

Probabilistic walks with Graeme

Sean Eddy

Molecular & Cellular Biology, and Applied Mathematics

HHMI and Harvard University

eddylab.org



Cambridge
Celts FC

The Perse Upper School

A1134

A1134

A1134

A1134

A1134

Queen Edith's Way RED CROSS

Long Road Sixth
Form College

Cambridge Academy for
Science and Technology

Fawcett Primary School

Cancer Research UK
Cambridge Institute

Cambridge Blood
Donor Centre

Nightingale
Recreation
Ground

The Green Man

CPDC Cambridge
Professional...

MRC Laboratory of
Molecular Biology

The University of
Cambridge School of...

Addenbrooke's Hospital

Bell Cambridge

n Inn

Hudson's Ale House

Cambridge
Biomedical Campus

The Rosie Hospital

Cambridge Institute
of Public Health

TRUMPINGTON

The Bun Shop
(Trumpington Post...)

Hobson's Brook

Babraham Rd

Partners

Hauxton Rd
A1307

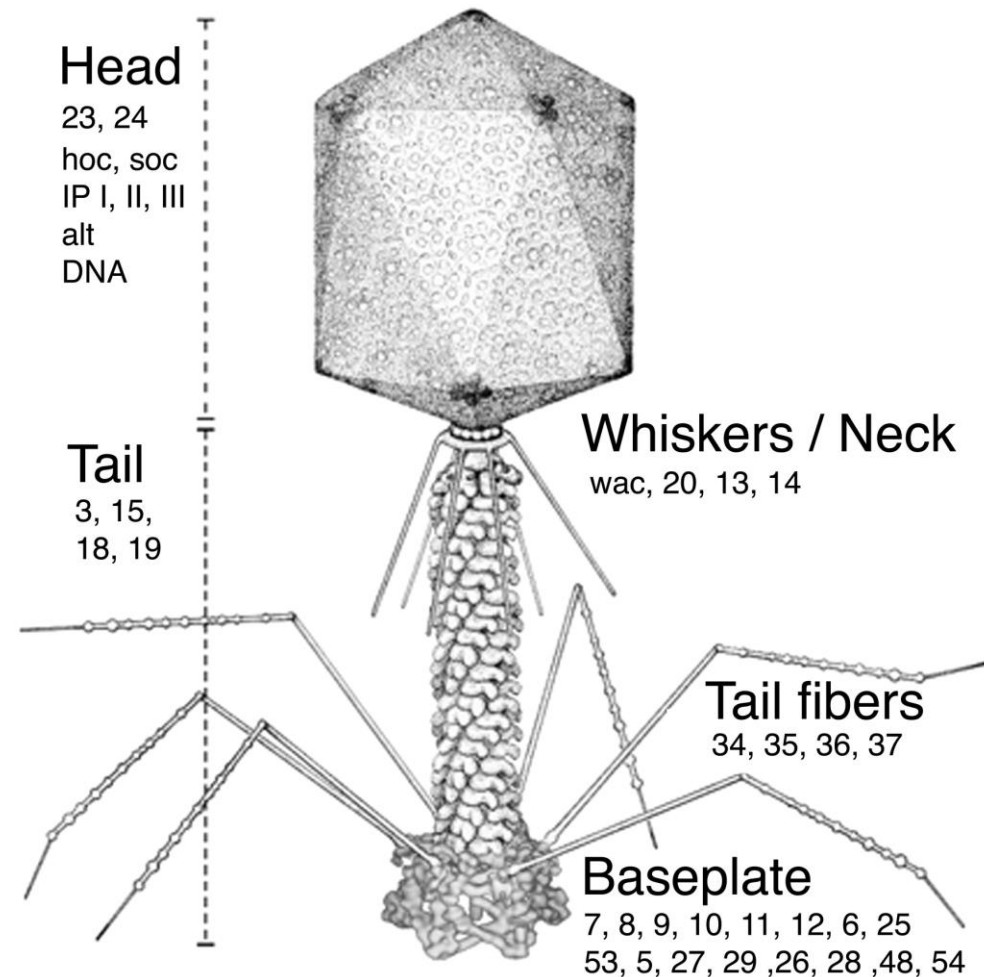
Abcam

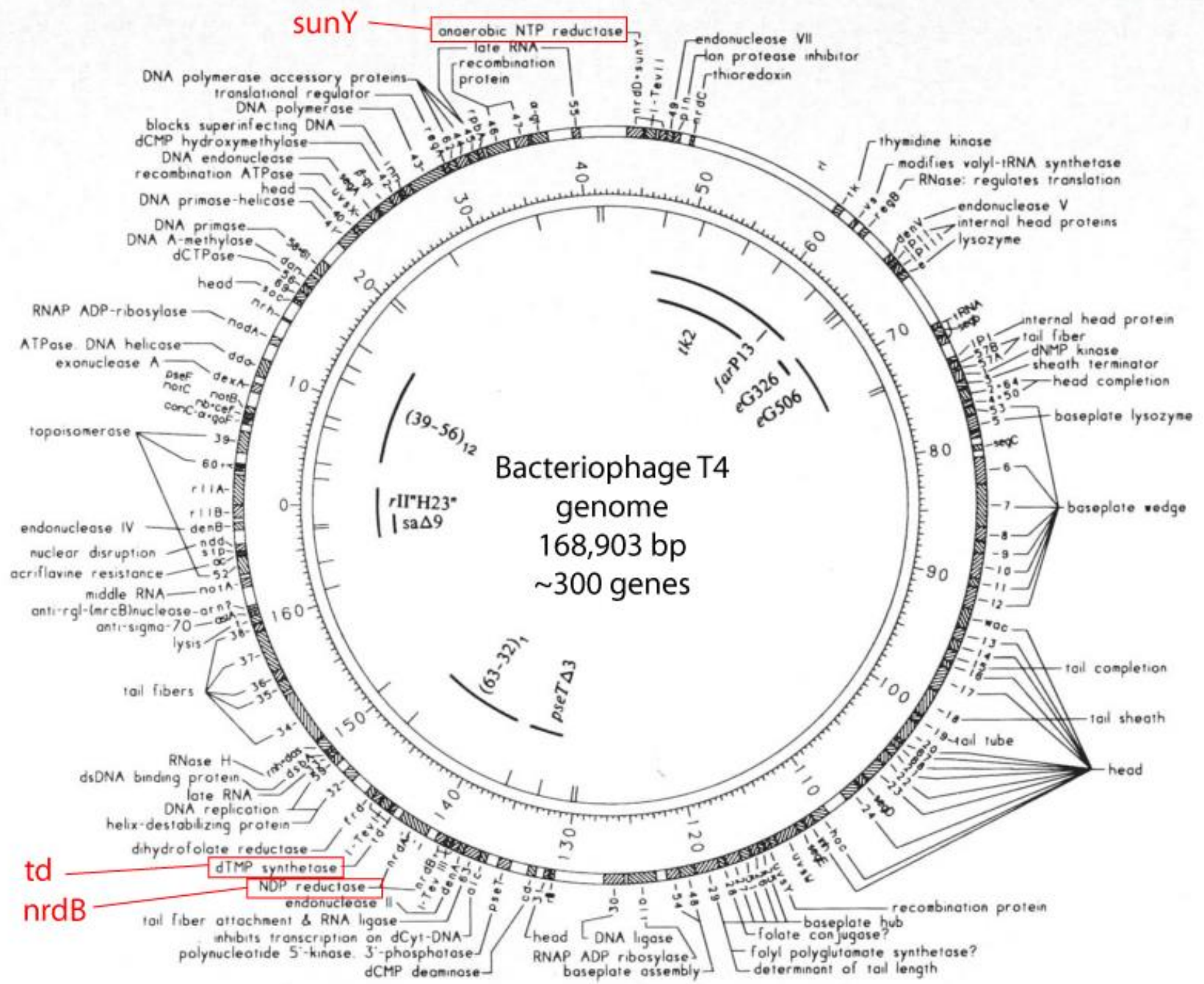
Multiple Self-Splicing Introns in Bacteriophage T4: Evidence from Autocatalytic GTP Labeling of RNA In Vitro

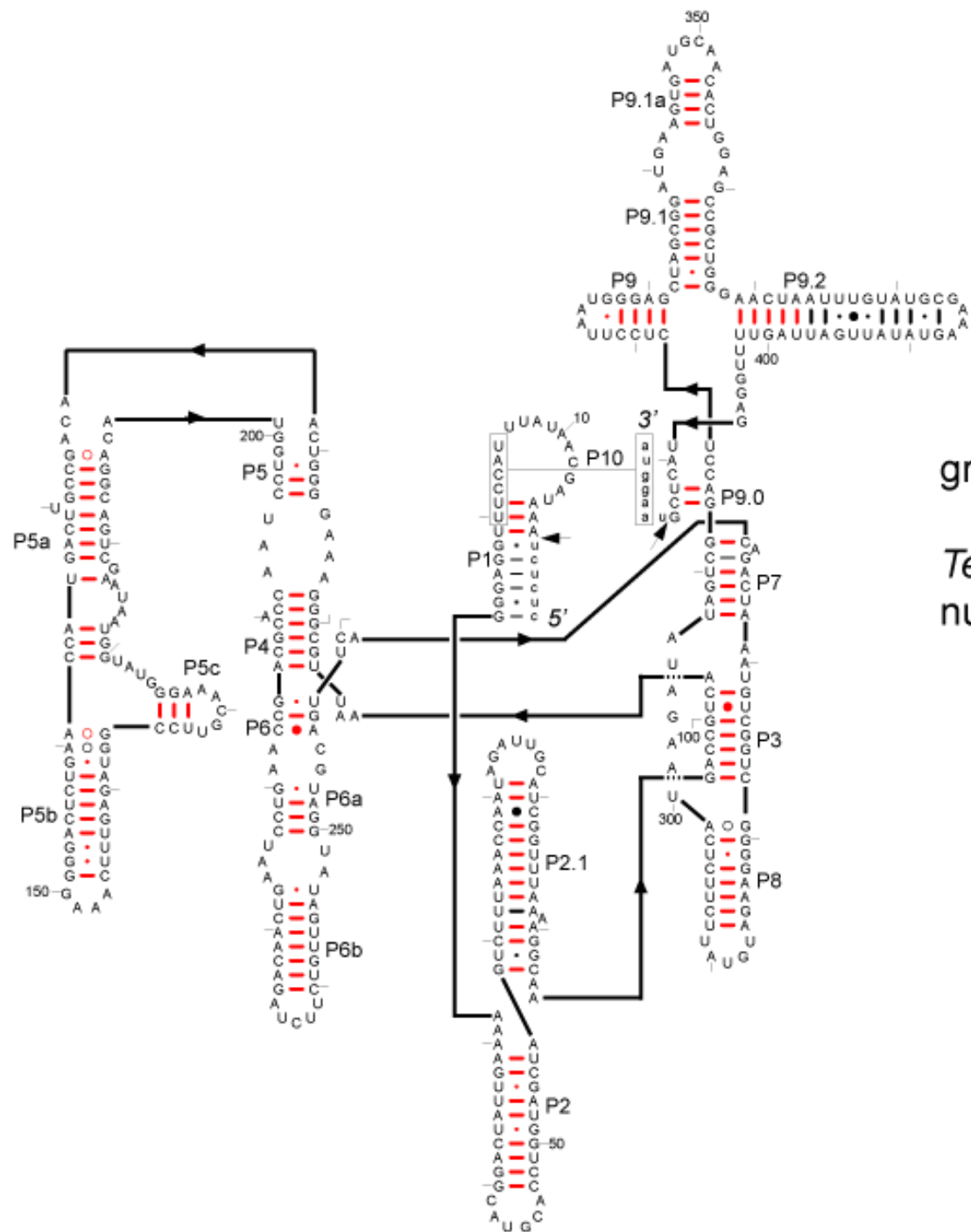
Jonatha M. Gott,^{*} David A. Shub,^{*}
and Marlene Belfort[†]

^{*}Department of Biological Sciences
State University of New York, Albany
Albany, New York 12222

[†]Wadsworth Center for Laboratories and Research
New York State Department of Health
Albany, New York 12201







group I self-splicing intron

Tetrahymena thermophila
nuclear LSU rRNA

source: Robin Gutell
Comparative RNA Website
<http://www.rna.icmb.utexas.edu/>

Hidden Markov Models in Computational Biology: Applications to Protein Modeling UCSC-CRL-93-32

Anders Krogh^{*†}, Michael Brown[†], I. Saira Mian[§],
Kimmen Sjölander[†], David Haussler[†]

[†] Computer and Information Sciences

[§] Sinsheimer Laboratories

University of California, Santa Cruz, CA 95064, USA.

email: krogh@nordig.ei.dth.dk, haussler@cse.ucsc.edu

August 17, 1993

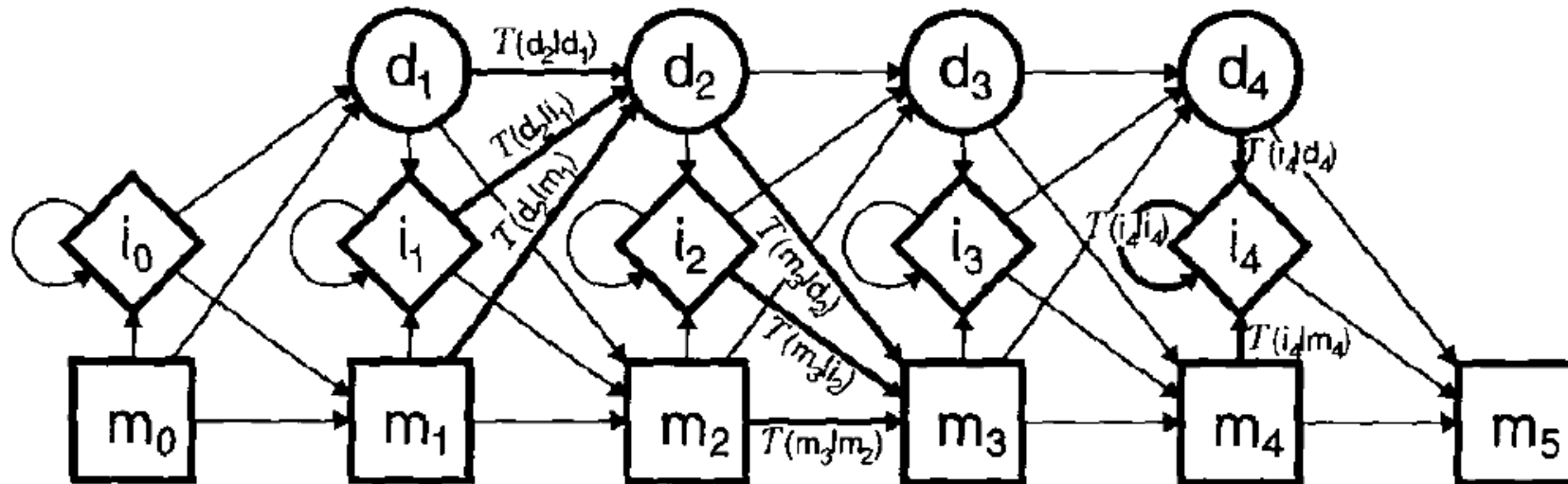
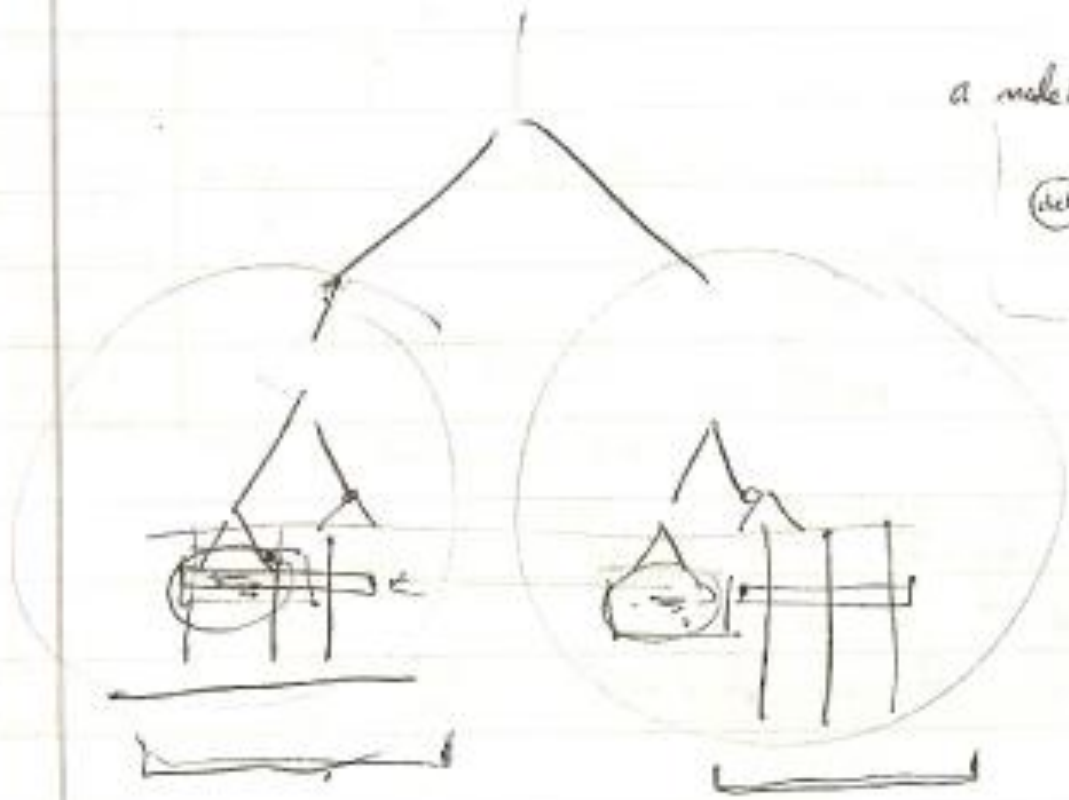


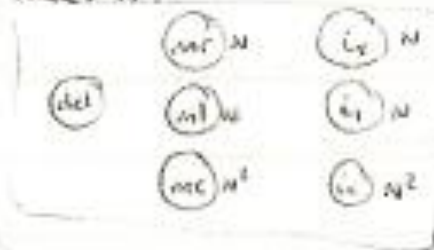
Figure 1. The model.

Attempt to repeat some analytic method that is considered unreliable and difficult until patience and hard work yield results similar to those published by the author. Pleasure derived from success, especially if it has come without the supervision of an instructor (that is, working alone), is a clear indication of aptitude for experimental work.

Santiago Ramon y Cajal (1916)



a model is:



100 P_g/state

7 substrates

49 transition P_g

4+8 = 48 amino acids

tRNA val has 34 ss

21 pairs

∴ needs 55 states

= 5500 P_g

16 s rRNA 1000 nt 600
1000 pairs = 500

1100 ~~states~~ states

experiment 1. construct a tRNA model by hand
use it to search sequences.

note

needs $7 \times 55 = 385$ states

$\times 70$

100

$\times 100$

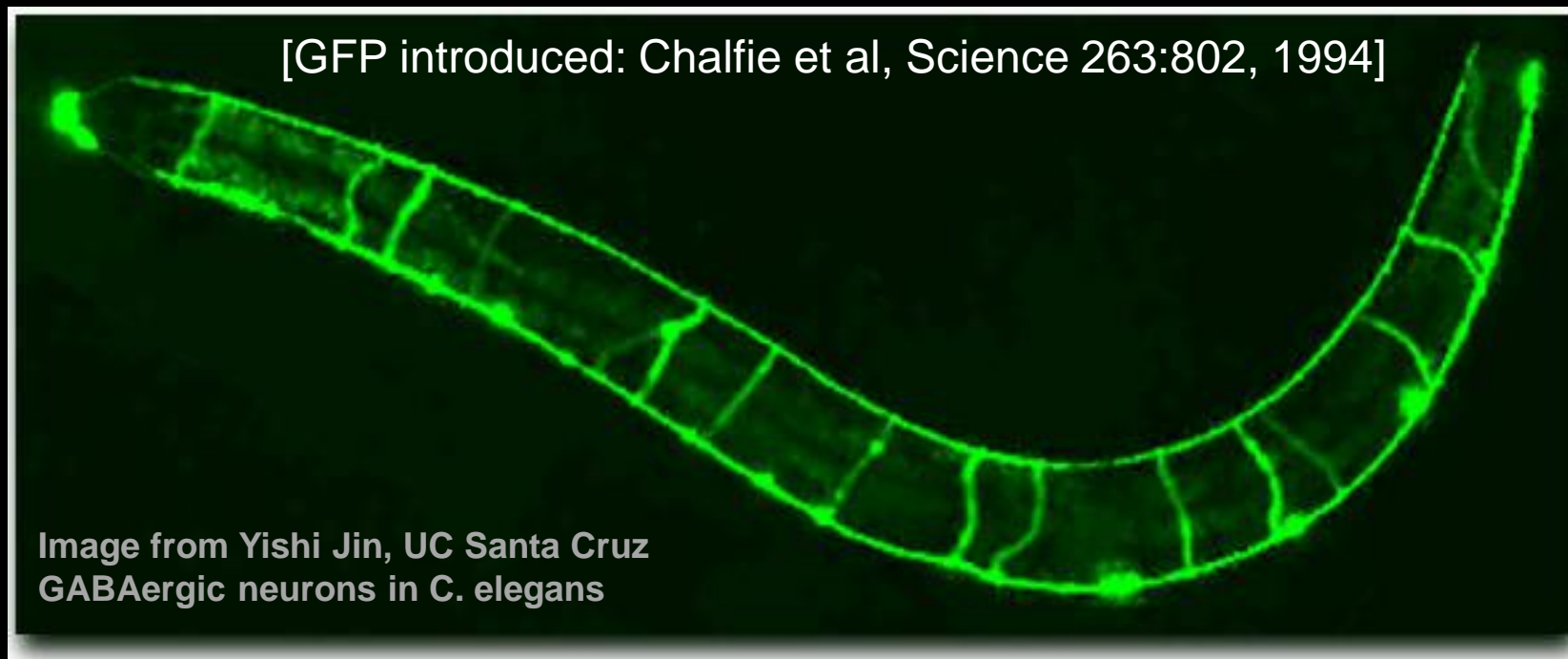
3.8 Mb

N^3 in space

I had proposed to develop reporter fusions to neural-specific promoters as a tool to visualize axonal processes in live animals, facilitating genetic screens...

I have started to play with... green fluorescent protein (GFP) from the jellyfish...

I found out at the worm meeting in June that Marty Chalfie's lab is onto the same idea. Chalfie has already obtained bright fluorescence in the axons of the six touch neurons...



Somewhat embarrassingly, my most productive work has been unrelated to the original proposal, resulting from some moonlighting as a computational biologist... and a lingering interest in RNA structure from my thesis work.... I invented a new kind of statistical model, related to HMMs, which can model the two-dimensional structure consensus of RNAs.

- progress report for my postdoc grant, 1993



HAMPSTEAD
GARDEN SUBURB

Lyttelton
Playing Fields

Golders Green
Crematorium

Hampstead Golf Club

Hampstead
Heath
Extension

Highgat

A598

Finchley Rd

GREEN

From sequence to RNA structure analysis

| Goal | HMM algorithm (sequence) | SCFG algorithm (RNA structure) |
|--|--|-----------------------------------|
| optimal alignment $P(\text{sequence} \mid \text{model})$ EM parameter estimation | Viterbi Forward Forward-Backward | CYK Inside Inside-Outside |
| memory complexity: | $O(MN)$ | $O(MN^2)$ |
| time complexity (general): | $O(M^2N)$ | $O(M^3N^3)$ |
| time complexity (as used): | $O(MN)$ | $O(MN^3)$ |

- we can analyze target sequences with secondary structure models;
- but the algorithms are computationally expensive.

Biological sequence analysis

Probabilistic models
of proteins and
nucleic acids

R. Durbin
S. Eddy
A. Krogh
G. Mitchison

CAMBRIDGE





The Plough

Cambridge North

Fen Ditton

Ditton Meadows

CHESTERTON

Stourbridge Common

Cambridge City Cemetery

Newmarket Road Park & Ride

Cambridge Museum Of Technology

McDonald's

Cambridge United

A1303

Tesco Superstore

B&Q

Currys PC World Featuring Carphone...

Marshall Aerospace and Defence Group

Marshall Ford Store Cambridge

Jesus Green

Midsummer Common

Homebase - Cambridge

Goldhams Common

Cambridge International Airport

The Round Church

Christ's College Cambridge

The Grafton Centre

TK Maxx

Asda Cambridge Superstore

Anglia Ruskin University

Mill Road Cemetery

Cambridge

Sedgwick Museum of Earth Sciences

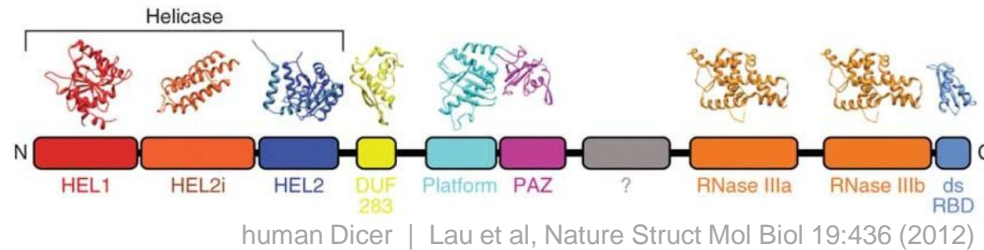
University Arms Hotel

Software tools for homology search and alignment

HMMER

protein homology search: profile HMMs
<http://hmmer.org>

55K lines source code
>40,000 downloads/yr



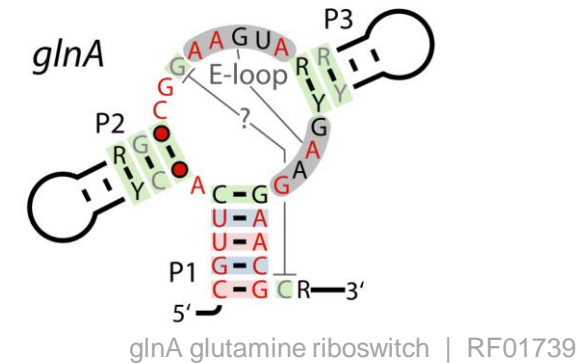
Pfam

17929 protein domain families
<http://pfam.xfam.org>

Infernal

RNA homology search: profile SCFGs
<http://eddylab.org/infernal>

90K lines source code
>10,000 downloads/yr



Rfam

3016 RNA structure families
<http://rfam.xfam.org>

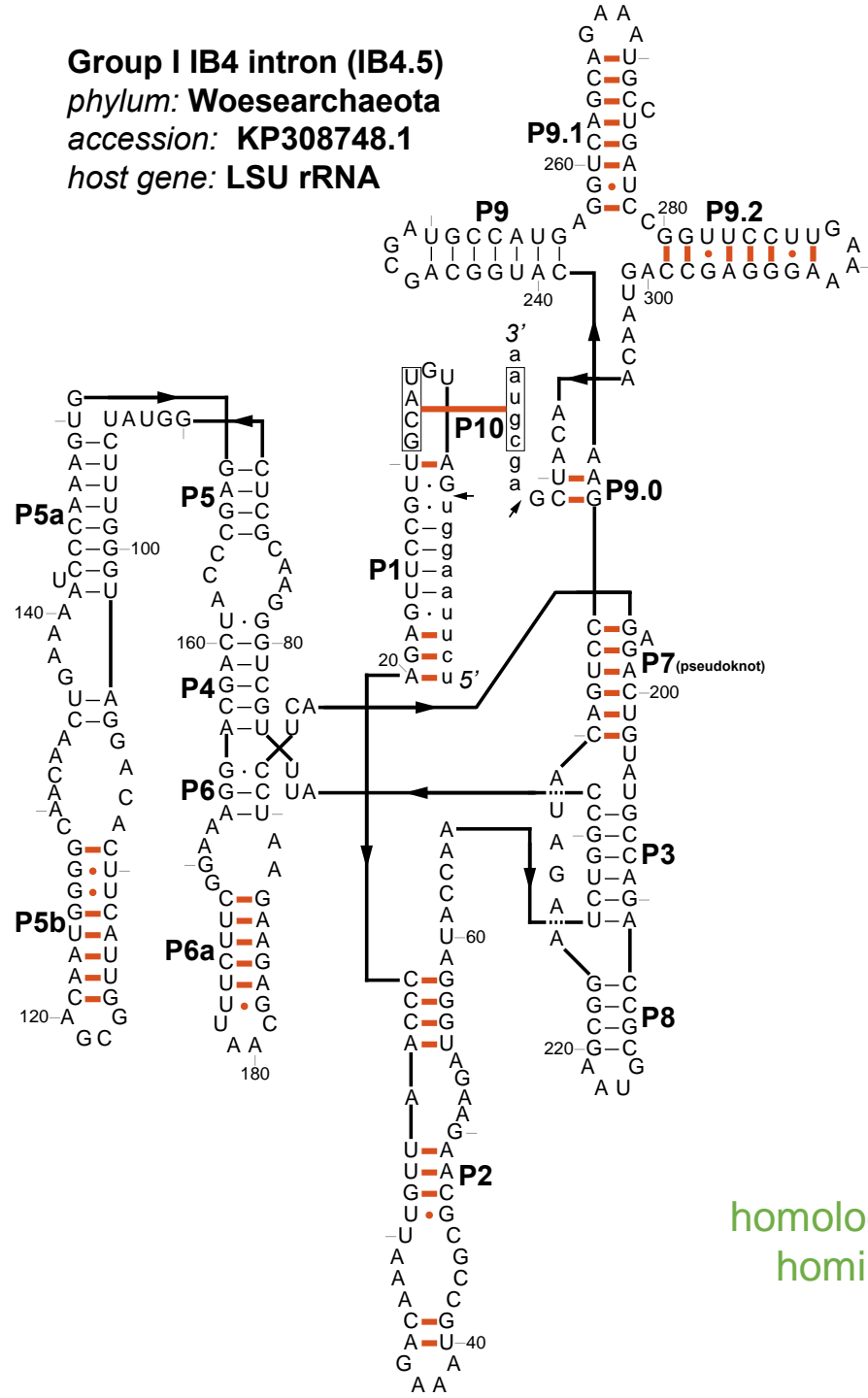
Xfam Consortium and HMMER server team: Rob Finn, Alex Bateman
European Bioinformatics Institute, Cambridge UK

Group I IB4 intron (IB4.5)

phylum: **Woesearchaeota**

accession: **KP308748.1**

host gene: **LSU rRNA**

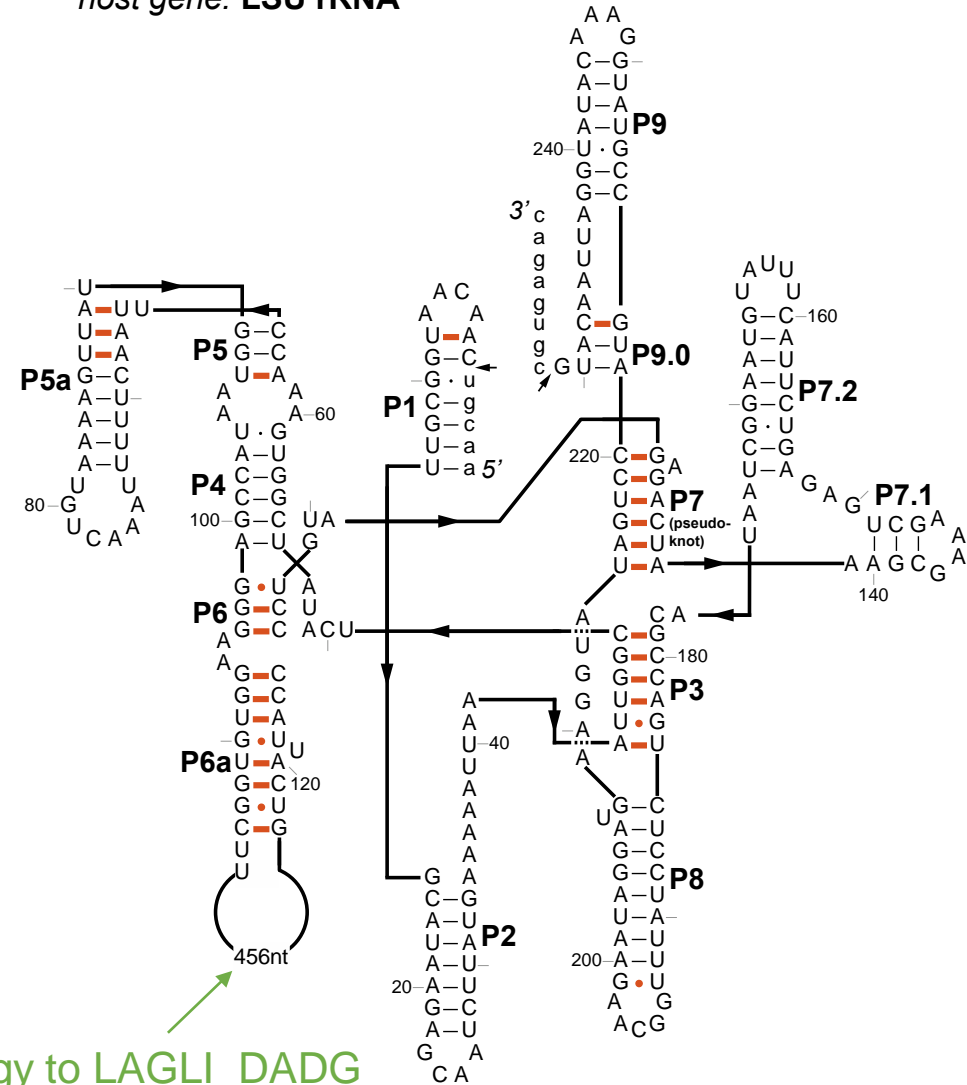


Group I IA3 intron (IA3.1)

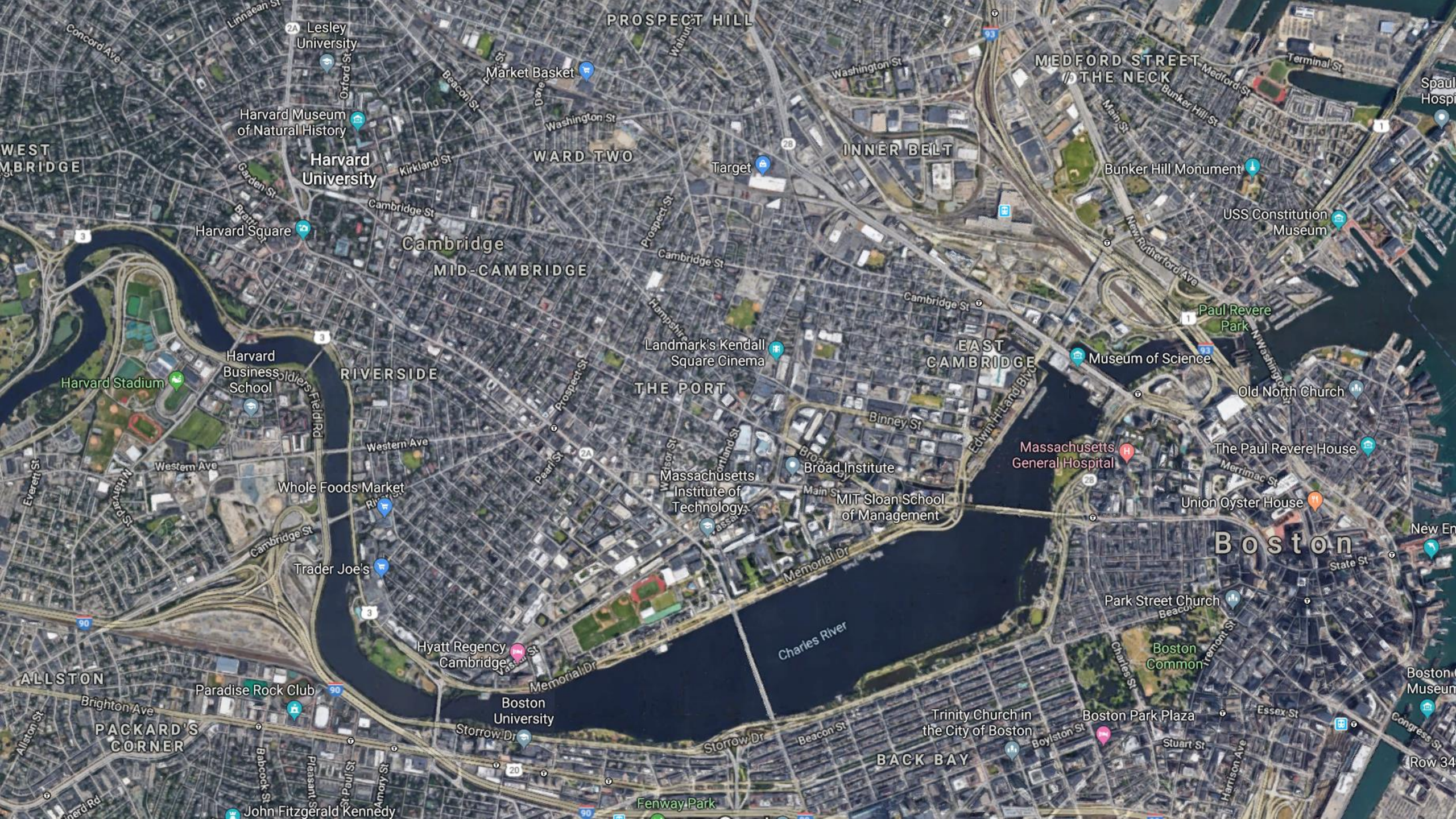
phylum: **Woesearchaeota**

accession: **CP010426.1**

host gene: **LSU rRNA**



homology to **LAGLI_DADG**
homing endonuclease
(PF00961)



PROSPECT HILL

MEDFORD STREET
(/ THE NECK

WEST
BRIDGE

Harvard
University

WARD TWO

INNER BELT

Cambridge

MID-CAMBRIDGE

Harvard
Business
School

RIVERSIDE

THE PORT

EAST
CAMBRIDGE

ALLSTON

PACKARD'S
CORNER

Boston
University

Charles River

Boston

BACK BAY

John Fitzgerald Kennedy

Fenway Park

Boston
Common

Row 34

“In my view, biology needs numbers; not after the fashion of physics, but in a good engineering, computational spirit.”

