# 10-601
# **Machine Learning**

## HMM applications in computational biology

# Central dogma

DNA

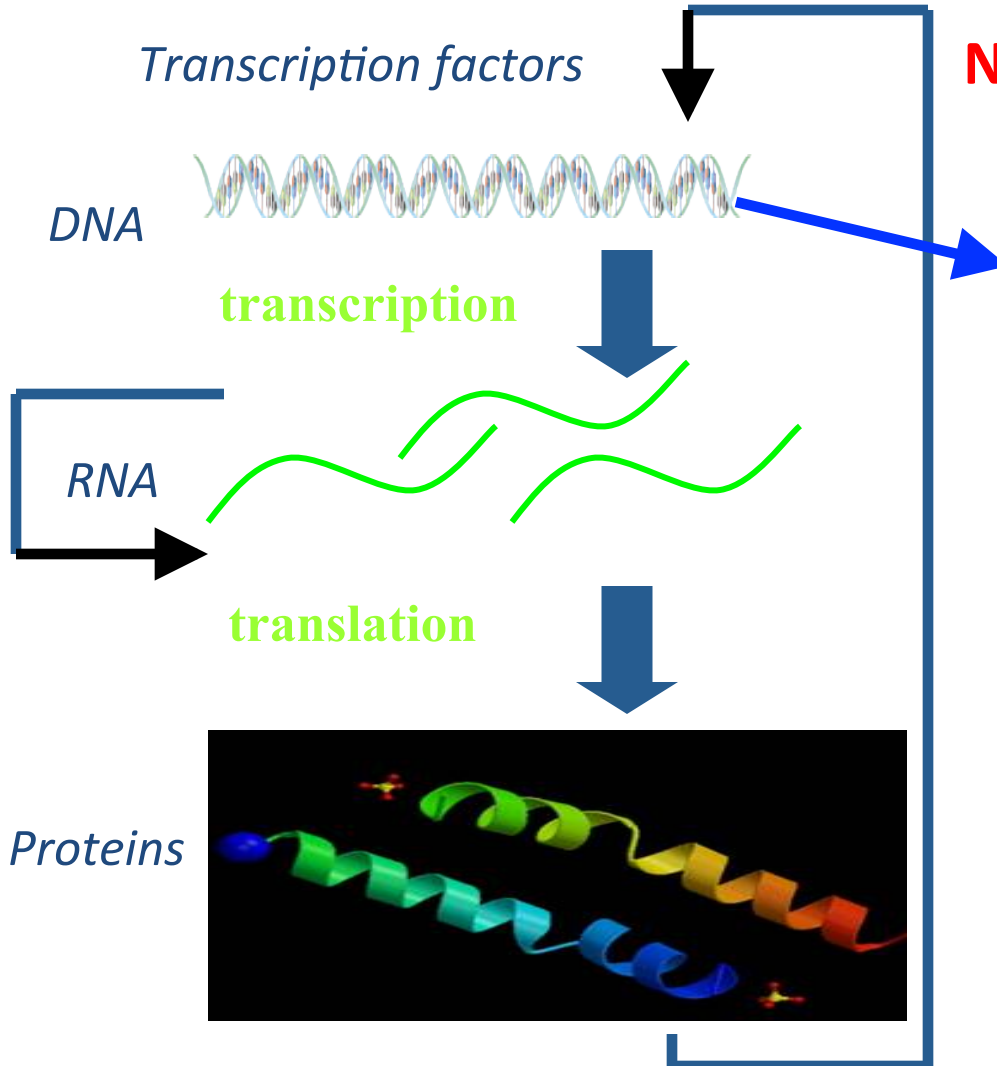transcription

mRNA

translation

Protein

CCTGAGCCAACTATTGATGAA

CCUGAGCCAACUAUUGAUGAA

PEPTIDE

# Biological data is rapidly accumulating

*Transcription factors*

*DNA*

**transcription**

*RNA*

**translation**

*Proteins*

**Next generation sequencing**



Growth of GenBank
(1982 - 2008)



Base Pairs
Sequences

# Biological data is rapidly accumulating

*Transcription factors*

*DNA*

**transcription**

*RNA*

**translation**

*Proteins*

**Array / sequencing technology**

# Biological data is rapidly accumulating

*Transcription factors*

*DNA*

**transcription**

*RNA*

**translation**

*Proteins*

**Protein interactions**

- 38,000 identified interactions
- Hundreds of thousands of predictions

# The New York Times

# Health

**Search Health** 3,000+ Topics

[                                      ] [Go]

**Inside Health**

**Research** | **Fitness & Nutritio**

# Company Unveils DNA Sequencing Device Meant to Be Portable, Disposable and Cheap

By ANDREW POLLACK
Published: February 17, 2012

DNA sequencing is becoming both faster and cheaper. Now, it is also becoming tinier.

A British company said on Friday that by the end of the year it would begin selling a disposable gene sequencing device that is the size of a USB memory stick and plugs into a laptop computer to deliver its
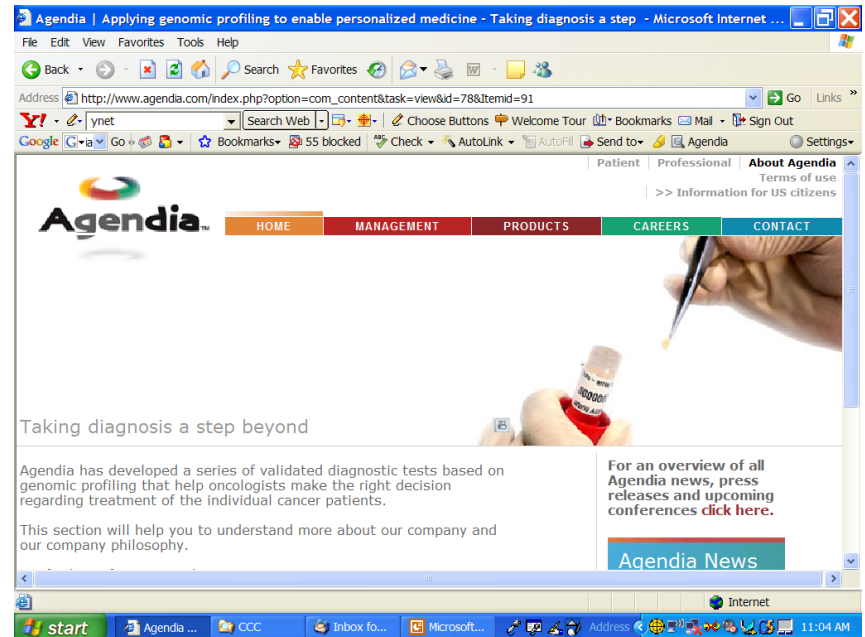
# FDA Approves Gene-Based Breast Cancer Test*

" MammaPrint is a DNA microarray-based test that measures the activity of 70 genes in a sample of a woman's breast-cancer tumor and then uses a specific *formula* to determine whether the patient is deemed low risk or high risk for the spread of the cancer to another site."



*Washington Post, 2/06/2007

...חדשות, ידיעות מהארץ והעול | N | עדי מקלט - 1156 קשר דף - ... | Error 0523: "The secon... | Latest Headlines | SquirrelMail | (Untitled) | (Untitled) | Pittsburgh, PA to Tole...

Positions: Fidelity Investments   ×   PLoS Computational Biology: Met... ×

Understand change with new tools for **epigenetics** research from **NEW ENGLAND BioLabs** Inc.

**New Software Section**
PLoS Computational Biology accepting presubmission inquiries

Login | Create Account | Feedback

**PLoS COMPUTATIONAL BIOLOGY**

a peer-reviewed open-access journal published by the Public Library of Science

Search articles...   GO   Advanced Search

Browse | RSS

**Home   Browse Articles   About   For Readers   For Authors and Reviewers**          **Journals    Hubs    PLoS.org**

**RESEARCH ARTICLE**                                    OPEN ACCESS

Download: PDF | Citation | XML

Print article

EzReprint New & improved!

# Metabolic Factors Limiting Performance in Marathon Runners

Article        Metrics        Related Content        Comments: 3

**Benjamin I. Rapoport**[1,2*]

1 M.D.– Ph.D. Program, Harvard Medical School, Boston, Massachusetts, United States of America, 2 Department of Electrical Engineering and Computer Science and Division of Health Sciences and Technology, Massachusetts Institute of Technology, Cambridge, Massachusetts, United States of America

To **add a note**, highlight some text. Hide notes

Make a general comment

**Jump to**
Abstract
Author Summary
Introduction
Results
Discussion
Methods
Acknowledgments

**Published in the** October 2010 Issue of PLoS Computational Biology

**Metrics** (i)

**Total Article Views: 74221**

**Average Rating**  (1 User Rating)
★ ★ ★ ★ ★  See all categories
Rate This Article
More

**Related Content**

**Related Articles on the Web**

## Abstract  Top

Each year in the past three decades has seen hundreds of thousands of runners register to run a major marathon. Of those who attempt to race over the marathon distance of 26 miles and 385 yards (42.195 kilometers), more than two-fifths experience

8

# Active Learning

nature
International weekly journal of science

Journal home > Archive > Letters to Nature > Abstract

## Journal content

+ **Journal home**
+ **Advance online publication**
+ **Current issue**
+ **Nature News**
+ **Archive**
+ **Supplements**
+ **Web focuses**
+ **Podcasts**
+ **Videos**

## Letters to Nature

## Functional genomic hypothesis generation and experimentation by a robot scientist

Ross D. King[1], Kenneth E. Whelan[1], Ffion M. Jones[1], Philip G. K. Reiser[1], Christopher H. Bryant[2], Stephen H. Muggleton[3], Douglas B. Kell[4] & Stephen G. Oliver[5]

1. Department of Computer Science, University of Wales, Aberystwyth SY23 3DB, UK
2. School of Computing, The Robert Gordon University, Aberdeen AB10 1FR, UK
3. Department of Computing, Imperial College, London SW7 2AZ, UK
4. Department of Chemistry, UMIST, P.O. Box 88, Manchester M60 1QD, UK

# Sequencing DNA



First human genome draft in 2001

Due to *accumulated errors*, we could only reliably read at most **100-200 nucleotides.**

# DARPA Shredder Challenge

# DARPA Shredder Challenge



DARPA Shredder Challenge
Puzzle 1
Image 1 of 1

# Shotgun Sequencing



Cloned genomes

Multiple genomes are sheared into variable sized segments

Unordered sequenced segments

Computational automated assembly

Resulting overlapping sequence segments. (The higher the coverage the better the quality of the sequencing.

ATGTTCCGATTAGGAAACCTATCTGTAACTGTTTCATTCAGTAAAAGGAGGAAATATAA

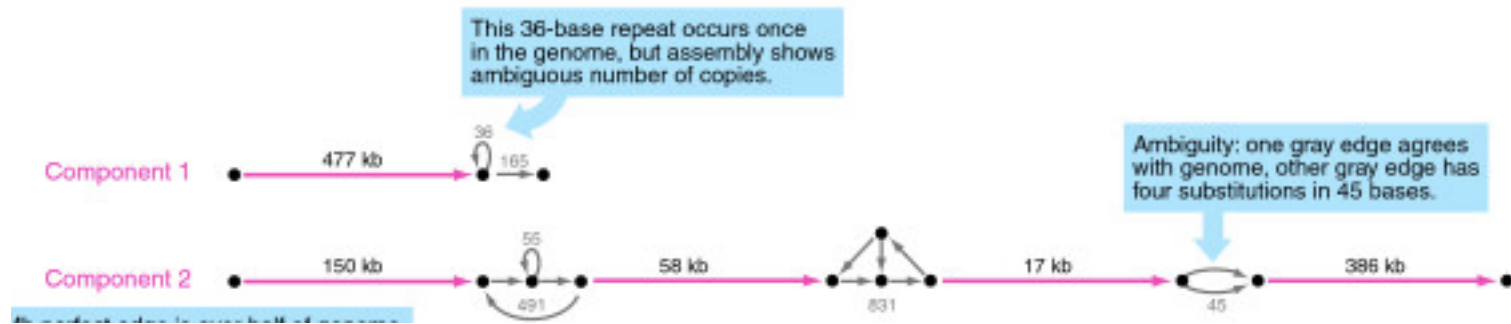Overlapping sequence segments combined to construct the genome consensus.

Wikipedia

# Caveats

- Errors in reading
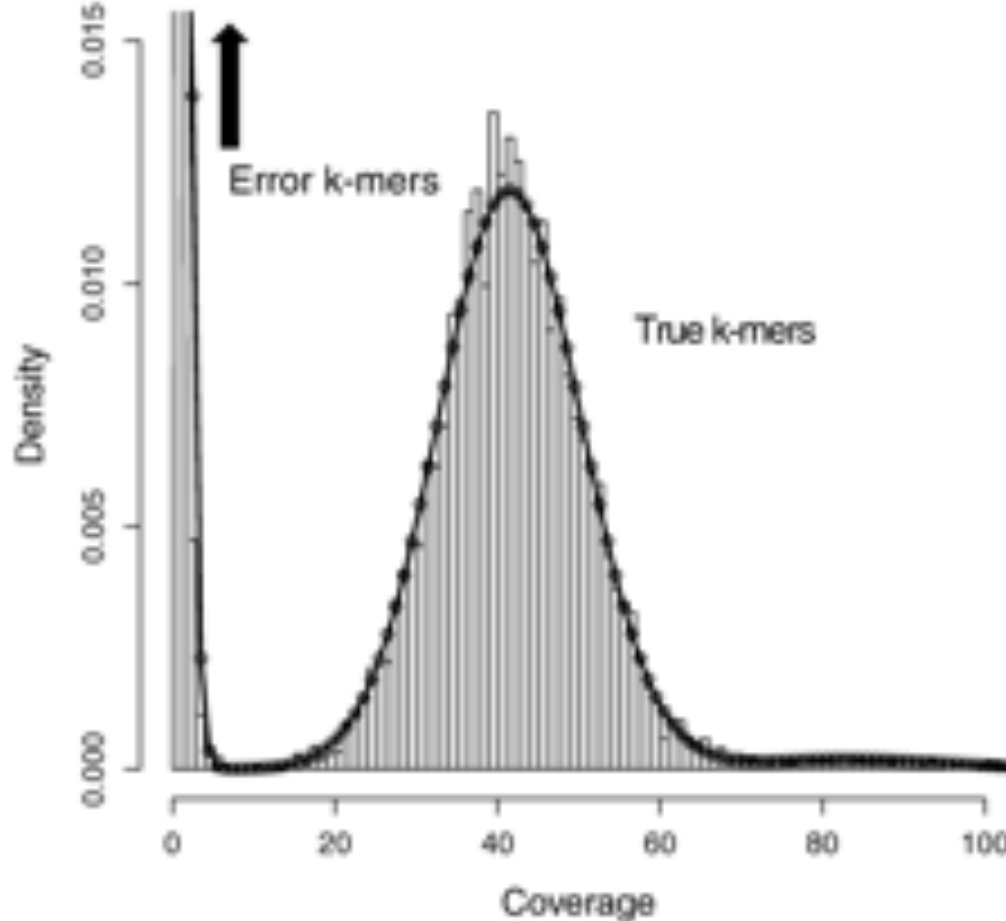- Non-trivial assembly task: repeats in the genome
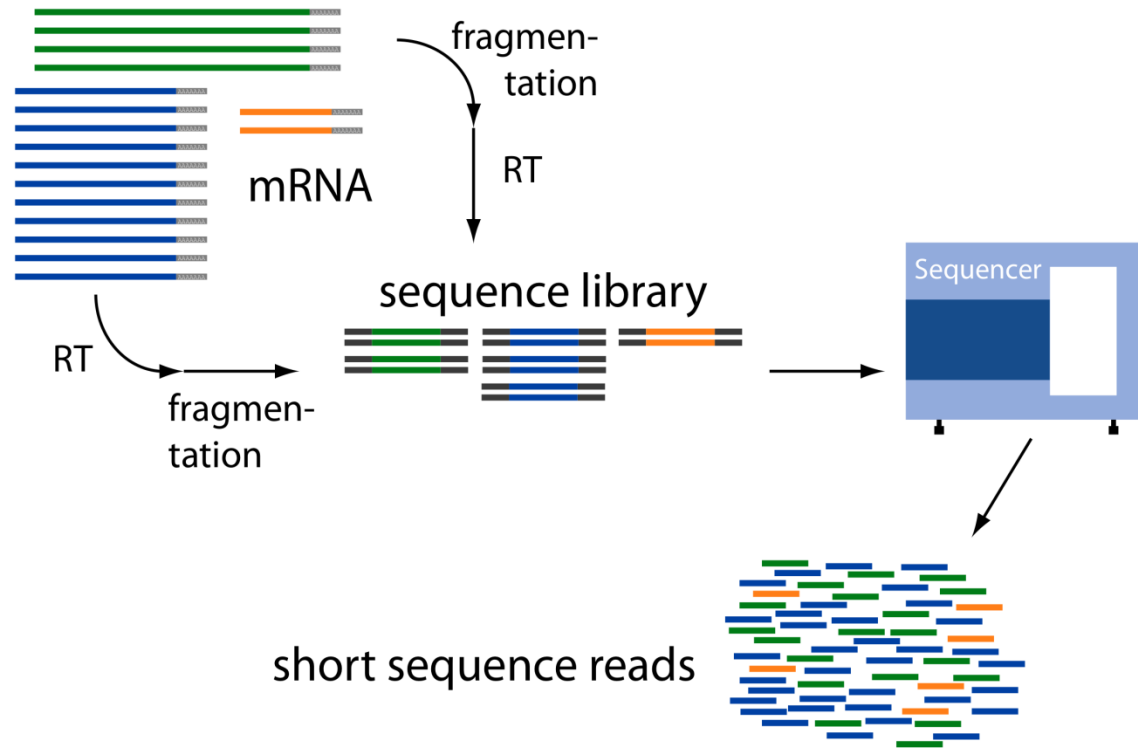


MacCallum et al., GB 2009

# Error Correction in DNA sequencing

• The fragmentation happens at random locations of the molecules. We expect all positions in the genome to have the same # number of reads

K-mers = substrings of length K of the reads. Errors create error k-mers.



Kellly et al., GB 2010

# Transcriptome Shotgun Sequencing (RNA-Seq)



@Friedrich Miescher Laboratory

Sequencing RNA transcripts.

*Reminder:*

- (mRNA) Transcripts are "expression products" of genes.
- Different genes having different expression levels so some transcripts are more or less abundant than others.
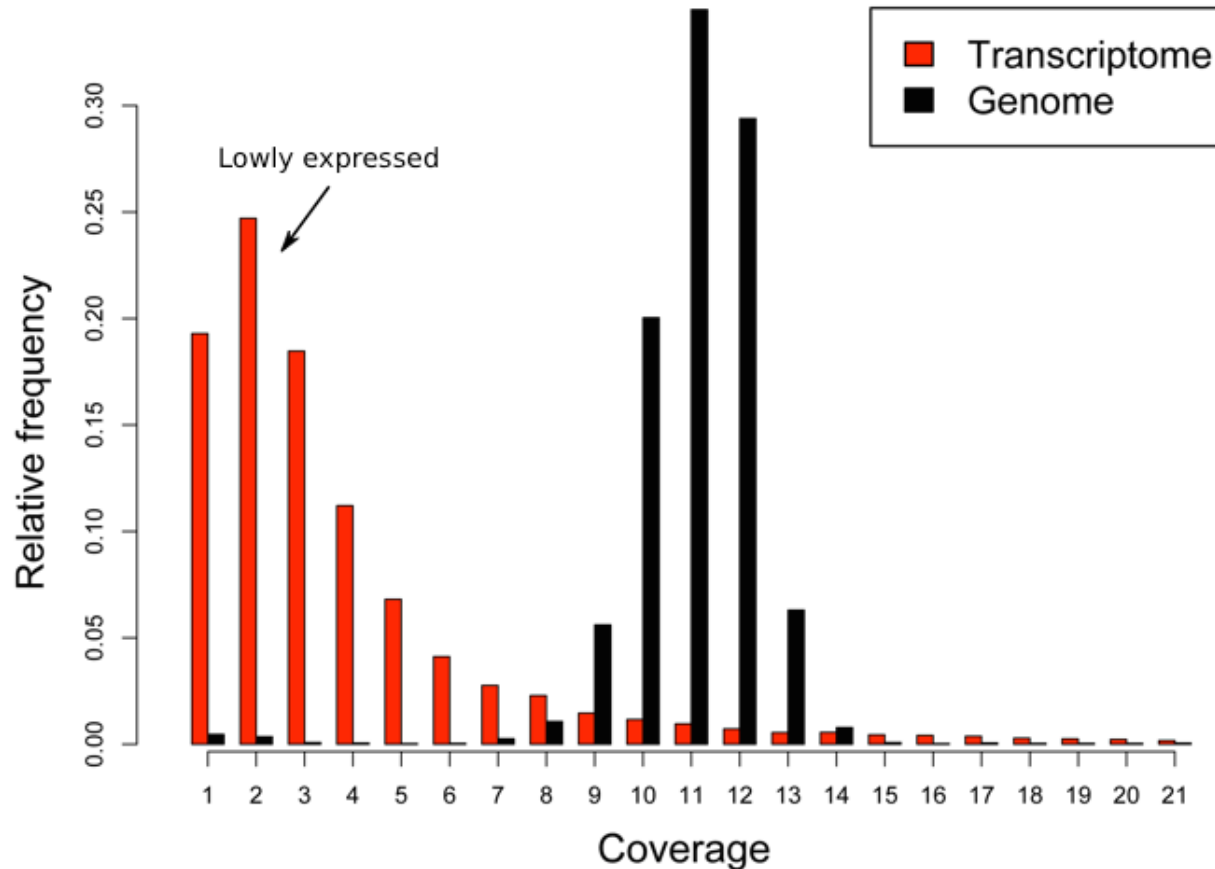
# Challenges

- Large datasets: 10-100 millions reads of 75-150 bps.
- Memory efficiency: Too time consuming to perform out-memory processing of data.

DNA Sequencing + **others :** alternative slicing, RNA editing, post-transcription modification.
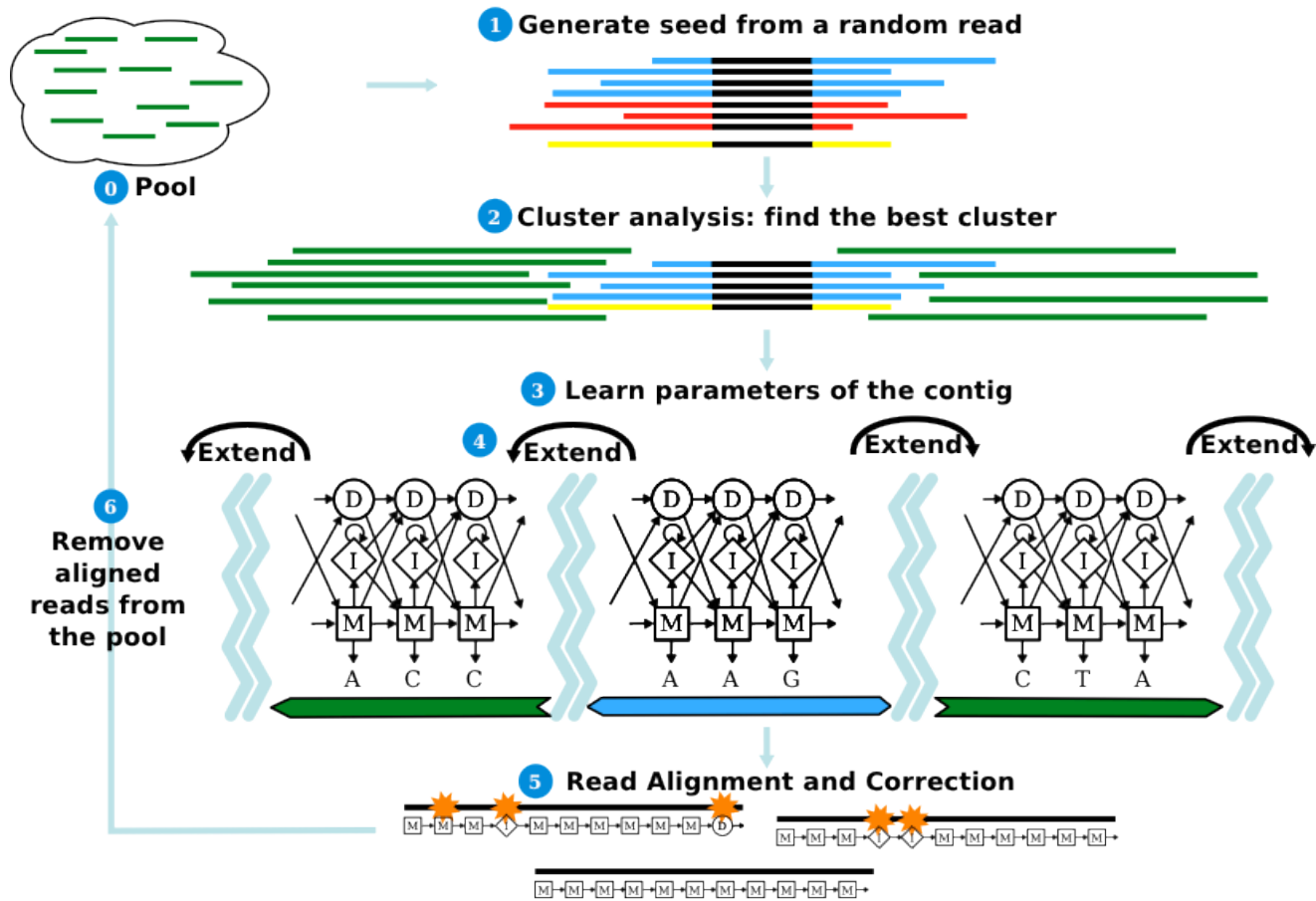
# Errors are non uniformly distributed

- Some transcripts are more prone to errors
- Errors are harder to correct in reads from lowly expressed transcripts

# SEECER
# Error Correction + Consensus sequence estimation for RNA-Seq data

# Key idea: HMM model

| Column | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Consensus | G | T | C | A | G | A | A | - | G | T | G | A | G | C | G | T | G | G | C | A | T | T | A | A | C | C | C | T | T | G | A | T | A | C | C | A | C | C | G | G | T | T | C | A | A | C | C |

**Read 1** G T C A G A A - G T G A G C G T G G C A T T A A C C C T T G A T A
40 40 40 40 40 40 **40** 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40

**Read 2** C A G A A [A] G T G A G C G T G G C A T T A A C C C T T G A T A C C
33 36 40 40 40 [34] 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40

**Read 3** C A G A A - G T G A G C G T G G C A T T A A C C C T T G A T A C C
40 40 40 40 40 **40** 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40

**Read 4** G T G A G C G T G G C A T T A A C C [-] T T G A T A C C A C C G G
40 40 40 40 40 40 18 40 40 40 40 40 40 33 40 30 40 22 [31] 40 40 40 40 40 40 40 40 40 40 40 40 40

**Read 5** T G G C A T T A A C C C T T G A T A C C A C C G G T T C A A
40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 37 40 40 40 40 40 40 40 40 19 40 40 40 40

**Read 6** G G C A T T A [C] C C C T T G A T A C C A C C G G T T C A A C C
22 40 40 17 40 40 40 [33] 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 27 40 40 40 40 40 40 40

Salmela et al., Bioinformatics 2011

The way sequencers work:
- Read letter by letter sequentially
- Possible errors: Insertion , Deletion or Misread of a nucleotide
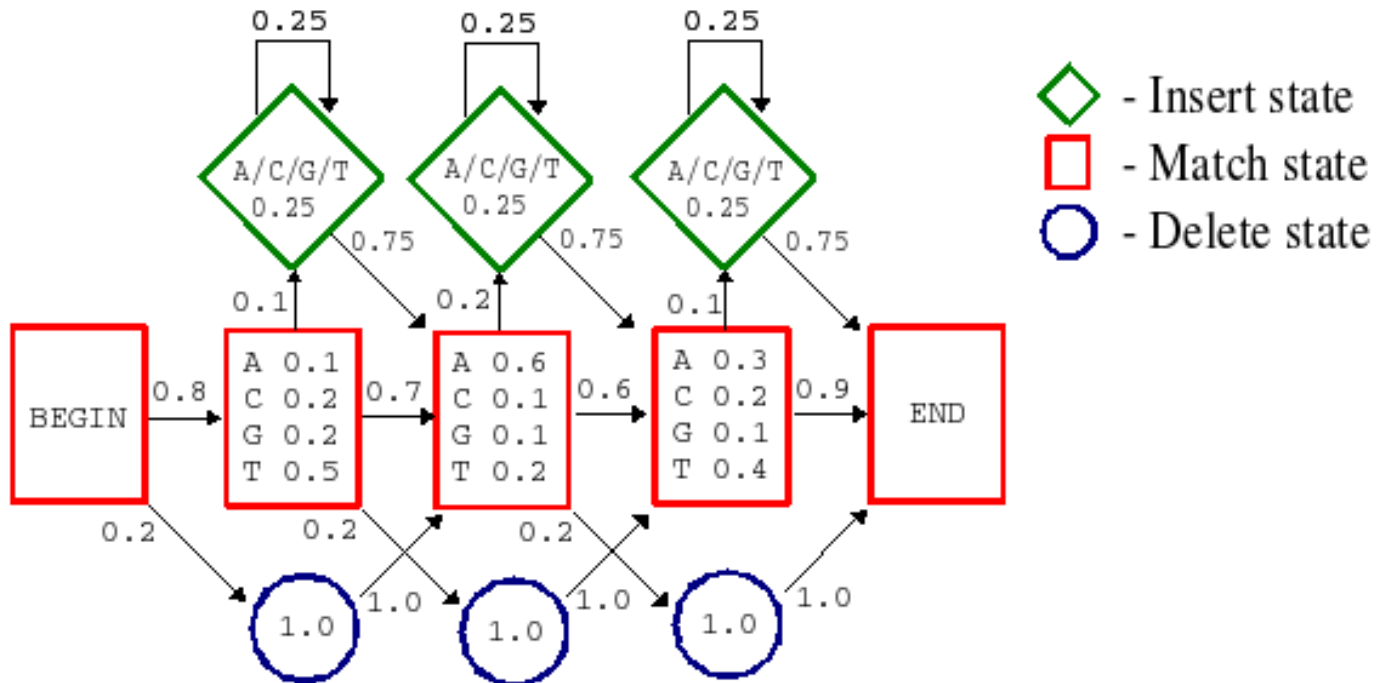
```
Column    1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
Consensus G  T  C  A  G  A  A  -  G  T  G  A  G  C  G  T  G  G  C  A  T  T  A  A  C  C  C  T  T  G  A  T  A  C  C  A  C  C  G  G  T  T  C  A  A  C  C

Read 1    G  T  C  A  G  A  A  -  G  T  G  A  G  C  G  T  G  G  C  A  T  T  A  A  C  C  C  T  T  G  A  T  A
          40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40

Read 2          C  A  G  A  A [A] G  T  G  A  G  C  G  T  G  G  C  A  T  T  A  A  C  C  C  T  T  G  A  T  A  C  C
                33 36 40 40 40[34]40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40

Read 3          C  A  G  A  A  -  G  T  G  A  G  C  G  T  G  G  C  A  T  T  A  A  C  C  C  T  T  G  A  T  A  C  C
                40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40

Read 4                            G  T  G  A  G  C  G  T  G  G  C  A  T  T  A  A  C  C [-] T  T  G  A  T  A  C  C  A  C  C  G  G
                                  40 40 40 40 40 18 40 40 40 40 40 40 33 40 30 40 22[31]40 40 40 40 40 40 40 40 40 40 40 40

Read 5                                              T  G  G  C  A  T  T  A  A  C  C  C  T  T  G  A  T  A  C  C  A  C  C  G  G  T  T  C  A  A
                                                    40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 37 40 40 40 40 40 40 40 19 40 40 40 40

Read 6                                                 G  G  C  A  T  T  A [C] C  C  C  T  T  G  A  T  A  C  C  A  C  C  G  G  T  T  C  A  A  C  C
                                                       22 40 40 17 40 40 40[33]40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 40 27 40 40 40 40 40 40 40
```

HMM diagram:

Insert states (green diamonds), each: A/C/G/T 0.25, self-loop 0.25, outgoing 0.75.

BEGIN (Match) → 0.8 → Match state:
A 0.1
C 0.2
G 0.2
T 0.5

→ 0.7 → Match state:
A 0.6
C 0.1
G 0.1
T 0.2

→ 0.6 → Match state:
A 0.3
C 0.2
G 0.1
T 0.4

→ 0.9 → END

Insert transitions: 0.1, 0.2, 0.1 (into insert states); 0.75 out.

Delete states (blue circles): 1.0 each; BEGIN → 0.2; Match → Delete 0.2; Delete → 1.0.

Legend:
◇ – Insert state
□ – Match state
○ – Delete state

# Building (Learning) the HMMs and Making Corrections (Inference)

Learning = Expectation-Maximization
Inference = Viterbi algorithm

**Seeding**:

Guessing possible reads using k-mer overlaps.
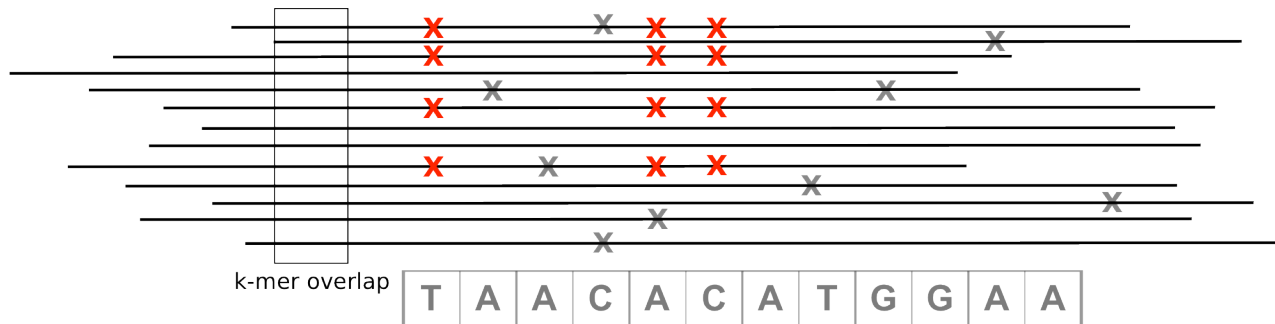
Constructing the HMM from these reads.

**Speed up:**

The k-mer overlaps yield approximate multiple alignments of reads.

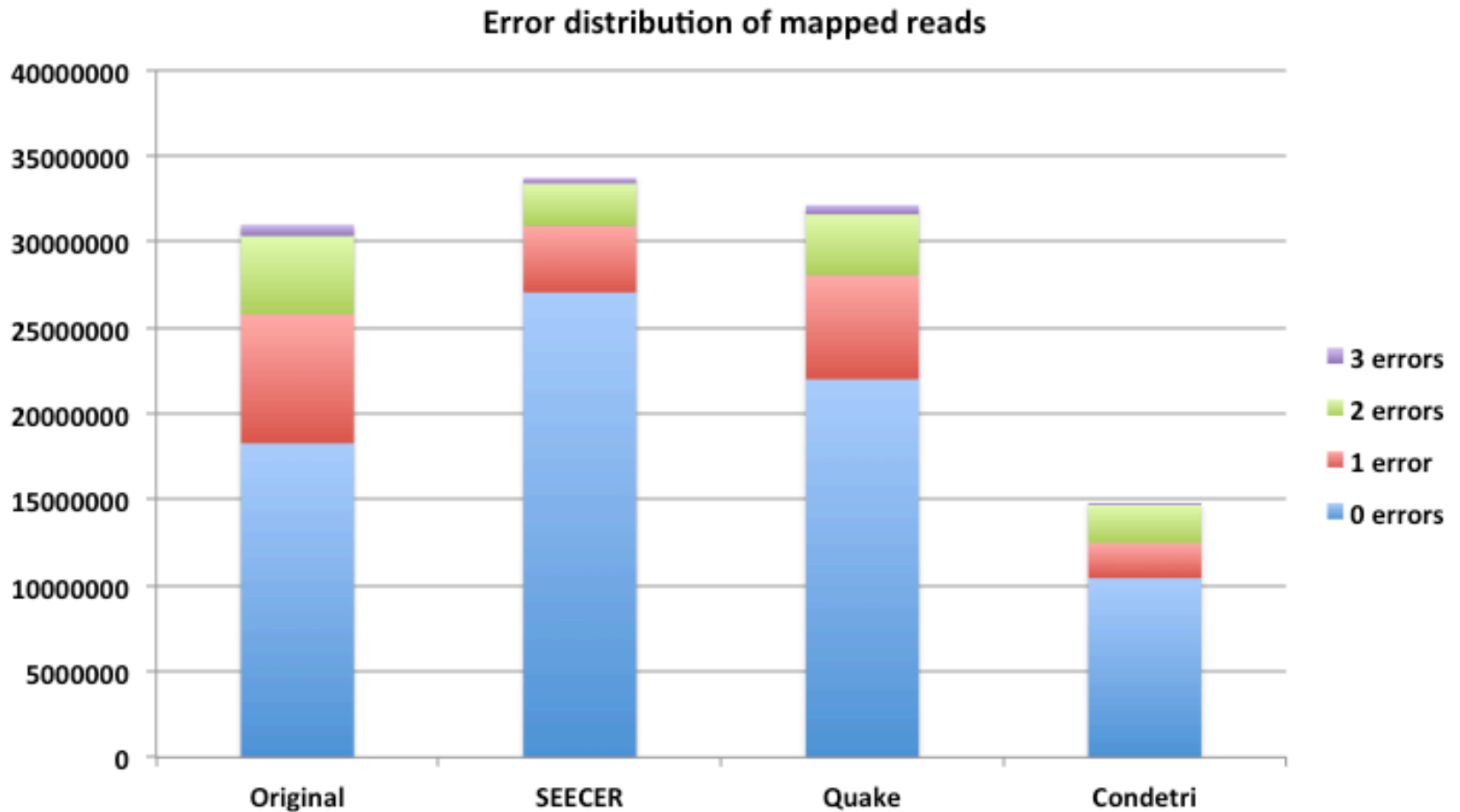We can learn HMM parameters from this directly.

# Clustering to improve seeding



k-mer overlap

| T | A | A | C | A | C | A | T | G | G | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|

Real biological differences should be supported by a set of reads with similar mismatches to the consensus

1. Clustering positions with mismatches to identify clusters of correlated positions.
2. Build a similarity matrix between these positions.
3. Use Spectral clustering to find clusters of correlated positions.
4. Filter reads have mismatches in these clusters.

# Comparison to other methods



Error distribution of mapped reads

# Using the corrected reads, the assembler can recover **more** transcripts



80% length reconstruction

# Things that work

- Approximate learning to speed up on large datasets.

- In real world, one technique is not enough. A solution involves using many techniques.

- Precision and Recall are trade-offs.

# Central dogma

Different regulators control the information flow from DNA to protein

**TF**

Transcription factors (TFs) bind to DNA and activate genes

**CCTGAGCCAACTATTGATGAA**

**DNA**

↓ transcription

**mRNA**

**CCUGAGCCAACUAUUGAUGAA**

**miR**

Micro RNAs (miRs) bind to mRNA to down regulate their expression

↓ translation

Protein

**PEPTIDE**

# Integrating expression and protein-DNA interaction data

Lee *et al Science* 2002

Bar-Joseph *et al Nature Biotechnology* 2003

# Methods for reconstructing networks in cells



Amit et al *Science* 2009

Venancio et al *Genome Biology* 2009

Gerstein et al *Science* 2010

# Key problem: Most high-throughput data is static

<u>Time-series measurements</u>

<u>Static data sources</u>

Sequencing

motif

CHIP-chip

microarray

PPI

Time

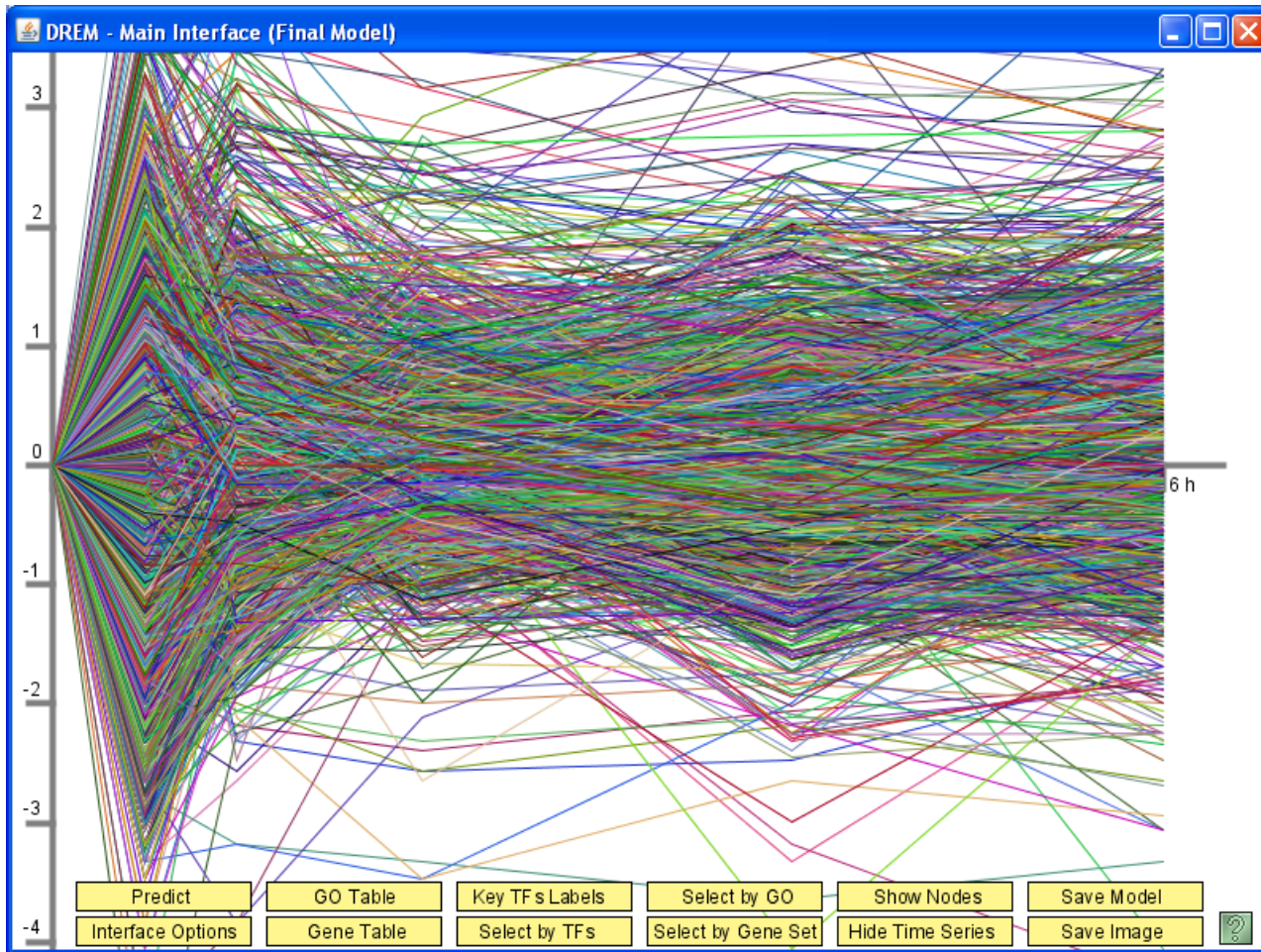# DREM: Dynamic Regulatory Events Miner

**a  Time Series Expression Data**

Expression Level

1

0

-1

0h    1h    2h    3h

time

**b  Static TF-DNA Binding Data**

TF A

TF B

TF C    TF D

**a** Time Series Expression Data

Expression Level

0h    1h    2h    3h    time

**b** Static TF-DNA Binding Data

TF A

TF B

TF C    TF D

**c** Model Structure

Expression Level

TF A
TF B

TF A

TF B

TF C
TF D

0    1h    2h    3h    time

**d** IOHMM Model

1

0.95

0.05

?

0.1

?

0.9

1

# Things are a bit more complicated: Real data

# A Hidden Markov Model

**Hidden States**

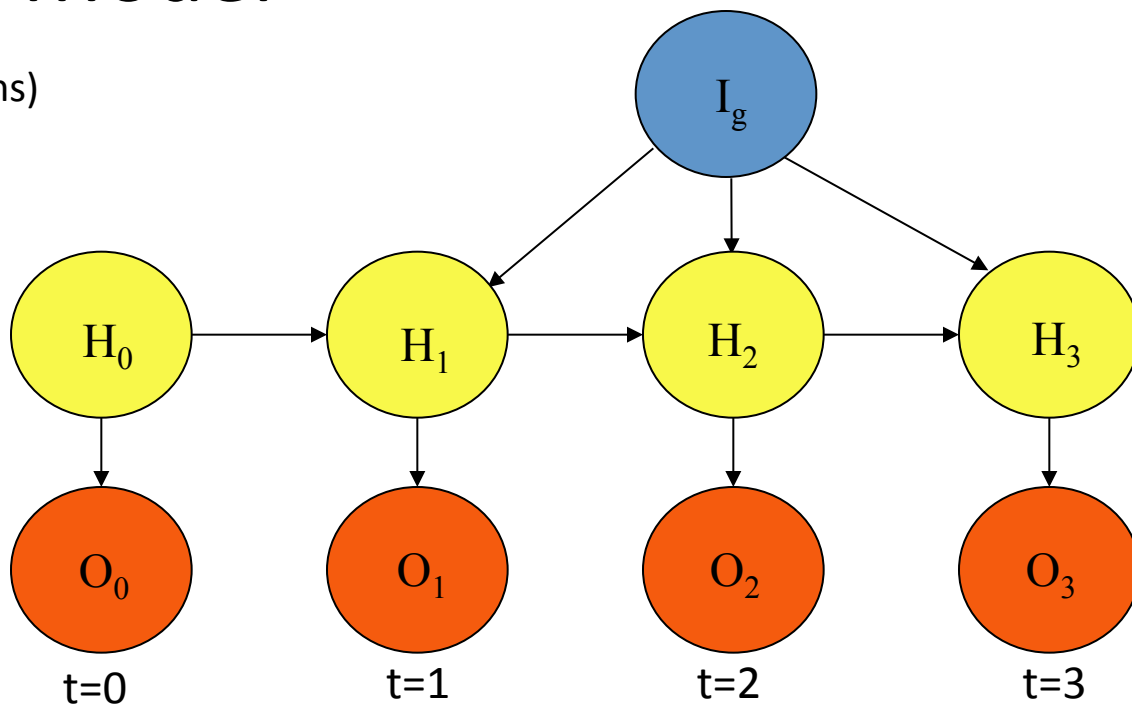**Observed outputs (expression levels)**



t=0    t=1    t=2    t=3

$$L(H,O;\Theta) = \prod_{i=1}^{n}\left[\prod_{t=1}^{T} p(O_t(i)\,|\,H_t(i))\right]\left[\prod_{t=2}^{T} p(H_t(i)\,|\,H_{t-1}(i))\right]$$

Schliep et al *Bioinformatics* 2003

# Input – Output Hidden Markov Model

Input (Static TF-gene interactions)

$I_g$

Hidden States (transitions between states form a tree structure)

$H_0 \quad H_1 \quad H_2 \quad H_3$

Emissions (Distribution of expression values)

$O_0 \quad O_1 \quad O_2 \quad O_3$

t=0     t=1     t=2     t=3

Log Likelihood       But how do we express these conditional probabilities?

$$r(G|M) = \sum_{g \in G} \log \sum_{q \in Q} \prod_{t=1}^{n-1} f_{q(t)}(o_g(t)) \prod_{t=1}^{n-1} P(H_t = q(t) | H_{t-1} = q(t-1), I_g)$$

Sum over all genes

Sum over all paths $Q$

Product over all Gaussian emission density values on path

Product over all transition probabilities on path
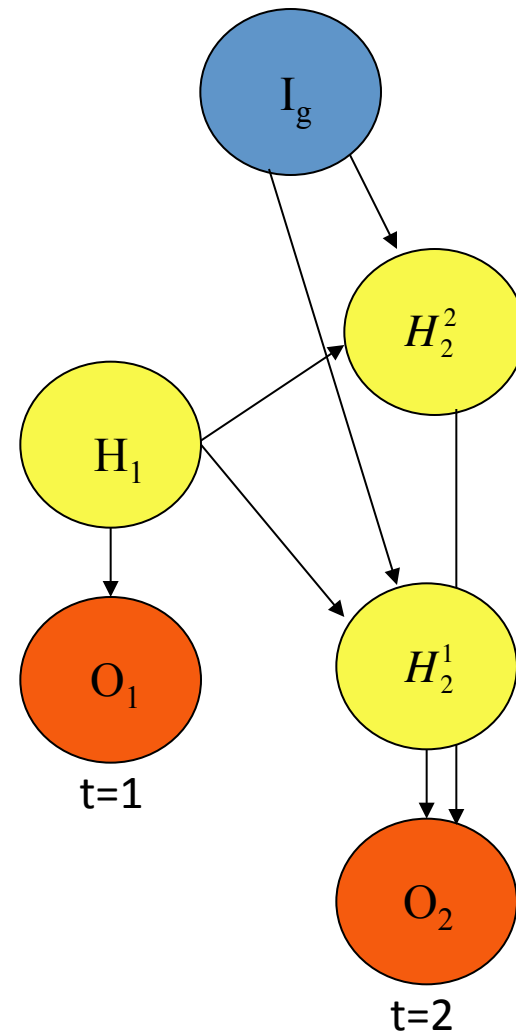
# Input-Output Hidden Markov Model

## learning the transition probabilities

How do compute $P$ for a state with 2 children?
We can write it as a logistic regression classification problem!

$$P(H_2^1 = q(2) \mid H_1 = q(1), I_g) = ?$$

$$P(H_2^2 = q(2) \mid H_1 = q(1), I_g) = ?$$

# Input-Output Hidden Markov Model

## learning the transition probabilities
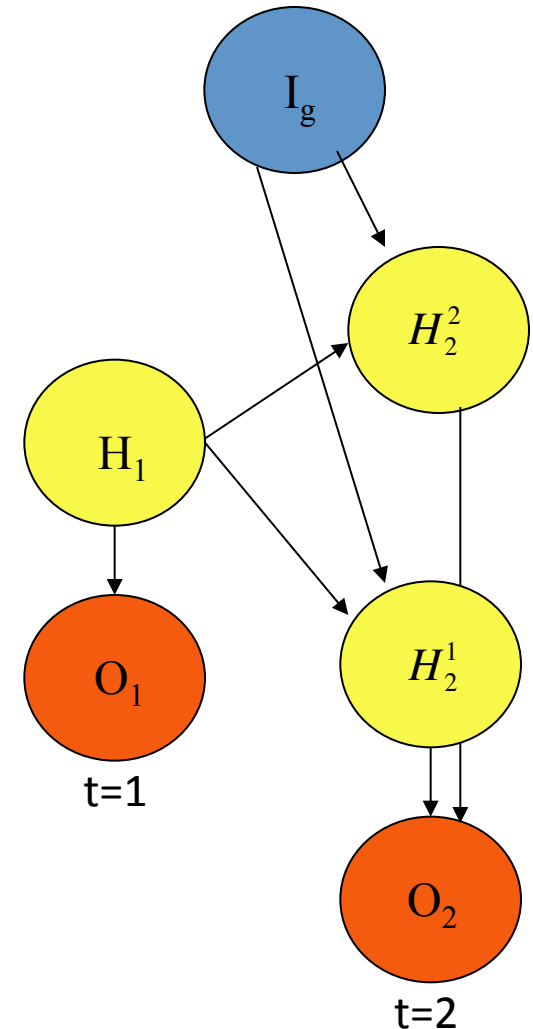
How do compute *P* for a state with 2 children?

We can write it as a logistic regression classification problem!

$$P(H_2^1 = q(2) \mid H_1 = q(1), I_g) = \frac{1}{1 + \exp(w_0 + \sum_r I_{g,r} w_r)}$$
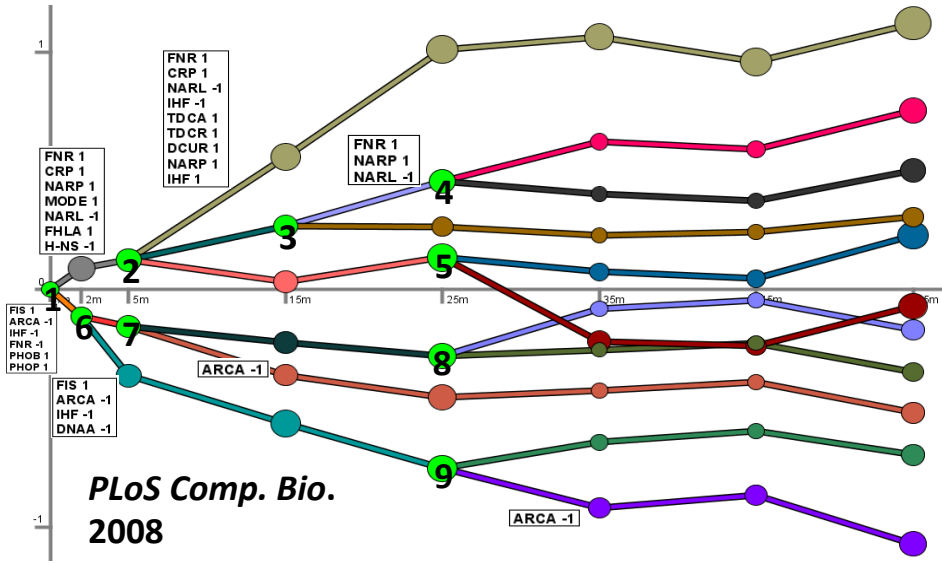
Sum over all regulators
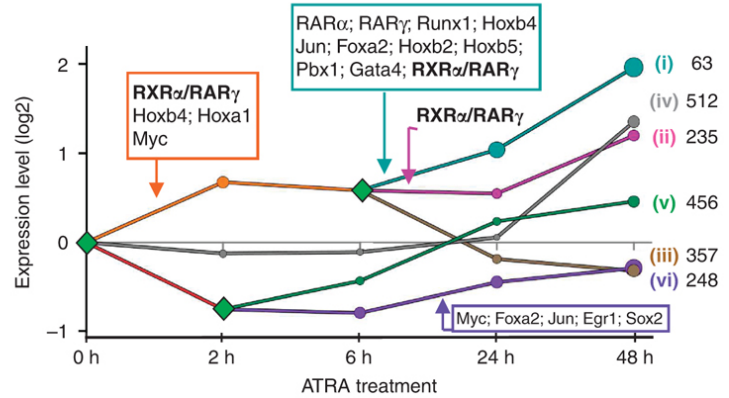
Optimize $w_i$'s with a logistic regression classifier

$$\hat{W} = \underset{W}{\operatorname{argmax}} \left[ l(W) + \ln p(W) \right]$$

likelihood with parameters *W*

$L_1$-penalty to promote sparsity



$I_g$

$H_2^2$

$H_1$

$O_1$

$H_2^1$

$O_2$

t=1

t=2

**E. coli. response**

FNR 1
CRP 1
NARL -1
IHF 1
TDCA 1
TDCR 1
DCUR 1
NARP 1
IHF 1

FNR 1
NARP 1
NARL -1

FNR 1
CRP 1
NARP 1
MODE 1
NARL -1
FHLA 1
H-NS -1

FIS 1
ARCA -1
IHF -1
FNR -1
PHOB 1
PHOP 1

FIS 1
ARCA -1
IHF -1
DNAA -1

ARCA -1

ARCA -1

*PLoS Comp. Bio.*
**2008**

**Stem cells differentiation**

RARα; RARγ; Runx1; Hoxb4
Jun; Foxa2; Hoxb2; Hoxb5;
Pbx1; Gata4; **RXRα/RARγ**

**RXRα/RARγ**
Hoxb4; Hoxa1
Myc

Myc; Foxa2; Jun; Egr1; Sox2

Expression level (log2)

(i) 63
(iv) 512
(ii) 235
(v) 456
(iii) 357
(vi) 248

0 h    2 h    6 h    24 h    48 h
ATRA treatment

*Nature MSB* **2011**

**Fly development**      *Science* **2010**

**Mouse Immune response**

IRF7

*Genome Research* **2010,**
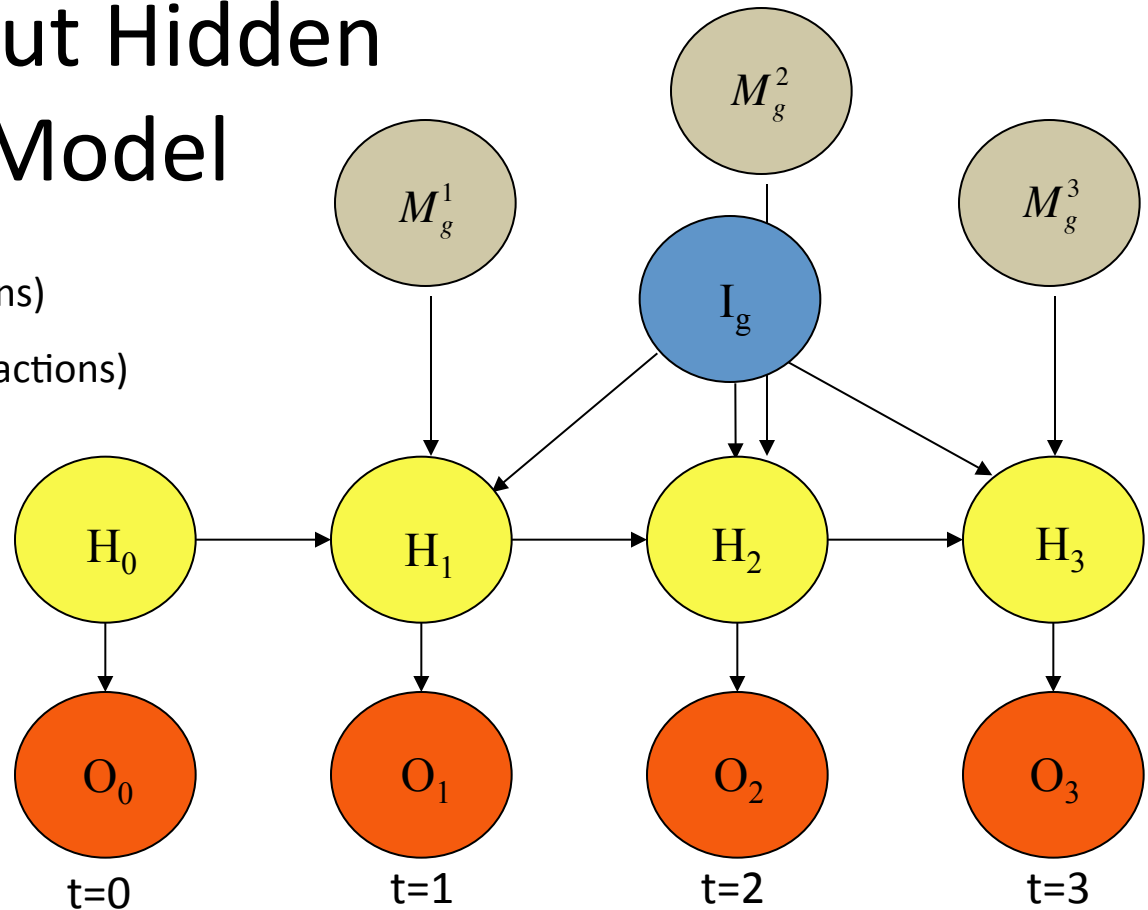**PLoS ONE 2011**

# mirDREM

# Input – Output Hidden Markov Model

Input (**Static** TF-gene interactions)

(**Dynamic** miR-gene interactions)

Hidden States (transitions between states form a tree structure)

Emissions (Distribution of expression values)

$M_g^1$   $M_g^2$   $M_g^3$

$I_g$

$H_0$   $H_1$   $H_2$   $H_3$

$O_0$   $O_1$   $O_2$   $O_3$

t=0   t=1   t=2   t=3

Log Likelihood

But how do we express these conditional probabilities?

$$r(G|M) = \sum_{g \in G} \log \sum_{q \in Q} \prod_{t=1}^{n-1} f_{q(t)}(o_g(t)) \prod_{t=1}^{n-1} P(H_t = q(t)|H_{t-1} = q(t-1), I_g, M_g^t)$$

Sum over all genes

Sum over all paths $Q$

Product over all Gaussian emission density values on path

Product over all transition probabilities on path

# Input-Output Hidden Markov Model

learning the transition probabilities

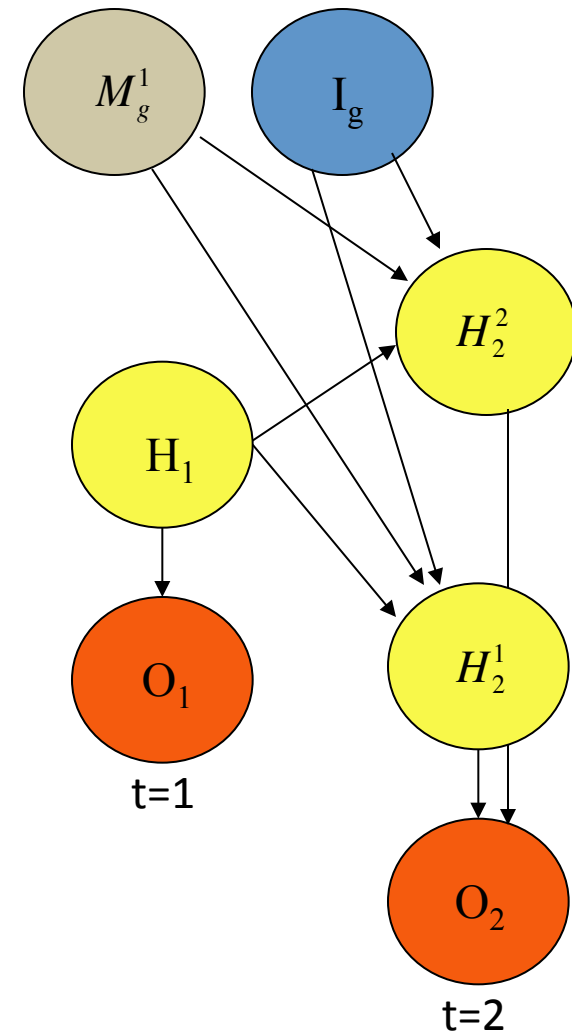How do compute *P* for a state with 2 children?

We can write it as a logistic regression classification problem!

$$P(H_2 = q(2) \mid H_1 = q(1), I_g) = \frac{1}{1 + \exp(w_0 + \underbrace{\sum_r I_{g,r} w_r}_{} + \underbrace{\sum_m M^1_{g,r} w_m}_{})}$$

Sum over all TFs     Sum over all miRs

Optimize $w_i$'s with a **constrained** logistic regression classifier

$$\hat{W} = \underset{W}{\mathrm{argmax}} \left[ \underbrace{l(W)}_{} + \underbrace{\ln p(W)}_{} \right]$$

likelihood with parameters *W*     L$_1$-penalty to promote sparsity

s.t. $w_i \leq 0, i \in miRNAs$



$M^1_g$    $I_g$

$H^2_2$
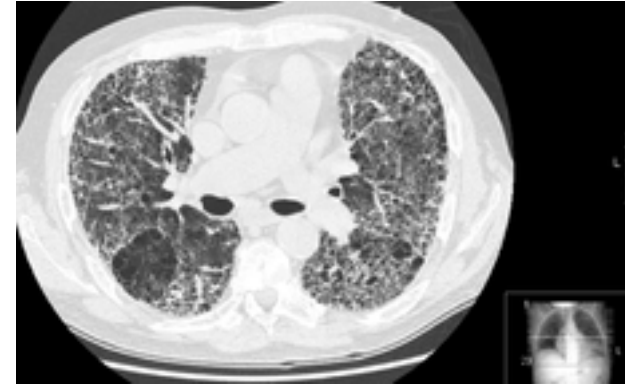
$H_1$

$O_1$

$H^1_2$

t=1

$O_2$

t=2

# Application to mouse data to understand human lung disease

Idiopathic pulmonary fibrosis (IPF)

-100,000 people are affected (USA)

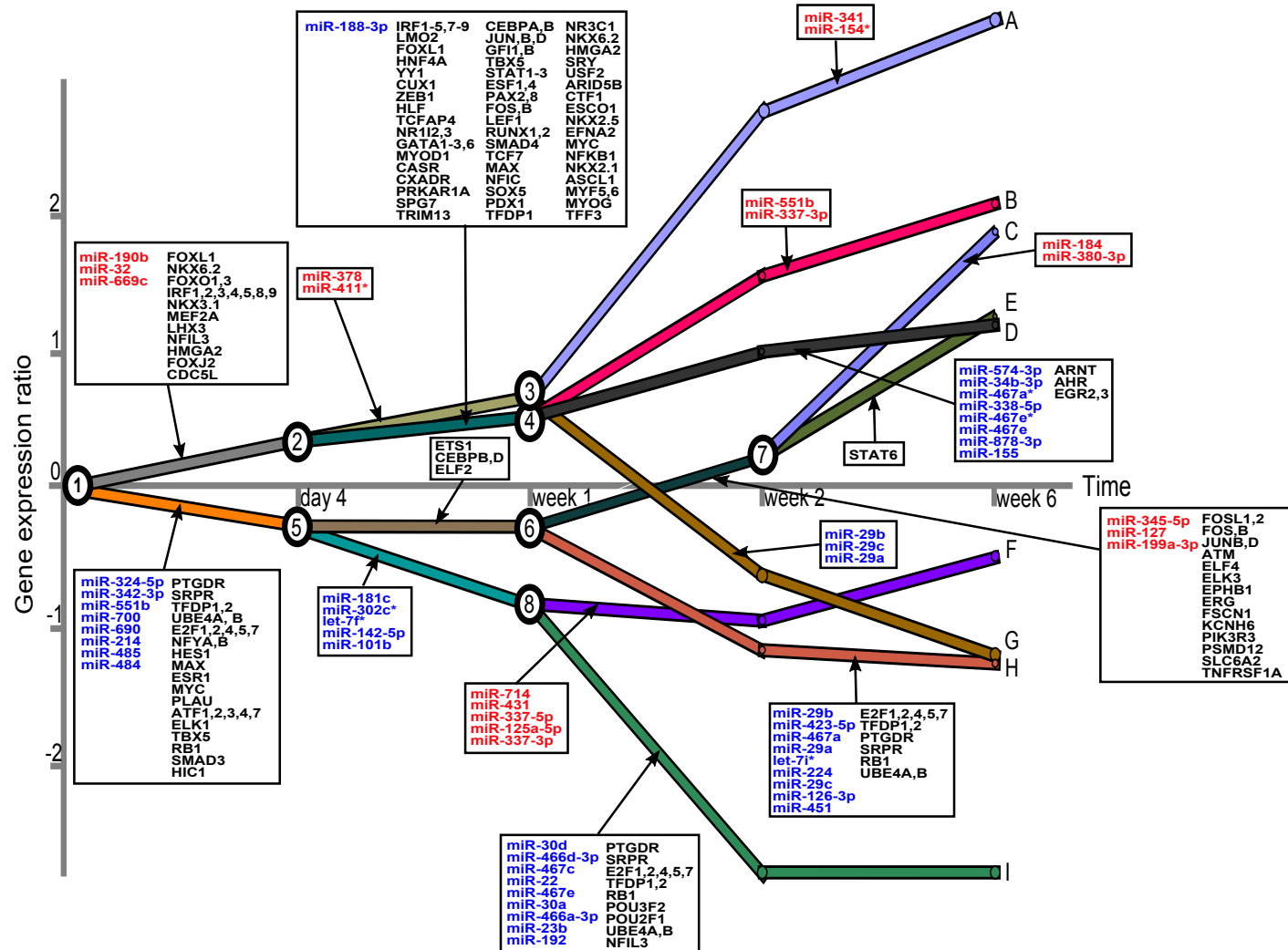-about 30,000 new cases each year

-50% death rate after 3 years



Extensive lung fibrosis (source:wikipedia)

-pathways for lung development appear activated in reverse direction during the disease

# Joint dynamic network for lung development



- 22 out 56 miRNAs predicted by the method are differentially expressed in patient cohorts with the IPF lung disease
- different to development miR-30d is down regulated in IPF patients