# Biosimplicity: Engineering Simple Life

**Tom Knight**

MIT Computer Science and
Artificial Intelligence Laboratory

*Simplicity is the ultimate sophistication*
*-- Leonardo da Vinci*

# Measuring Complexity

- Number of components
- Size of description
- Size of functional model

Have you ever thought... about whatever man builds, that all of man's industrial efforts, all his calculations and computations, all the nights spent over working draughts and blueprints, invariably culminate in the production of a thing whose sole and guiding principle is the ultimate principle of simplicity?
In any thing at all, perfection is finally attained, not when there is no longer anything to add, but when there is no longer anything to take away.
    -- Antoine de Sainte Exupery, in "Wind, Sand, and Stars"
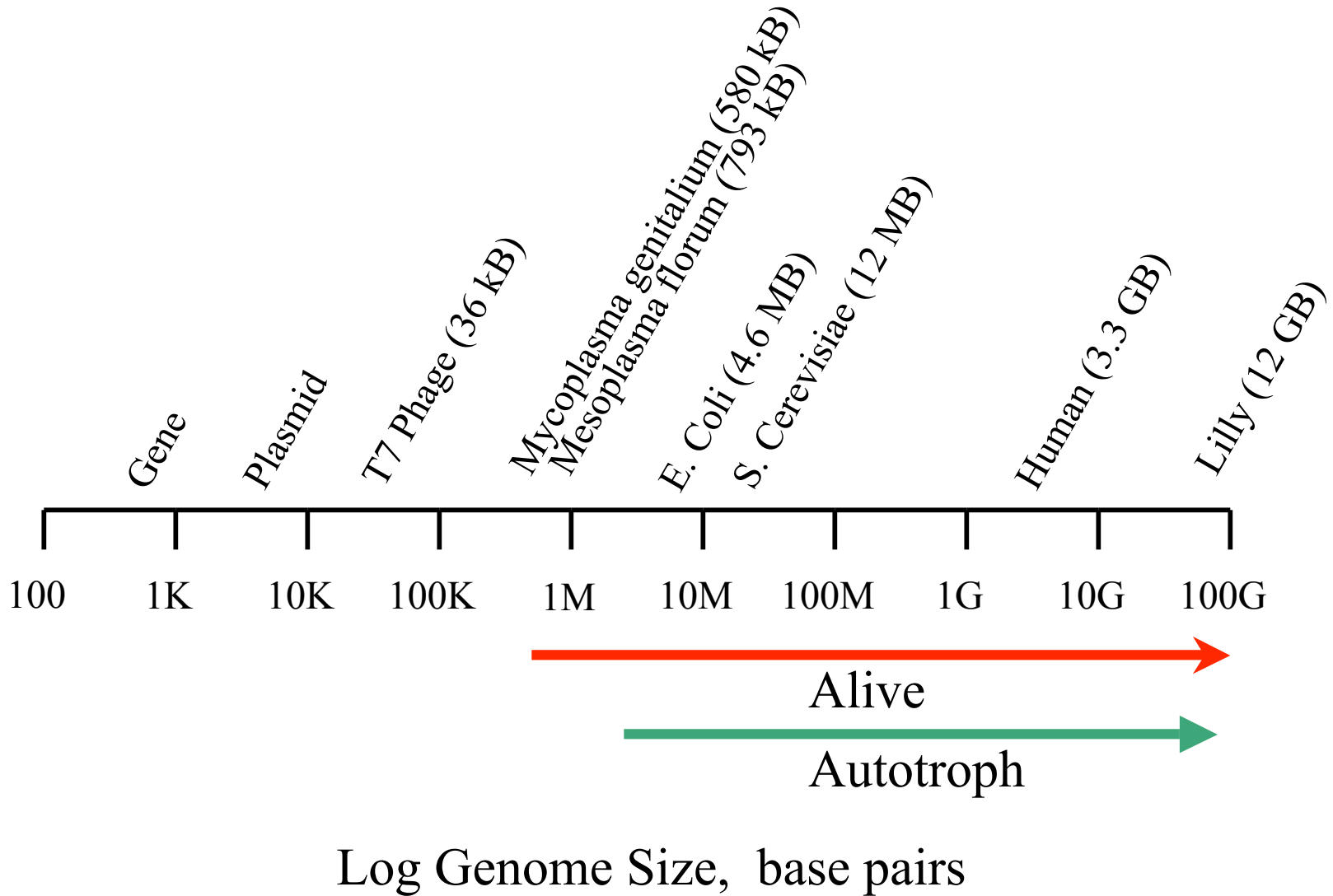
# Engineered Simple Organisms

- modular

- understood

- malleable

- low complexity

- Start with a simple existing organism

- Remove structure until failure

- Rationalize the infrastructure

- Learn new biology along the way

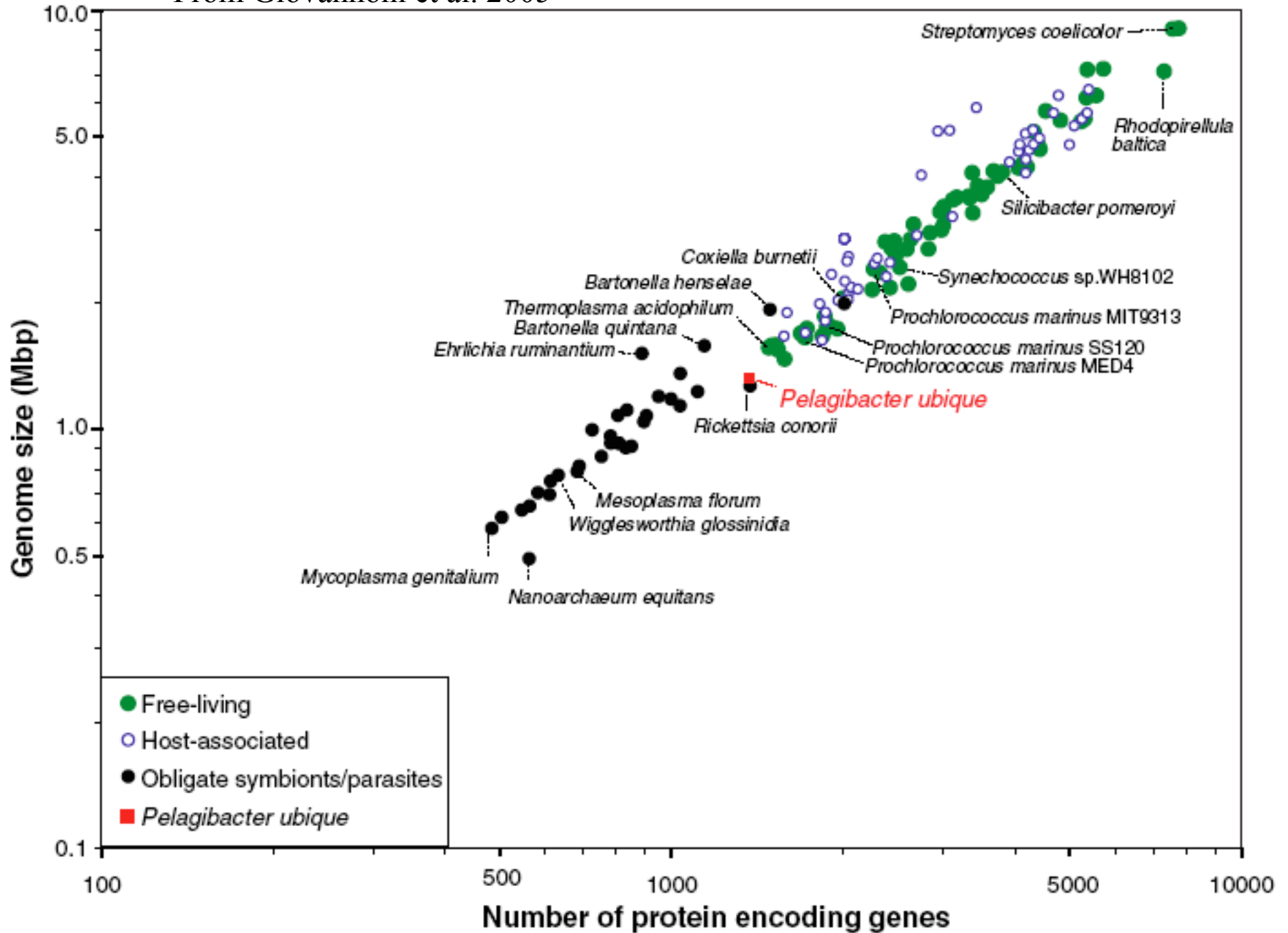**The chassis and power supply for our computing**

# Some history...

- Confusion over what PPLO/Mycoplasma were
  - ➤ "The Microbe of pleuorpneumonia" Nocard 1896
- 1932 isolation of "PPLO." Koch postulates.
- 1958 Klieneberger-Nobel identifies them as free living bacterial species
- Morowitz 1962 SciAm: "the smallest living cell"
- 1980 Gilbert effort to sequence *M. capricolum*
- 1982 Morowitz "complete understanding of life"
- 1996 Fraser et al. *M. genitalium* sequence
- 1999 Hutchison et al. Minimal genome set for *M. genitalium*

# Relative Complexity



Gene  Plasmid  T7 Phage (36 kB)  Mycoplasma genitalium (580 kB)  Mesoplasma florum (793 kB)  E. Coli (4.6 MB)  S. Cerevisiae (12 MB)  Human (3.3 GB)  Lilly (12 GB)

100  1K  10K  100K  1M  10M  100M  1G  10G  100G

Alive

Autotroph

Log Genome Size, base pairs

From Giovannoni et al. 2005

Figure axes — Genome size (Mbp) vs Number of protein encoding genes

Labeled data points:
- Streptomyces coelicolor
- Rhodopirellula baltica
- Silicibacter pomeroyi
- Coxiella burnetii
- Bartonella henselae
- Thermoplasma acidophilum
- Bartonella quintana
- Ehrlichia ruminantium
- Synechococcus sp.WH8102
- Prochlorococcus marinus MIT9313
- Prochlorococcus marinus SS120
- Prochlorococcus marinus MED4
- Pelagibacter ubique
- Rickettsia conorii
- Mesoplasma florum
- Wigglesworthia glossinidia
- Mycoplasma genitalium
- Nanoarchaeum equitans

Legend:
- Free-living
- Host-associated
- Obligate symbionts/parasites
- Pelagibacter ubique

# Choosing an organism

- Safe
  - ➤ BSL-1 organism        -- insect commensal
- Un-regulated
  - ➤ Not a crop plant or domesticated animal pathogen
- Fast growing
  - ➤ 40 minute doubling time
  - ➤ vs. six hours for M. genitalium
- Convenient to work with
  - ➤ Facultative anaerobe
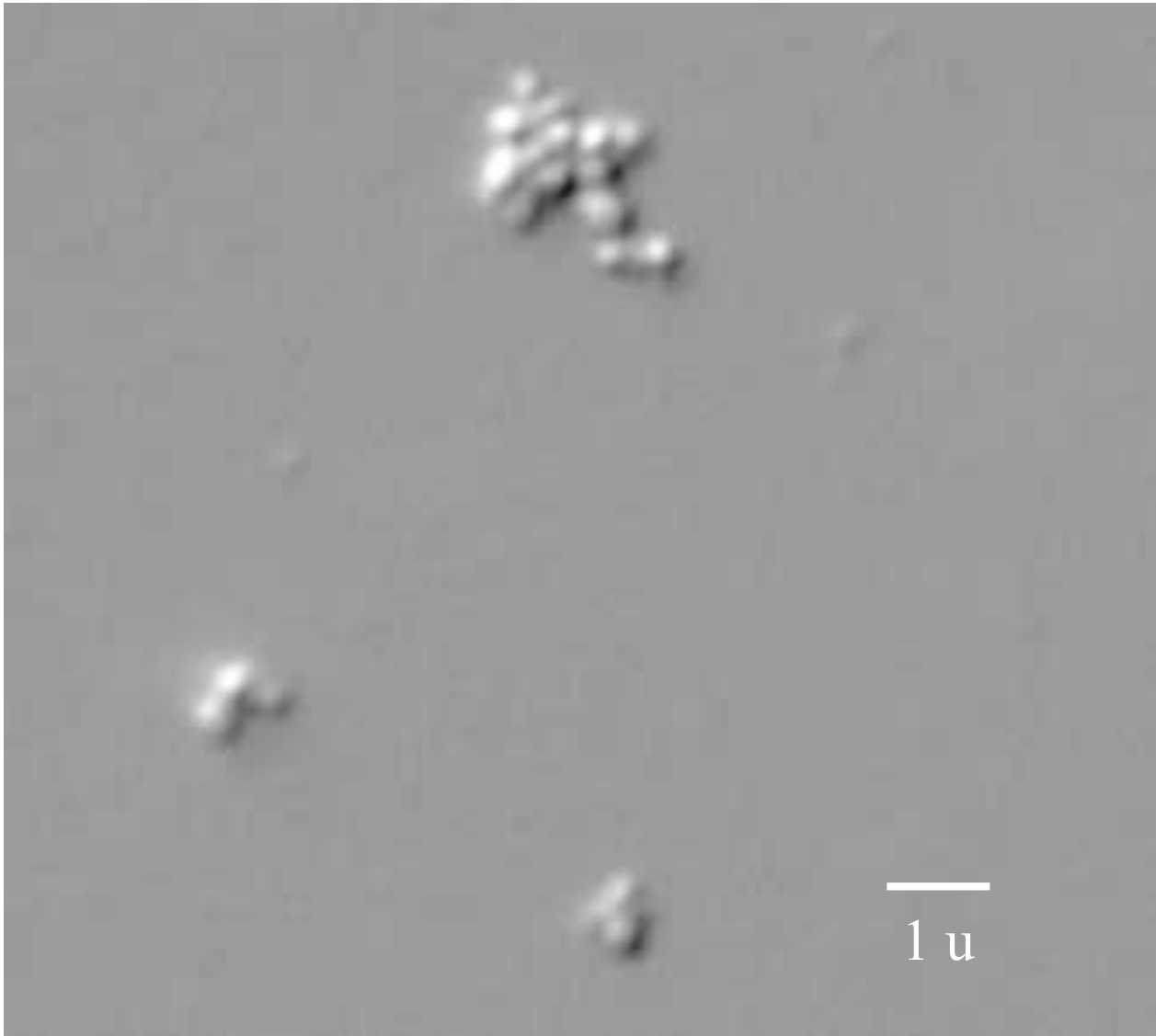- Small genome
- Known sequence
- Complete annotation

# The Mollicute Bibliome

Complete collection of mycoplasma related papers:

- 6,418 and counting

- All books and book chapters also

- Endnote & Refworks

- Downloaded .pdfs for articles > 1995

- Scanned articles and books, OCR with Abbyy Finereader

- Plans for a Google appliance search engine

- Collaboration for "shallow semantic" understanding

-  people.csail.mit.edu/tk/mfpapers/      user=meso, pass=meso

200 Micron

# Mesoplasma florum

# Tomographic TEM



Courtesy
Jensen Lab,
Caltech

# Culture Medium 1161

- Beef Heart infusion
- 4% Sucrose
- Fresh yeast culture broth
- 20% horse serum
- Penicillin
- Phenol red

# Defined medium

Sodium phosphate
KCl
Magnesium sulfate

Glycerol
Spermine

Nicotinic acid
Thiamine
Riboflavin
Pyridoxamine
Thioctic acid
Coenzyme A

Amino acids (minus asp, glu)
Guanine
Uracil
Thymine
Adenine

Glucose

BSA
Palmitic acid
Oleic acid

# Synthesis vs. Import

- Mycoplasma import virtually all small biochemical molecules
- Each import is done with a specific membrane protein – some are capable of importing a class
- Complexity is reduced if the import is simpler than the synthesis
- Example of the opposite:
  - Glutamine → Glutamic acid
  - Asparagine → Aspartic acid

# Sequencing the genome

- First sequenced small portions of the genome to test that we had the correct species
  - Compared the results to Genbank entries
  - Sequenced PTS system gene, identical to reported sequence
  - Sequenced 16S rRNA (unreported)
  - Discovered identical to Mesoplasma entomophilum 16S rRNA sequence – probably the same species
    - Genbank entry
- Measured genome size with pulse field gels
- Sequenced 12% of genome to see what we were up against
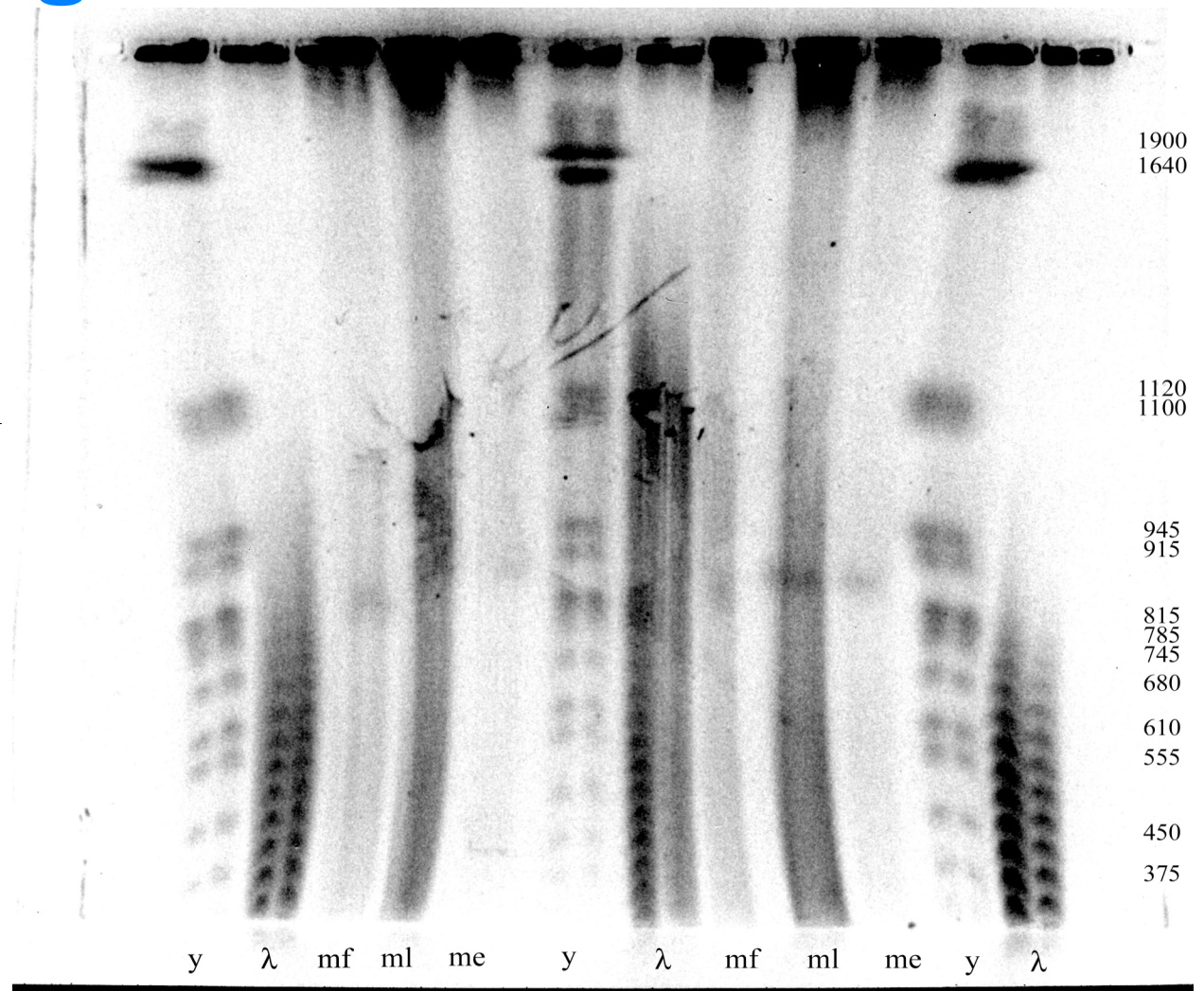
# PFGE of Mesoplasma genomic DNA

y
yeast marker

I
lambda marker

me
Mesoplasma entomophilum

mf
Mesoplasma florum

ml
Mesoplasma lactucae

1900
1640
1120
1100
945
915
815
785
745
680
610
555
450
375

y  λ  mf  ml  me  y  λ  mf  ml  me  y  λ

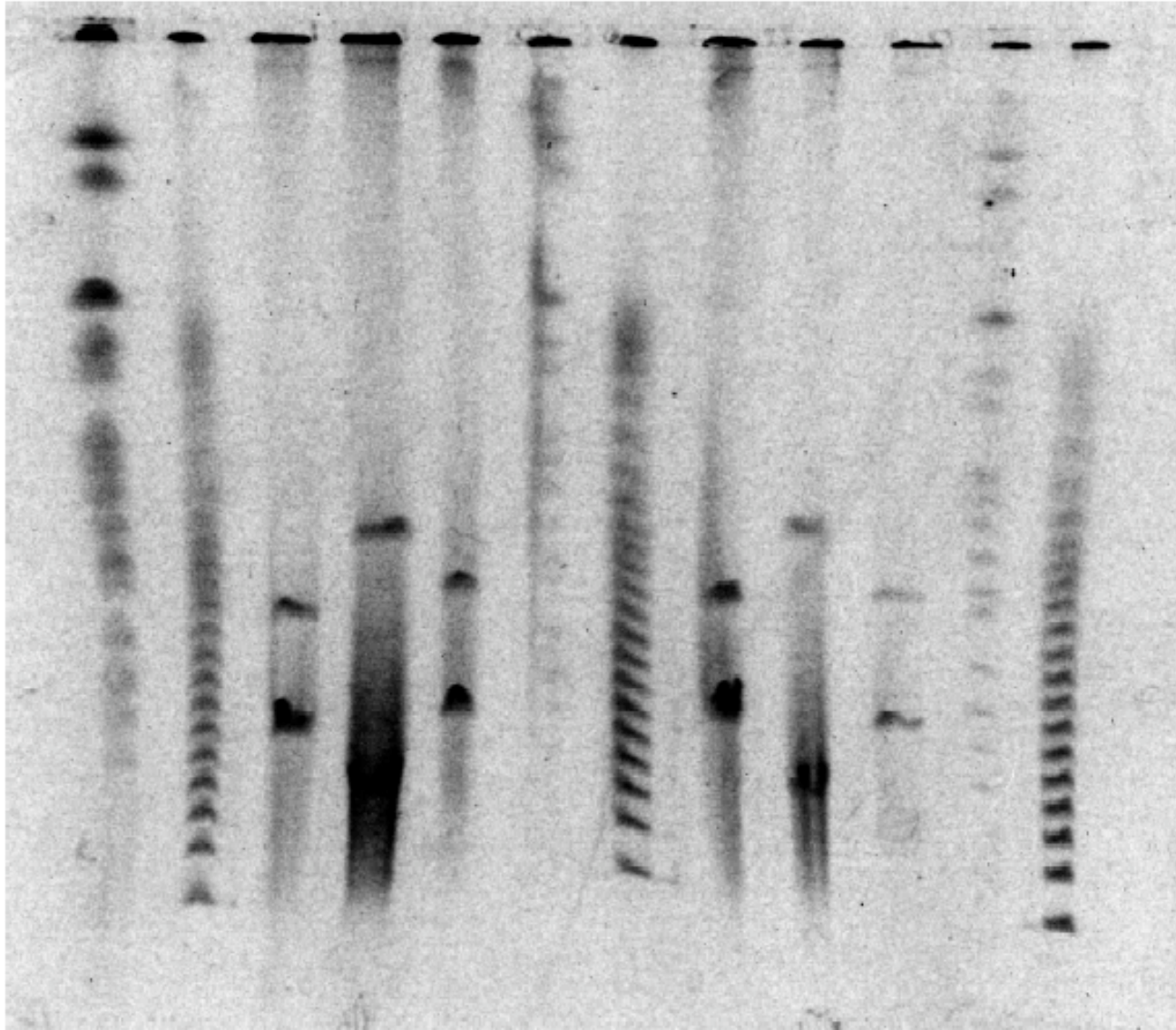1.2% agarose  9C   6V/cm  Ramped 90 - 120 sec  48 hours

# I-CeuI digestion

- Special restriction enzymes cut only at 23S rRNA sites

```
5´  ..TAACTATAACGGTC CTAA^GGTAGCGA..3´

3´  ..ATTGATATTGCCAG^GATT CCATCGCT..5´
```

Calculate the number of rRNA sites in the genome from the number of cut fragments
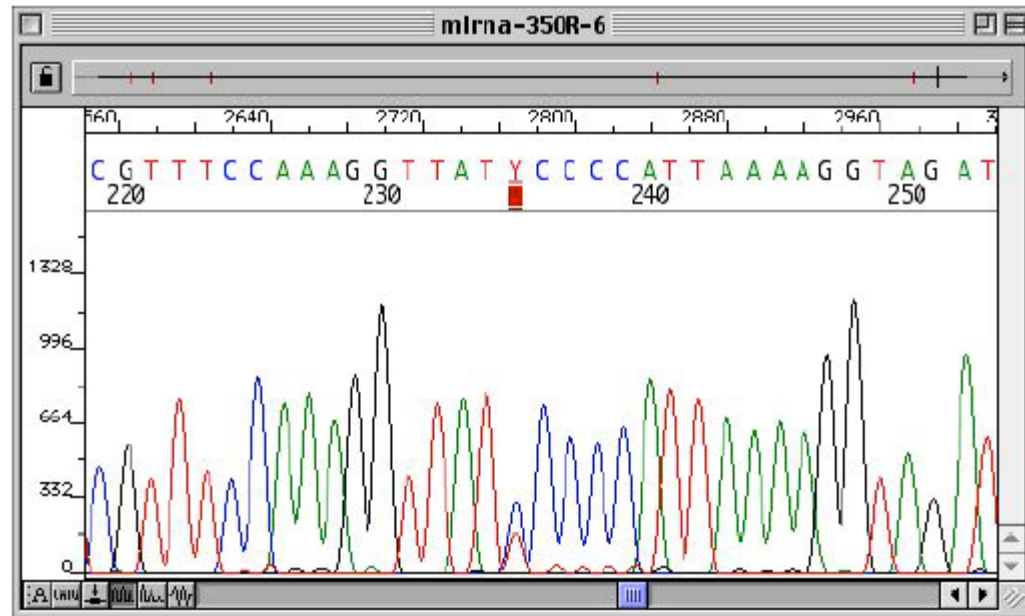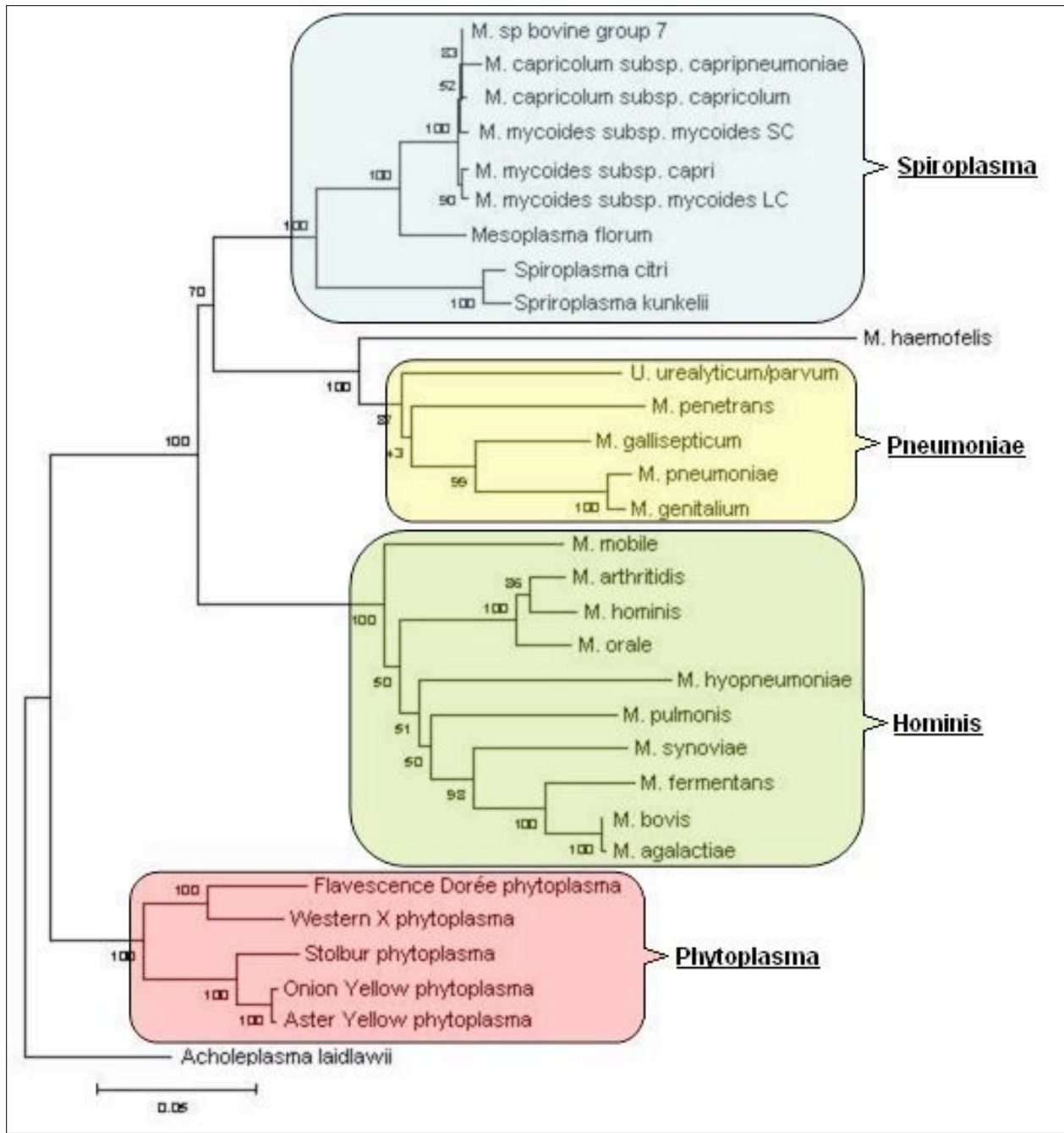
# I-CeuI digests

# rRNA sequences

- Degenerate primers

| Primer | F/R | Sequence 5' to 3' |
|--------|-----|-------------------|
| U1 | F | gtt tga tcc tgg ctc agg ayd aac g |
| U2 | F | čca rac tcc tac ggr agg cag c |
| U5 | R | čtt gtg cgg gyy ccc gtc aat tc |
| U8 | R | ḡaa agg agg trw tcc ayc csc ac |
| 1300 | F | ṭaa tcg cga atc agc tat gtc |
| 350 | R | ṭgc ttc atc aga ctt tcg tcc |
| 1000 | F | ṭgg agg tta aca ttg ata cag g |
| 1150 | R | čat gat gat ttg acg tca tcc |

- The two rRNA sequences are different:
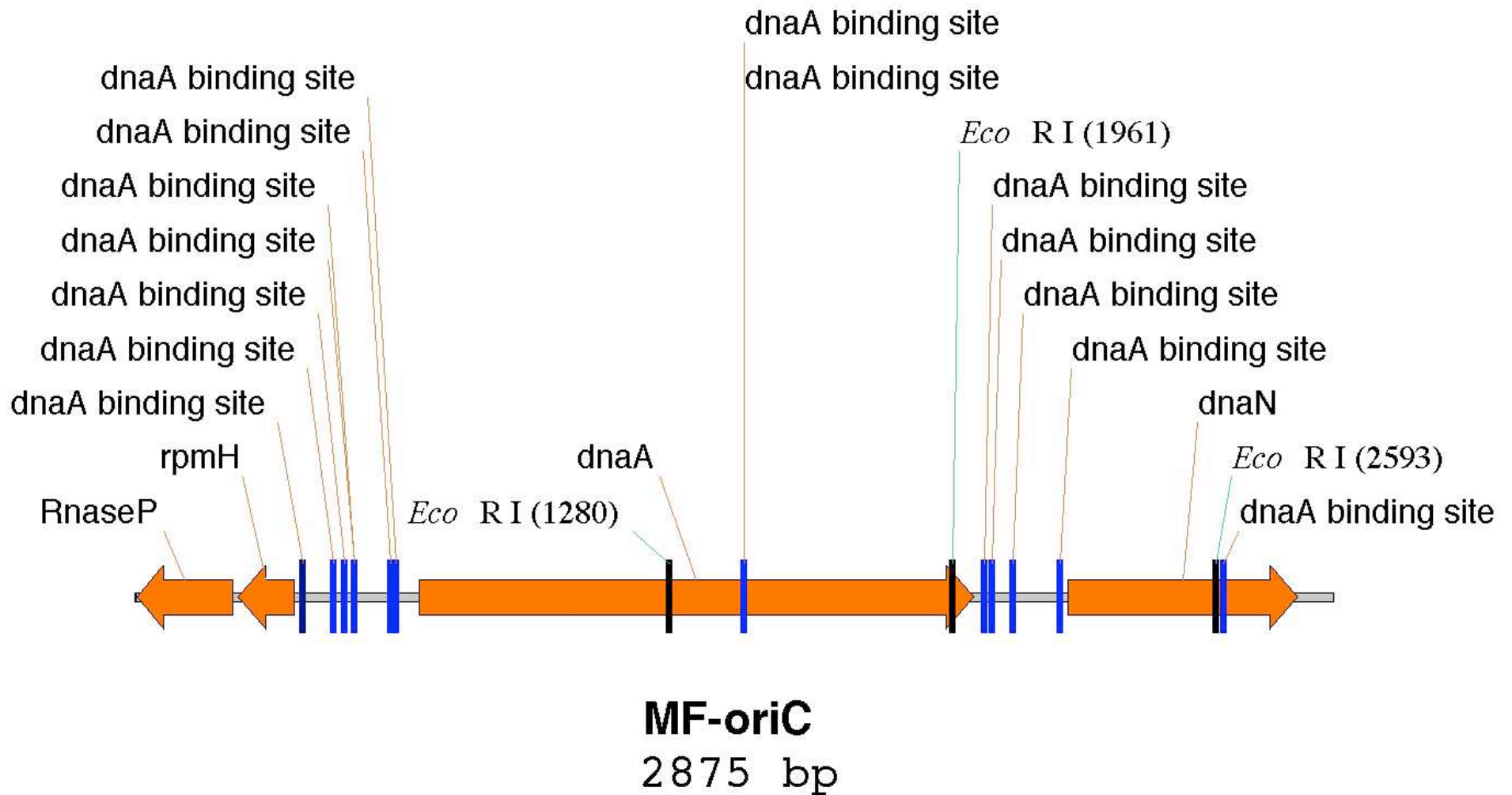
# Library creation

- Randomly cut genomic DNA with EcoRI
- Shotgun cloned into pUC18 vector
- Sequenced the inserts  (0.1 – 8 Kb)

- Sheared the genomic DNA with a needle
- End repaired
- Cloned into defective lambda phage vector
- Packaged vector into phage heads
- Infected E. coli cells with phage
  - ➤ ~ 40 Kb inserts

# What we learned from partial sequence

- Almost all "old friends"
- Little or no extra junk
- Inter-gene sequences small (-4 to 30 bp)
- Little transcriptional control
- High AT vs. GC content – 27% average GC
  - ➢ 20% in typical genes, even less in control regions
  - ➢ 40% in rRNA regions

# Origin of Replication



MF-oriC
2875 bp

# Whitehead Agreement

- Whitehead agreed in January 2002 to sequence the organism
- Estimated to take about two hours of time on their sequencers
  - ➢ "Sure, we can do it Tom, but what do we do with the rest of the day after the coffee break?"
- "How many other organisms like this are there?"
  - ➢ 300
- "Why don't we sequence them all?"
- Good draft available November 2002
- Gaps closed July 2003 -- final November 2003

# Gap closure

- 9 gaps remained
- Long range PCR
- Primer walking
- One difficult sequence
  - Poly A region 16-17 bp long
  - Sequencing stuttered
  - Reprime with aaaaaaaaaaaaaaag

- Repeat region 186 bp in surface lipoprotein
  - Give up on accurate sequence, PCR for length
- Final assembly verification
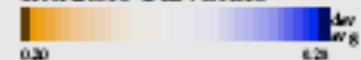
# Genome characteristics

- 793281 base pairs
- 26.52% G + C
- 682 protein coding regions
  - UGA for tryptophan
  - No CGG codon or corresponding tRNA
  - Classic circular genome
  - oriC, terminator region, gene orientation
- 39 stable RNAs
  - 29 tRNAs
  - 2x 16S, 23S, 5S
  - RNAse-P, tmRNA, SRP

# Standard Motifs

- -10 present usually very highly conserved
  - ➢ Often preceded by a "TG" 1-2 bp upstream
- Seldom a conserved -35 region
- RBS is standard Shine-Dalgarno
- Alternate RBS matches complementary region of 16S rRNA UAACAACAU (Loechel 91)
- Standard stem-loop terminators with loop TTAA
  - ➢ 6-8 poly T tail in forward direction
- dnaA box TTATCCACA
- Four ribo-box sequences (Thi, Ile, Val, Guanine)

Intrinsic Curvature

0.20 ... 0.28

Stacking Energy

-7.40 ... -6.32

Position Preference

0.14 ... 0.19

Annotations:

CDS +
CDS -
rRNA
tRNA

Global Direct Repeats

0.00 ... 7.00

Global Inverted Repeats

0.00 ... 7.00

GC Skew

-0.12 ... 0.12

Percent AT

0.20 ... 0.80

Resolution: 318

http://www.cbs.dtu.dk/
Center for Biological Sequence Analysis

GENOME ATLAS

M. florum
793,139 bp

# Understand the metabolism

- Identify major metabolic pathways by finding critical genes coding for known enzymes
- Predict necessary enzymes which may not have been found
- Evaluate the list of unknown function genes for candidates
- Build the major metabolic pathway map of the organism
- Consider elimination of entire pathways

# Identified Metabolic Pathways in *Mesoplasma florum*

G. Fournier
02/23/04

PTS II System

glucose  sucrose  trehalose  xylose  beta-glucoside  unknown  fructose

ribose ABC transporter

ATP Synthase Complex

ATP ← ADP

sn-glycerol-3-phosphate ABC transporter

Pentose-Phosphate Pathway

chitin degradation

Glycolysis

glucose-6-phosphate

Lipid Synthesis

fatty acid/lipid transporter

glyceraldehyde-3-phosphate

L-lactate, acetate

acetyl-CoA

unknown substrate transporters

ribose-5-phosphate

PRPP

cardiolipin/ phospholipids

membrane synthesis

phospholipid membrane

Purine/Pyrimidine Salvage

niacin?

Pyridine Nucleotide Cycling

variable surface lipoproteins

xanthine/uracil permease

hypothetical lipoproteins

DNA Polymerase

RNA Polymerase

Electron Carrier Pathways

NAD+

NADP

K+, Na+ transporter

competence/ DNA transport

DNA    RNA

degradation

ribosomal RNA    transfer RNA

Flavin Synthesis

FMN, FAD

malate transporter?

metal ion transporter

hypothetical transmembrane proteins

riboflavin?

NADPH    NADH

cobalt ABC transporter

Signal Recognition Particle (SRP)

Ribosome

tRNA aminoacylation

Formyl-THF Synthesis

phosphonate ABC transporter

protein secretion (ftsY)

messenger RNA

met-tRNA formylation

phosphate ABC transporter

Export    proteins

THF?

formate/nitrate transporter

protein translocation complex (Sec)

degradation

amino acids

intraconversion?

spermidine/putrescine ABC transporter

oligopeptide ABC transporter

putrescine/ornithine APC transporter

arginine/ornithine antiporter

glutamine ABC transporter

alanine/Na+ symporter

glutamate/Na+ symporter

lysine APC transporter

unknown amino acid ABC transporter

Amino Acid Transport

# How Simple is this?

- Missing cell wall, outer membrane
- Missing TCA cycle
- Missing amino acid synthesis
- Missing fatty acid synthesis
- One sigma factor
- Small number of dna binding proteins
- One insertion sequence, probably not active
- One restriction system (Sau3AI-like)
- CTG/CAG methylation  (function?)
- Evidence for shared protein function
  - ➢ MDH/LDH (Pollack 97 Crit rev microbiol 23:269)

# Minimal is not always simple

- Shared function of parts
- Overlapped genes
- Tradeoff of import vs. synthesis

- Example:
  - Television set design
  - Shared deflection coil, high voltage power supply, isolated filament supply
  - Three functions, a single circuit, a difficult engineering, modeling, debugging, and repair task

- *how many genes have multiple functions*

# DNA Methylation

- Bisulfite conversion of genomic DNA
- Sequencing of converted DNA to identify methylated C positions

- Results: GATC sequences, as expected
- Unexpected: CAG and CTG sequences

# Current work

- Array experiments
  - Close species
  - Transcriptional units, pseudo-genes
  - TRASH
- Protein species by LC/LC/MS/MS
- Elimination of  the restriction system
- Plasmid system
  - pBG7AU based
- Recombination system
  - Positive/negative selection
- Yeast chromosome transfer
- Genome edits to reduce size
- Genome edits to modularize
- Genome edits to eliminate complexity
- Use as a construction chassis

# Reconstruct the Genome

- Use recombination techniques to edit the genome
- Eliminate unnecessary genes
- Remove overlaps
- Standardize promoters, ribosomal binding sites
- Identify transcriptional and translational regulators
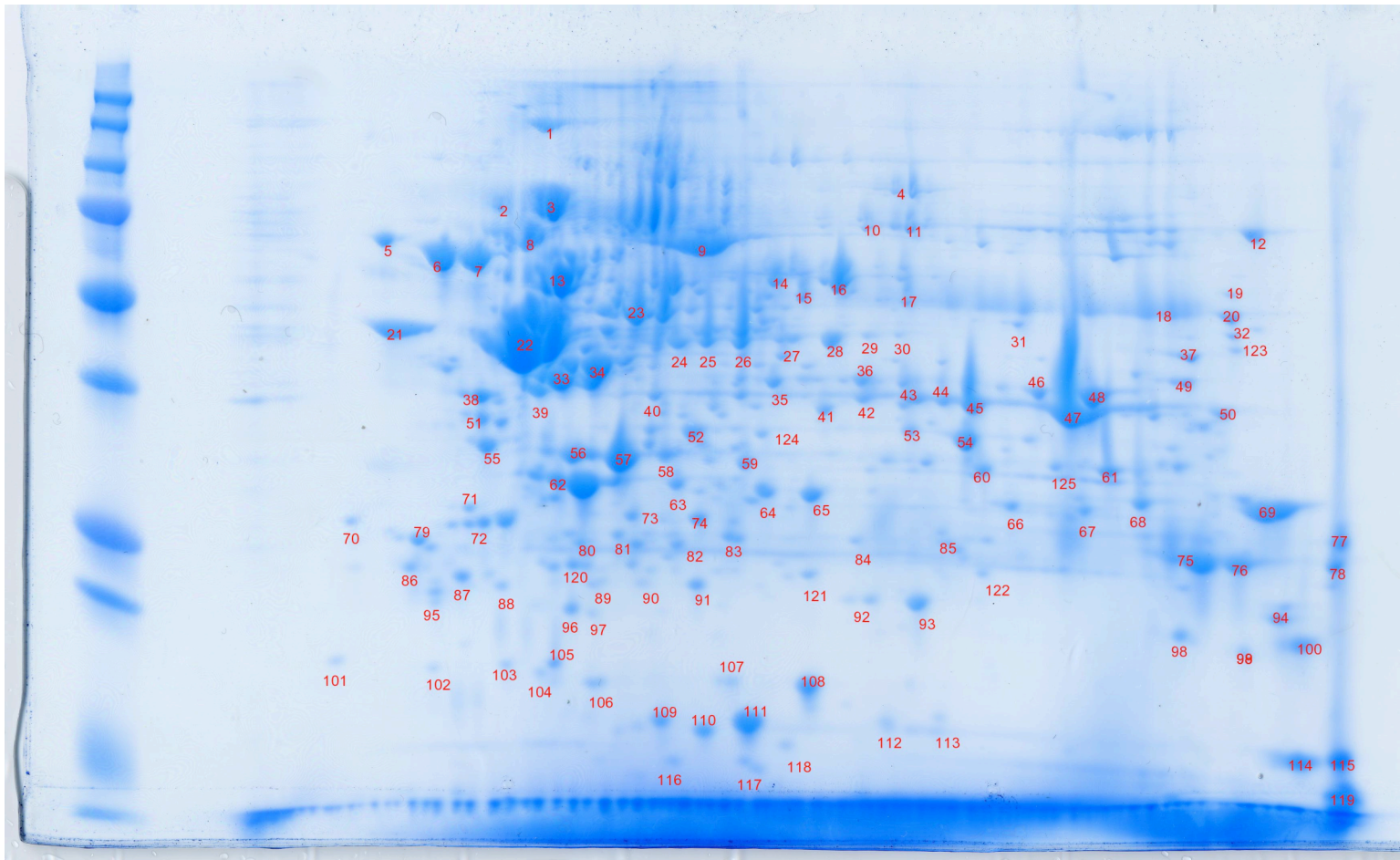- Recode proteins to use a reduced portion of the coding space

*The code is 4 billion years old, it's time for a rewrite*

# YAC mutagenesis

- Bring up YAC technology
  - Spheroplasts, old YAC plasmid sequencing
- Triple transform with MF chromosome and
  - pRML1    (Spencer 92)
  - pRML2
  - Genome inactive except for ARS, telomeres, selection markers
- Use yeast recombination systems for genome editing
- Isolate and recircularize YAC to form a new genome
- Lipid encapsulate genome into vesicles
- Fuse vesicles with genome-killed wild type cells

# Proteome

- Collaboration with Steve Tannenbaum / Yingwu Wang
- 2-D gels + MS spot ID
- LC/LC/MS/MS ID of trypsin digests

# Riboswitch Analysis

- Collaboration with Ron Breaker / Adam Roth

- Discovery of unique riboswitches specific for GTP rather than dGTP

- Found in no other sequenced genomes

- Analysis of close relatives under way

# Engineer plasmids

- No known plasmids for this class of organism
- Renaudin has made plasmids using the chromosomal OriC as the replicative element
  - Lartigue 03
- M. mycoides has pADB201 and pBG7AU rolling circle plasmids similar to pE194  (1082 bp)
- We know the antibiotic sensitivities and have working resistance genes

# Kit Part the genome

- Make Biobrick parts from each gene, tRNA, promoter, other part-like genome element

- Attempt to develop techniques for recombining parts into coherent modules

- Develop techniques for assembling and modeling the resulting structures

# Thanks to...

- Harold Morowitz
- Greg Fournier
- Gail Gasparich
- Bob Whitcomb
- Eric Lander
- Bruce Birren
- Nicole Stange-Thomann
- George Church
- Roger Brent
- Grant Jensen
- Yingwu Wang

- Nick Papadakis
- Ron Weiss
- Drew Endy
- Randy Rettberg
- Austin Che
- Reshma Shetty
- MIT Synthetic biology working group
- DARPA, NTT, NSF, Microsoft

# Synthetic Biology

- An alternative to understanding complexity is to remove it
- This complements rather than replaces standard approaches
- Engineering synthetic constructs will be easier
  - Enabling quicker more facile experiments
  - Enabling deeper understanding of the basic mechanisms
  - Enabling applications in nanotechnology, medicine and agriculture

*Simplicity is the ultimate sophistication*
*-- Leonardo da Vinci*