- Powerful approach
  - exploits relationships between motifs
  - uses **modules** for prediction
- **Conservation** is useful
- **Oligo counting** works surprisingly well

- Future directions:
  - Integrate binding data
  - Compare against established methods
  - **Visualize** modules

# Acknowledgements

▶ **Gunnar Rätsch** and AG Rätsch/MLB at FML

▶ **Wolfgang Busch, Jan Lohmann** and the JLab

▶ **Oliver Kohlbacher**, Div. SBS at University of Tübingen

▶ Detlef Weigel and Dept. VI at MPI for Dev. Biology

# Thank you for your attention!