

# **Pathogen Genome Data**

## **My life in sequences**



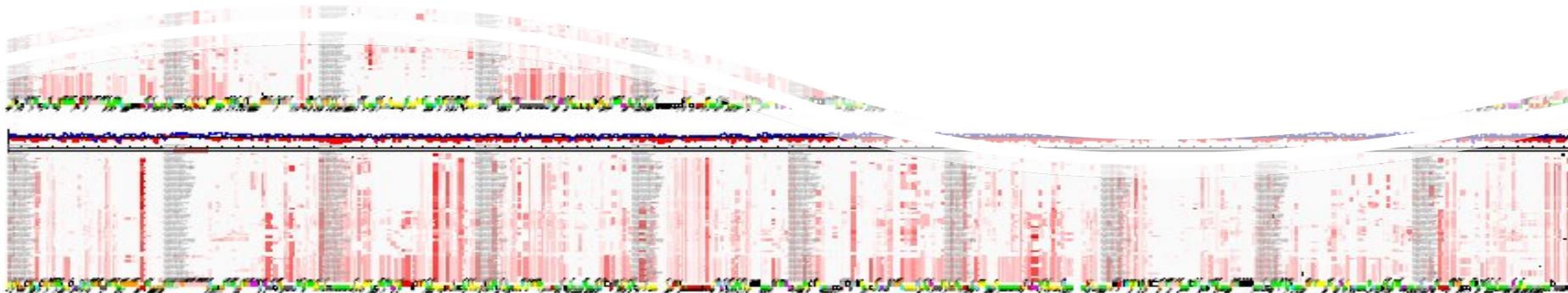
**The James  
Hutton  
Institute**

Leighton Pritchard



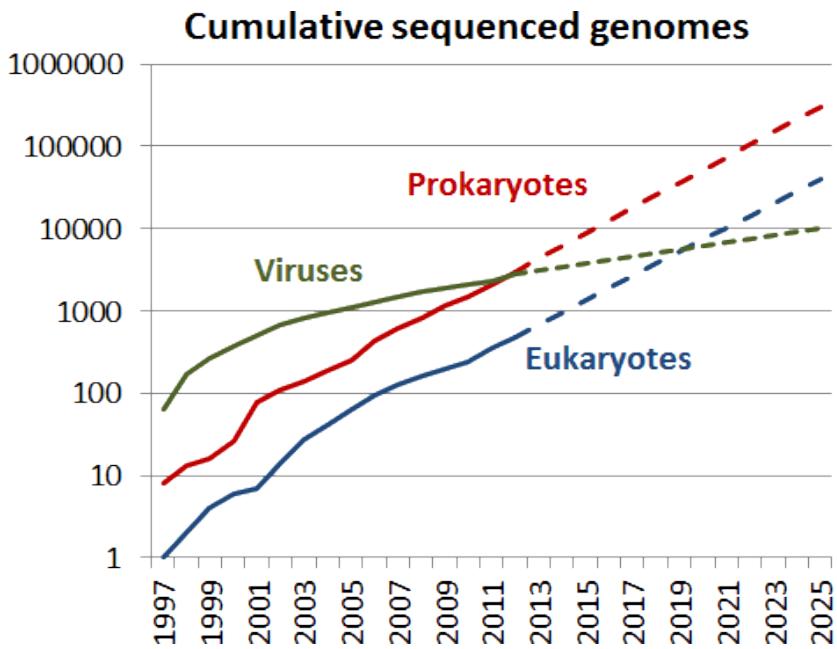
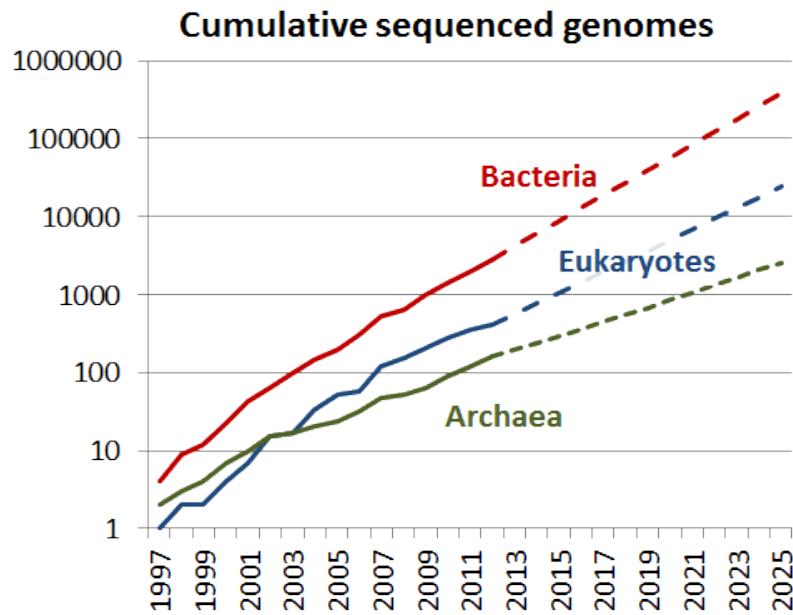
# Introduction

- 2003: *Erwinia carotovora* subsp. *atroseptica*
- 2009: *Phytophthora infestans*
- 2013: *Dickeya* spp.

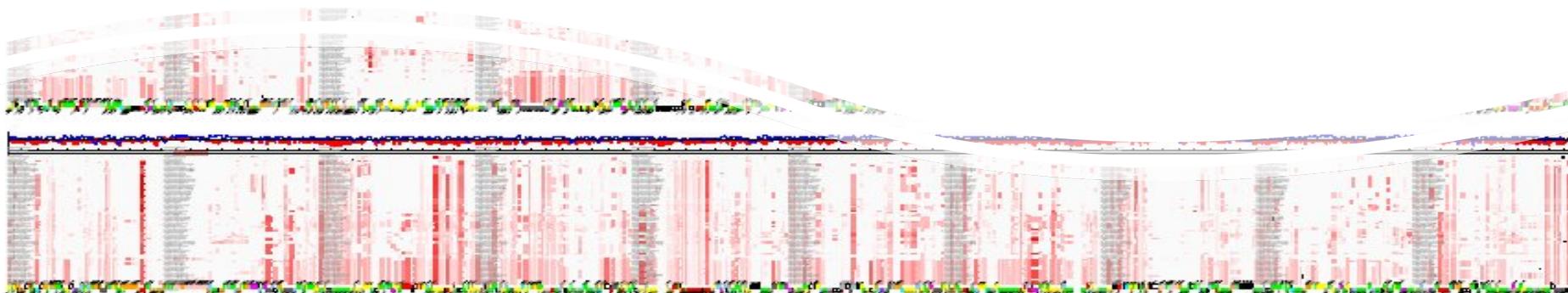




# Introduction



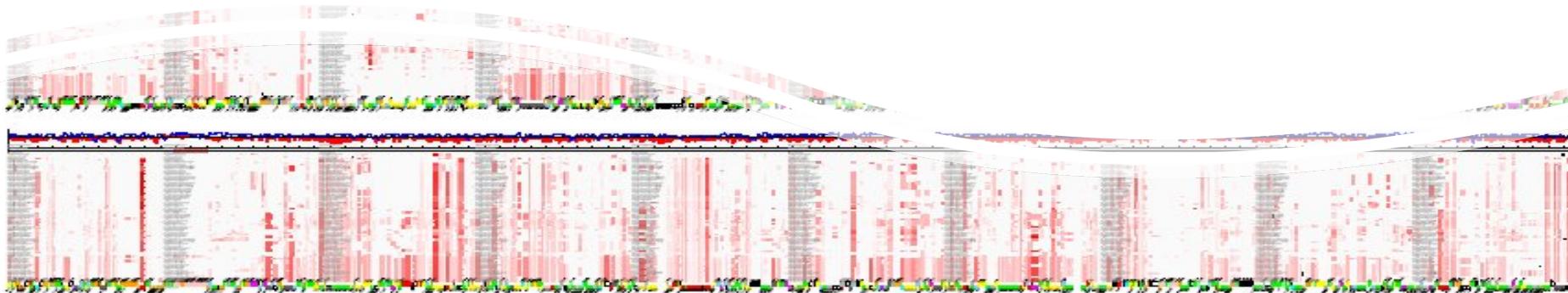
Figures and code from: <http://sulab.org/2013/06/sequenced-genomes-per-year/>





# Introduction

- As the number of sequenced genomes rises...
  - The collective power of genomics to explain biology increases...
  - The impact of a single genome falls...
  - The impact of a single genome paper falls faster...
  - The time available for any scientist to annotate and examine an individual genome in detail falls even faster than that!
- I discovered this the hard way...





# Introduction

- What can pathogen genome data do for you?

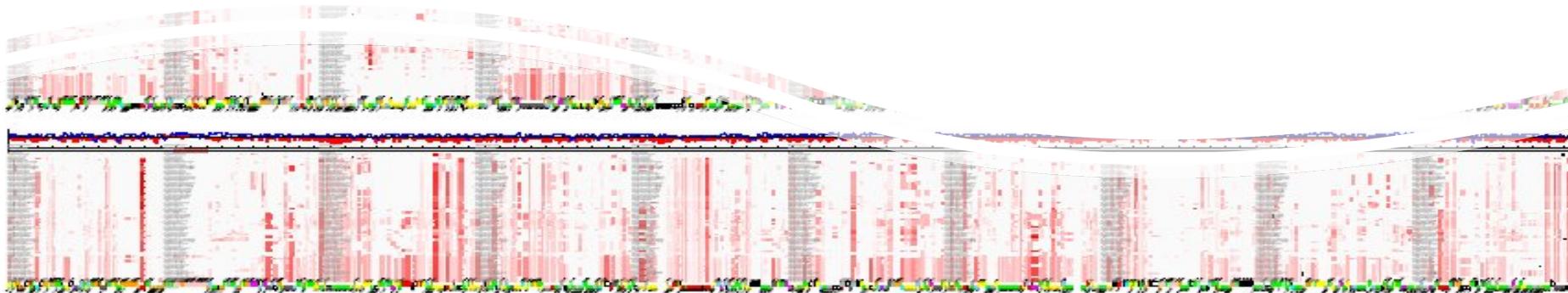
Catalogue of genome components

Differences in genomic complement: hypotheses for function

Characteristic signals for diagnosis and epidemiology

Parts lists and wiring diagrams for systems and synthetic biology

- How do we get to biology from genome data?





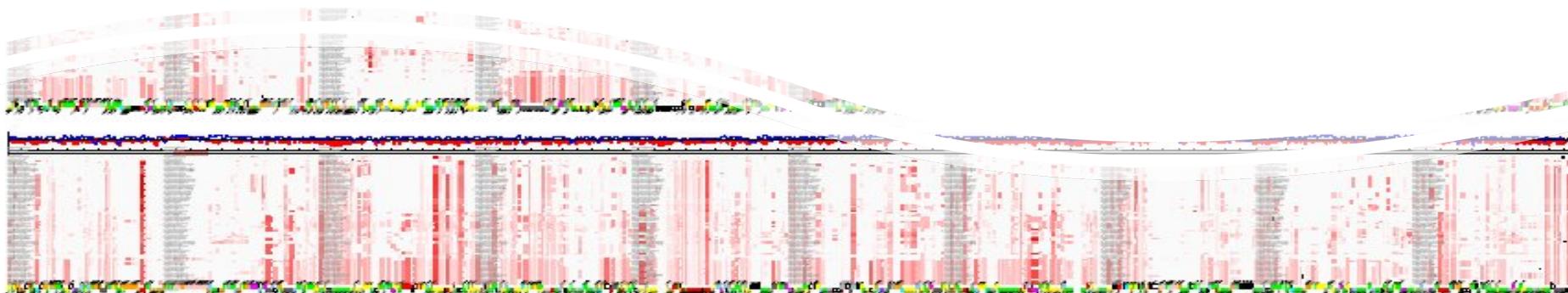
# 2003: *Erwinia carotovora* subsp. *atroseptica*

## Genome sequence of the enterobacterial phytopathogen *Erwinia carotovora* subsp. *atroseptica* and characterization of virulence factors

K. S. Bell<sup>\*†</sup>, M. Sebaihia\*, L. Pritchard<sup>†</sup>, M. T. G. Holden\*, L. J. Hyman<sup>†</sup>, M. C. Holeva<sup>†</sup>, N. R. Thomson\*, S. D. Bentley\*, L. J. C. Churcher\*, K. Mungall\*, R. Atkin\*, N. Bason\*, K. Brooks\*, T. Chillingworth\*, K. Clark\*, J. Doggett\*, A. Fraser\*, Z. Hance\*, H. Hauser\*, K. Jagels\*, S. Moule\*, H. Norbertczak\*, D. Ormond\*, C. Price\*, M. A. Quail\*, M. Sanders\*, D. Walker\*, S. Whitehead\*, G. P. C. Salmond<sup>‡</sup>, P. R. J. Birch<sup>†</sup>, J. Parkhill<sup>\*§</sup>, and I. K. Toth<sup>†§</sup>

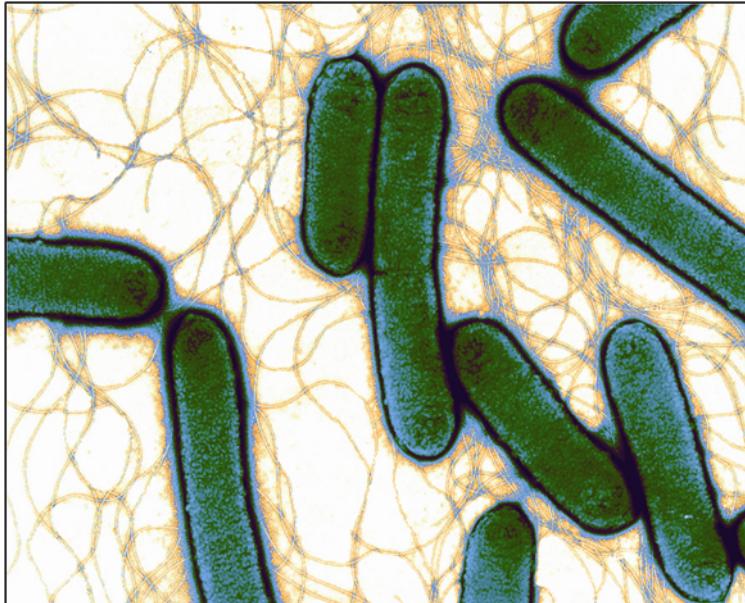
\*The Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, United Kingdom; <sup>†</sup>Plant–Pathogen Interactions Programme, Scottish Crop Research Institute, Invergowrie, Dundee DD2 5DA, United Kingdom; and <sup>‡</sup>Department of Biochemistry, Tennis Court Road, Cambridge University, Cambridge CB2 1QW, United Kingdom

Bell *et al.* (2004) Proc Natl Acad Sci USA [doi:10.1073/pnas.0402424101](https://doi.org/10.1073/pnas.0402424101).



# *Pectobacterium atrosepticum*

- Plant pathogenic enterobacterium
- *E. carotovora* subsp. *atroseptica* reclassified as *Pectobacterium atrosepticum*
- Causes blackleg (stem rot), rotting of stored tubers
- Major rot symptoms due to PCWDEs



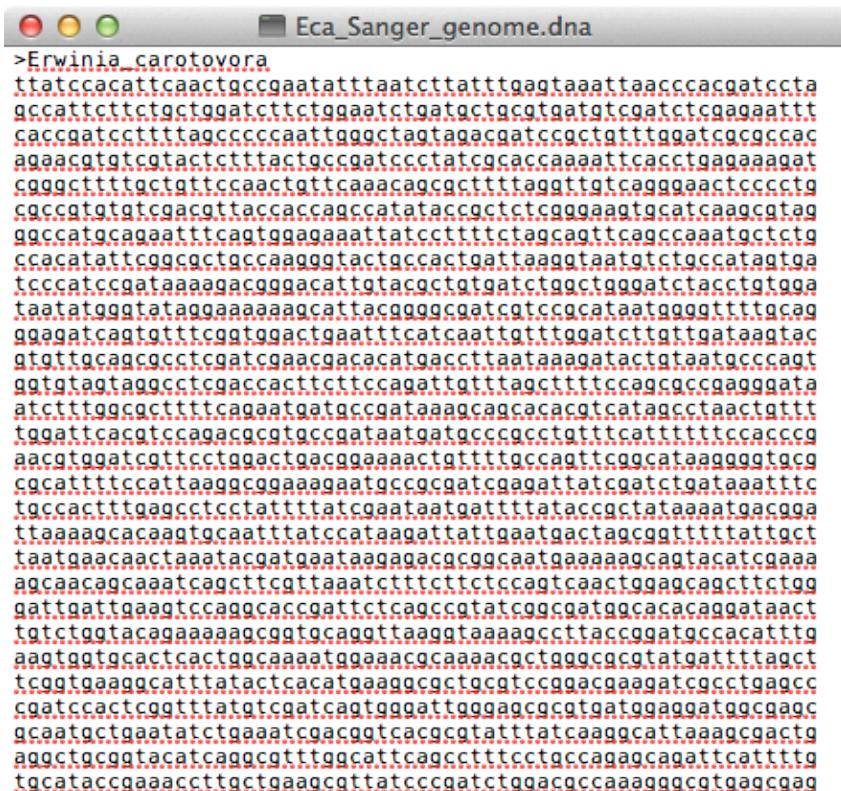


# *Pectobacterium atrosepticum*

- First plant pathogenic enterobacterium to be sequenced
- £250,000 collaboration between SCRI, University of Cambridge, and WT Sanger Institute
- WGS to 10.2X coverage, 106,500 reads
- Repeats and gaps bridged and sequenced directly
- A single complete high quality 5Mbp circular chromosome
- Published 2004, PNAS (32 authors)  
*Bell et al. (2004) Proc Natl Acad Sci USA* [doi:10.1073/pnas.0402424101](https://doi.org/10.1073/pnas.0402424101).

# *Pectobacterium atrosepticum*

- Having a genome sequence is useful, but...



```
>Erwinia_carotovora
ttatccacattcaactgccqaatattaatcttatttgagtaaaattAACCCACGATCTA
GCCATTCTCTGTGGATCTTCTGGAAATCTGATGTCGCGTGTGTCGATCTCGAGAAATT
CACCAGATCCTTTAGCCCCCAATTGGCTAGTAGACGATCCCGTGGATCQCQCAC
AGAACAGTGTGTACTCTTACTGCCATCCATCGCACCAAAATTACCGTGGAGAAAT
CGGGCTTTGCTGTTCCAACCTGTTCAAACAGCGCTTTAGGTTGTCAAGGGAAACTCCCTG
CGCCATGTGTGCGACGTTACCAACCGGCATATAACCGCTCTCGGGAGTGCGATCAAGCGTAG
GGCCATGCAGAATTTCGTTGGAGAAATTATCCTTCTAGCTTCCAGGAAATATGCTCTG
CCACATATTGGCGCTGCCAAGGGTACTGCCACTGTTAAAGGTATGTCGCCATAGTGA
TCCCATCCGATAAAAGACGGGACATTGTAACGTGTGATCTGGCTGGGATCTACCTGTGGA
TAATATGGGTATAGGAAAAAAAGCATTACGGGGCGATCGTCCGCATAATGGGGTTTCGAG
GGAGATCAGTGTTCGGTGGACTGAATTTCATCAATTGTTGATCTTGTGATAAGTAC
GTGTTGAGCQCCTCQATCQAACQACACATGACCTTAATAAAAGATACTGTAATGCCAGT
GGTGTAGTAGGCGCTCGACCACTTCTTCAGATTGTTAGCTTTCCAGCGGCCQAGGGATA
ATCTTGGCGCTTTCAAGAAATGATGCCATAAAAGCAGCAGCACACTGTATAGCCTAACTGTT
TGGGATCAGCTCCAGAGCAGCGTGGCCATAATGATGCCCGCTGTCTCATTTCGCGACCCCG
AACGTTGGATCTTCTGACTGACGGAAACTGTGTTGCCAGGTGCGCATAAAGGGGTGCG
CGCATTTCCATTAAAGGCGGAAAGAATGCCCGATCGAGATTATGATCTGATAAAATTCTG
TGCCACTTGGAGCGCTCCATTGATCGAATAATGATTGTTATACCGCTATAAAATGACGGA
TAAAGGACACAAGTGCATAATTGATCCATAAAAGATTATTGAAATGACTAGCGGTTTATTGCT
TAATGAAACAACATAATACGATGAATAAGGAGACGCCGGAATGAAAAAGCAGTACATCQAAAG
AGCAACAGCAGGAAATCAQCTTCGTTAAATCTTCTCCAGTCAGCTGAGCGCAGCTTCGG
GATTGATTGAAAGTCCAGGCACCGATTCTCAGCGTATCGCGATGGCACACAGGATAACT
GGCTGGTACAGAAAAAGCGGTGCAAGGTTAAGGTTAAAAGGCCCTTACCGGATGCCACATTG
AAGTGGTGCACTCAGCTGCAAAATGGAAAACGCAAAACGCTGGGCGCGTATGATTTTAACT
TGGTGAAGGCAAGTATGATCAGTCAGTGGGGATTGGAGCGCGTGTAGGGAGATGGCGAGC
CGATCCACTCGGTTATGTCGATCGTGGGATTGGAGCGCGTGTAGGGAGATGGCGAGC
GCAATGCTGAATACTGGAAATCGACGGTCAAGCGTATTTATCAAGGGCATTAAAGCGACT
AGGCTGCGTACATCGGGCTTGGCATTCAAGCTTTCCTGCCAGAGCGAGATTCTTGTGG
TGCGATACCGGAAACCTTGTGAAGCGTTATCCCGATCTGGACGCCAAAGGGCGTGAAGCGAG
```

- Relive the annotation with Artemis!
- Eca\_Sanger\_genome.dna
- Eca\_Sanger\_annotation.gbk
- Manual annotation by Sanger Pathogen Sequencing Unit:**

gene ID: FASTA, BLASTN, BLASTX

genecalls: ORPHEUS, GLIMMER

domain ID: PFAM, SIGNALP, TMHMM

metabolism: KEGG

ncRNA: RFAM

**Literature search and comparison**

- ≈6 people, ≈6 months



# A Word About Annotation

- Annotation quality is critical to downstream biology!
- **Annotation is *curation*, not cataloguing!**  
(but may have no budget on a project)  
(not a new situation: Enzyme Nomenclature Committee never funded...)
- Automation from curated data is the **ONLY** game in town,  
but you can't propagate information that does not exist!  
(≈30% of metabolic activity has no known gene associated with it)
- Core, well-curated resources are essential  
(Dirty little fact: biocurators can spend as much time “de-annotating” literature-based annotations as entering new data)



# Another Word About Annotation

- Biocuration (i.e. maintenance of reliable annotation corpus) also requires:
  - software platforms and interfaces (web-based is good...)
  - controlled vocabularies and ontologies (GO, PAMGO, SO, ...)
  - computing and network infrastructure, IT and sysadmin support
- **Community annotation:** “Many potential parents, but no-one wants to look after the children”
- **Key problem:** lack of direct incentives to curate (one’s own) data.
- Solution: ELIXIR and centralised resources, e.g. EBI (Ensembl etc.)



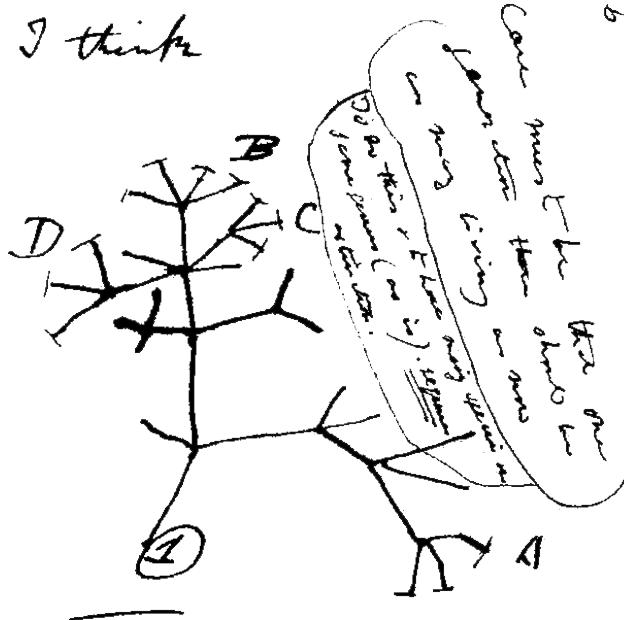
# Comparative Genomics

- Having a single genome is useful, but...  
...the real power of genomics is *comparative genomics*
  
- **The combination of genomic data and comparative and evolutionary biology to address questions of genome structure, evolution and function.**

# Why Comparative Genomics?

- Genomes describe heritable characteristics
- Related organisms share ancestral genomes
- Functional elements encoded in genomes are common to related organisms
- Functional understanding of model systems (*E. coli*, *A. thaliana*, *D. melanogaster*) can be transferred to non-model systems on the basis of genome comparisons
- Genome comparisons can be informative, even for distantly-related organisms

I think

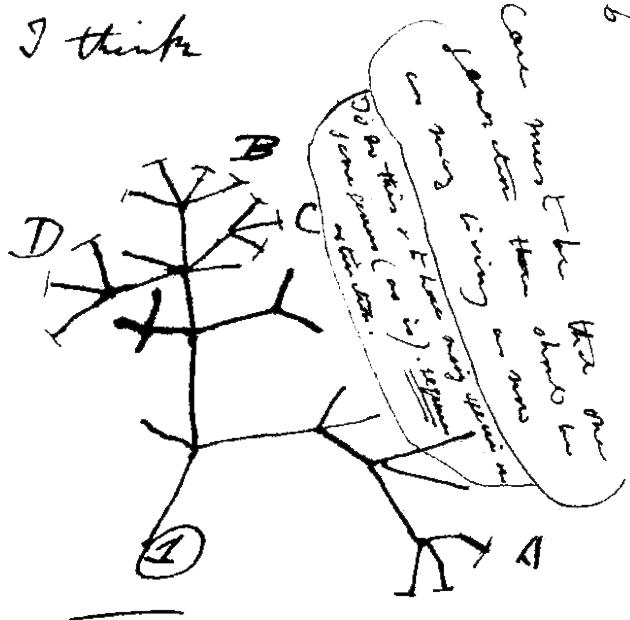


Then between A + B. various  
ways of relation. C + B. the  
finest gradation. B + D  
rather greater distinction.  
Then genome would be  
formed. - binary relation

# Why Comparative Genomics?

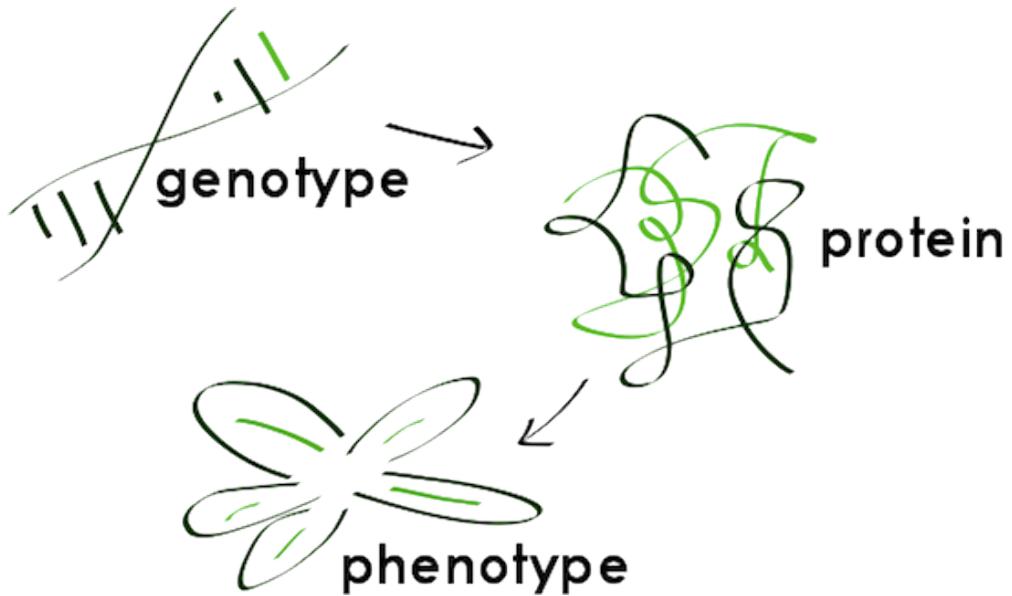
- Genomic differences can underpin phenotypic (morphological, physiological, host range, virulence) differences
- If phenotypes or other organism-level properties are known, comparison of genomes may give mechanistic or functional insight into differences (e.g. environmental persistence).
- Genome comparisons aid identification of functional elements on the genome.
- Studying genomic changes reveals evolutionary processes and constraints.

I think



Then between A + B. various  
ways of relation. C + B. the  
first gradation, B + D  
rather greater distinction.  
Then genera would be  
formed. - binary relation

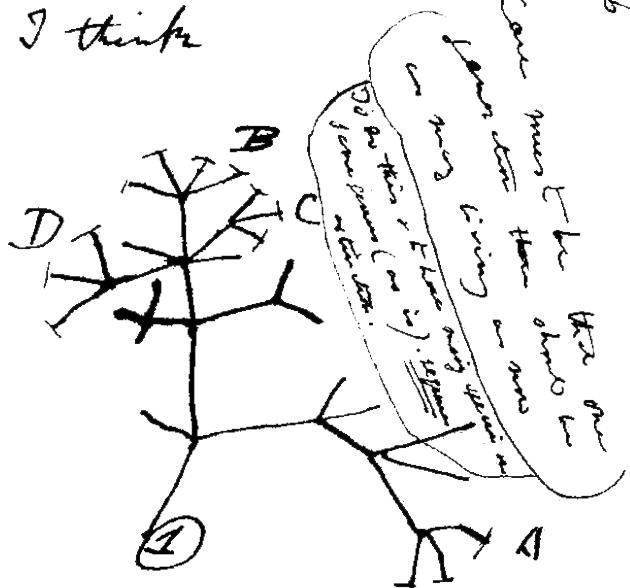
# Caveats of Comparative Genomics



- **BUT: the genome isn't everything**
  - **Context:** epigenetics, tissue differentiation, mesoscale systems, etc.
  - **Phenotypic plasticity:** responses to temperature, stress, environment, etc.

(therefore Systems Biology, of which more later...)

I think



Thus between A + B. various  
taxa & relation. C + B. the  
first gradation, B + D  
rather greater distinction  
Thus genome would be  
formed. - binary relation



# Levels of Genome Comparison

- Bulk Properties
  - chromosome/plasmid counts and sizes,
  - nucleotide content,  $k$ -mers, etc.
- Whole Genome Sequence
  - sequence similarity
  - genomic structure and organisation (rearrangements), etc.
- Genome Features/Functional Components
  - numbers and types of features (genes, ncRNA, regulatory elements, etc.)
  - organisation of features (synteny, operons, regulons, etc.)
  - complements of features
  - selection pressure, etc.



# *Pba* Comparative Genomics Questions

- Just about enough sequenced genomes in 2002/3 to ask:
- What are the gene complement similarities and differences between *Pba* and:
  - other enterobacteria?
  - other plant pathogens?
  - animal pathogens?
  - other environmental bacteria?
- What is genomic evidence for pathogenicity determinants?
- Do pathogenicity determinants show evidence of LGT?
- Are there characteristic features of the *Pba* genome?

# Comparing Gene Features

- Given gene annotations for more than one genome, one can organise and understand relationships in multiple ways
  - Functional similarity (analogy)
  - Evolutionary common origin (homology, orthology, etc.)
  - Evolutionary/functional/family relationships (paralogy)

## DISTINGUISHING HOMOLOGOUS FROM ANALOGOUS PROTEINS

WALTER M. FITCH

### *Abstract*

Fitch, W. M. (Dept. Physiological Chem., U. Wisconsin, Madison 53706) 1970. *Distinguishing homologous from analogous proteins.* Syst. Zool., 19:99–113.—This work provides a means by which it is possible to determine whether two groups of related proteins have a common ancestor or are of independent origin. A set of 16 random

Terms first suggested by Fitch (1970) *Syst. Zool.* [doi:10.2307/2412448](https://doi.org/10.2307/2412448)  
Fitch (2000) *Trends Genet.* [doi:10.1016/S0168-9525\(00\)02005-9](https://doi.org/10.1016/S0168-9525(00)02005-9)



# Attack of the –logues

- Important technical terms describing feature relationships
- **Homologues:** elements that are similar because they share a common ancestor (**NOTE: There are NOT degrees of homology!**)
- **Analogues:** elements that are (functionally?) similar, maybe through convergent evolution rather than common ancestry
- **Orthologues:** homologues that diverged through speciation
- **Paralogues:** homologues that diverged through duplication within the same genome
- (also **co-orthologues**, **xenologues**, etc.)



# Note on “Orthology”

- Frequently abused/misused as a term
- “Orthology” was defined as an evolutionary relationship, but now often bent into service as a functional descriptor
- Strictly defined only for two species or clades!
  - (cf. cluster definitions OrthoMCL, etc.)
- Orthology is not always transitive
  - (A is orthologue of C and B is orthologue of C does not necessarily imply A is an orthologue of B)

Storm & Sonnhammer (2002) *Bioinformatics*. [doi:10.1093/bioinformatics/18.1.92](https://doi.org/10.1093/bioinformatics/18.1.92)  
Jensen (2001) *Genome Biol.* [doi:10.1186/gb-2001-2-8-interactions1002](https://doi.org/10.1186/gb-2001-2-8-interactions1002)  
Fitch (2000) *Trends Genet.* [doi:10.1016/S0168-9525\(00\)02005-9](https://doi.org/10.1016/S0168-9525(00)02005-9)



# Finding “Orthologues”

Or, more pragmatically for hypothesis generation...

- **Finding evolutionary (and/or functional) equivalents of features across two or more organisms' genomes.**

# Why “orthologues”?

- Need to formalise the concept of *corresponding genes* across multiple organisms.
  - Evolutionary
  - Functional? (**“The Ortholog Conjecture”**)
- Many (>35) databases attempt to describe orthology relationships
  - [http://questfororthologs.org/orthology\\_databases](http://questfororthologs.org/orthology_databases)

## List of orthology databases

If you know of any other database, please edit this page directly or contact us.

1. [COGs/TWOGa/KOGs](#)
2. [COGs-COCO-CL](#)
3. [COGs-LOFT](#)
4. [eggNOG](#)
5. [EGO](#)
6. [Ensembl Compara](#)
7. [Gene-Oriented Ortholog Database](#)
8. [GreenPhyDB](#)
9. [HCOP](#)
10. [HomoloGene](#)
11. [HOGENOM](#)
12. [HOVERGEN](#)
13. [HOMOLENS](#)
14. [HOPS](#)
15. [INVHOGEN](#)
16. [InParanoid](#)
17. [KEGG Orthology](#)
18. [MetaPhOrs](#)
19. [MBGD](#)
20. [MDG](#)
21. [OMA](#)
22. [OrthoDB \(OrthoDB on Wikipedia\)](#)
23. [OrthologID](#)
24. [ORTHOLOGUE](#)
25. [OrthoInspector](#)
26. [OrthoMCL](#)
27. [Panther](#)
28. [PhIGs](#)
29. [PHOG](#)
30. [PhylomeDB](#)
31. [PLAZA](#)
32. [P-POD](#)
33. [ProGMap](#)
34. [Proteinortho](#)
35. [RoundUp](#)
36. [TreeFam](#)
37. [YOGY](#)

Dessimoz (2011) *Brief. Bioinf.* [doi:10.1093/bib/bbr057](https://doi.org/10.1093/bib/bbr057)

Chen & Zhang (2012) *PLoS Comp. Biol.* [doi:10.1371/journal.pcbi.1002784](https://doi.org/10.1371/journal.pcbi.1002784)

Altenhoff *et al.* (2012) *PLoS Comp. Biol.* [doi:10.1371/journal.pcbi.1002514](https://doi.org/10.1371/journal.pcbi.1002514)



# How to find “orthologues”?

- Many published methods and databases:
  - Pairwise between two genomes:
    - [RBBH](#) (aka BBH, RBH, etc.), [RSD](#), [InParanoid](#), [RoundUp](#)
  - Multi-genome:
    - Graph-based: [COG](#), [eggNOG](#), [OrthoDB](#), [OrthoMCL](#), [OMA](#), [MultiParanoid](#)
    - Tree-based: [TreeFam](#), [Ensembl Compara](#), [PhylomeDB](#), [LOFT](#)
- Methods may apply different - or *refined* - definitions of orthology, paralogy, etc.

Salichos *et al.* (2011) *PLoS One*. [doi:10.1371/journal.pone.0018755](https://doi.org/10.1371/journal.pone.0018755)

Trachana *et al.* (2011) *Bioessays* [doi:10.1002/bies.201100062](https://doi.org/10.1002/bies.201100062)

Kristensen *et al.* (2011) *Brief. Bioinf.* [doi:10.1093/bib/bbr030](https://doi.org/10.1093/bib/bbr030)



# Evaluating “Orthologue” Prediction

- Works the same way as for all prediction tools
  1. Define a “validation set” (gold standard), unseen by the prediction tool
  2. Make predictions with the tool
  3. Evaluate confusion matrix and performance statistics
    - Sensitivity
    - Specificity
    - Accuracy

Standard:	+ve	-ve
Predict +ve	TP	FP
Predict -ve	FN	TN

False positive rate	$FP/(FP+TN)$
False negative rate	$FN/(TP+FN)$
Sensitivity	$TP/(TP+FN)$
Specificity	$TN/(FP+TN)$
False discovery rate (FDR)	$FP/(FP+TP)$
Accuracy	$(TP+TN)/(TP+TN+FP+FN)$



# “Orthologue” Prediction Performance

- Performance varies by choice of method and interpretation of “orthology”
- Biggest influence is often annotation quality
- Relative performance varies with benchmark choice
- (clustering) RBH outperforms the more complex algorithms under many circumstances

Wolf *et al.* (2012) *Genome Biol. Evol.* [doi:10.1093/gbe/evs100](https://doi.org/10.1093/gbe/evs100)

Salichos *et al.* (2011) *PLoS One.* [doi:10.1371/journal.pone.0018755](https://doi.org/10.1371/journal.pone.0018755)

Altenhoff & Dessimoz (2009) *PLoS Comp. Biol.* [doi:10.1371/journal.pcbi.1000262](https://doi.org/10.1371/journal.pcbi.1000262)

Trachana *et al.* (2011) *Bioessays.* [doi:10.1002/bies.201100062](https://doi.org/10.1002/bies.201100062)

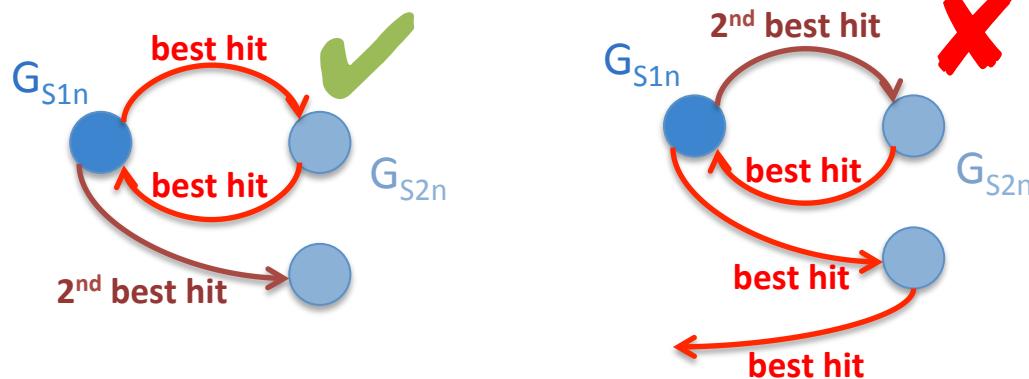


# Back To Our *Pba* Questions

- What are the gene complement similarities and differences between *Pba* and:
  - other enterobacteria?
  - other plant pathogens?
  - animal pathogens?
  - other environmental bacteria?
- What is genomic evidence for pathogenicity determinants?
- Do pathogenicity determinants show evidence of LGT?
- Are there characteristic features of the *Pba* genome?

# Pba Reciprocal Best Hits (RBH)

- RBH not necessarily orthologues, but “equivalent” features
- Compared to 64 completely sequenced bacterial genomes (most of those available at the time, also limited by computing power)
- Limited matches to:
  - >30% identity (above “twilight zone”)
  - >80% coverage (exclude domain-only matches)





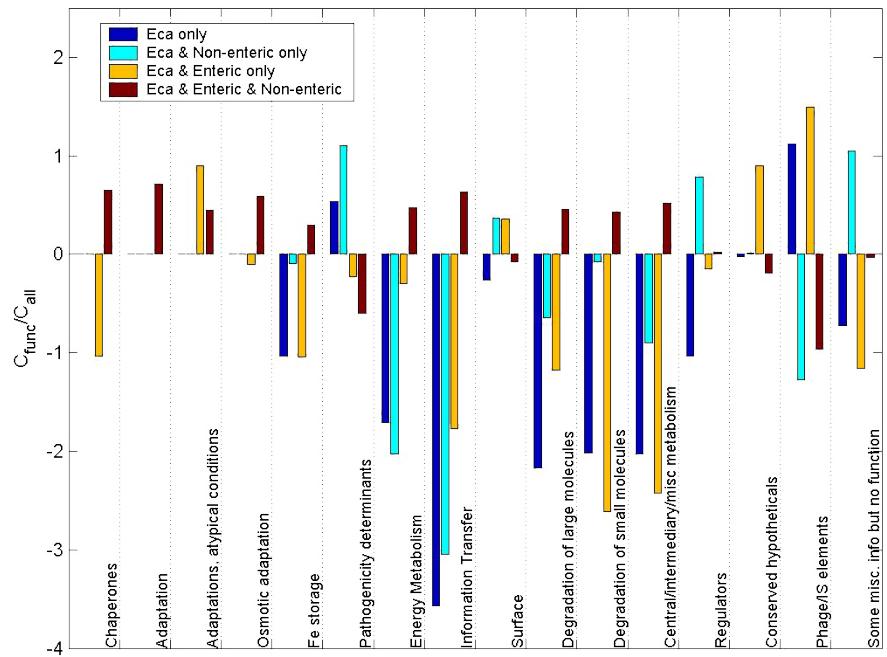
# *Pba* Reciprocal Best Hits (RBH)

- **Advantages:**
  - quick
  - easy
  - performs surprisingly well
- **Disadvantages:**
  - 1:1 matches miss paralogues
  - not good at identifying gene families or \*-to-many relationships without more detailed analysis.
  - no strong theoretical/phylogenetic basis.



# Functional Classification

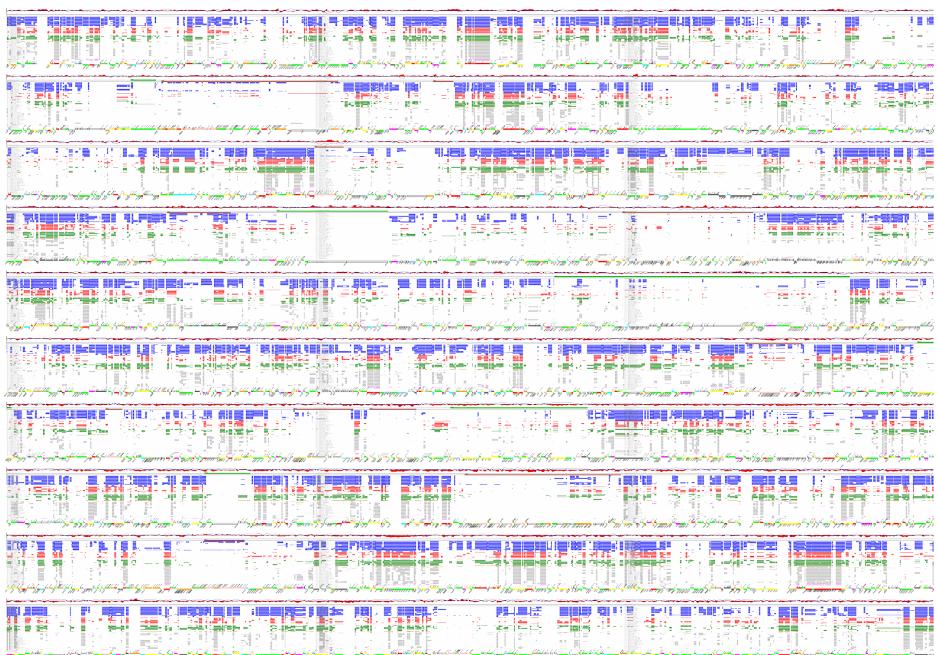
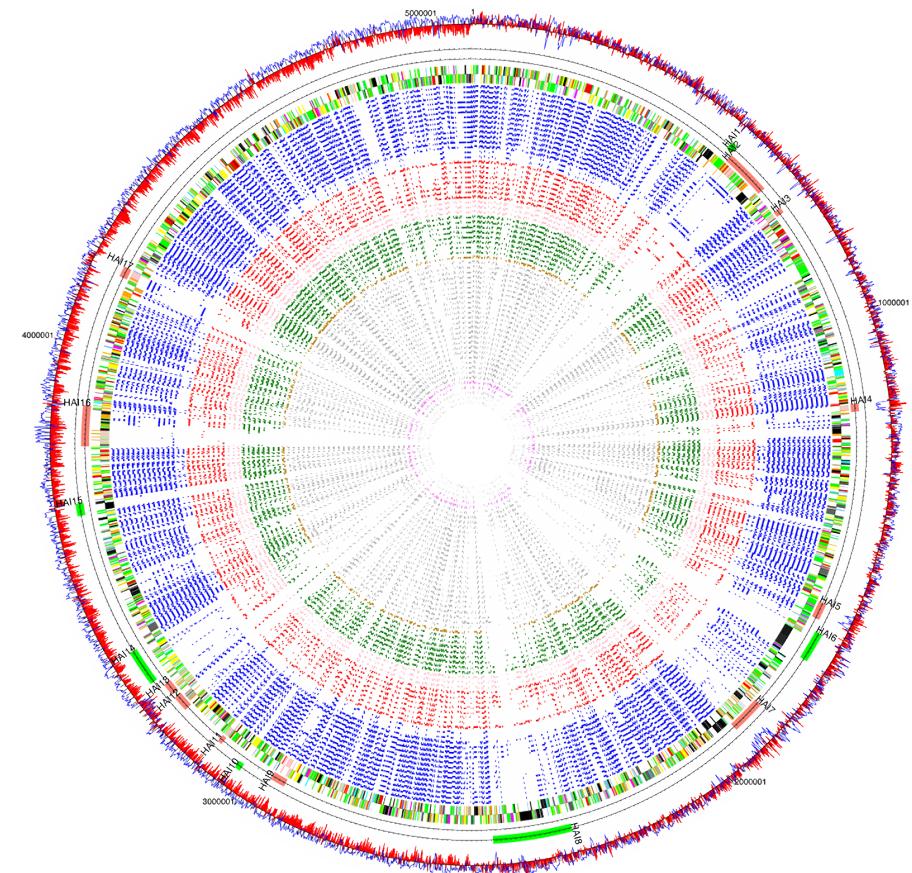
- Compare frequency of functional classifications in *Pba* and other organism classes (**NOTE:** annotation-dependent)
- Test statistically for enrichment/overrepresentation
- Many, many methods & programs for this (68 tools summarised in 2009 review)
- *Pba*-specific:
  - Pathogenicity determinants
  - Phage/IS elements
- *Pba* & non-enteric (LGT?):
  - Pathogenicity determinants
  - Surface proteins
  - Regulatory proteins





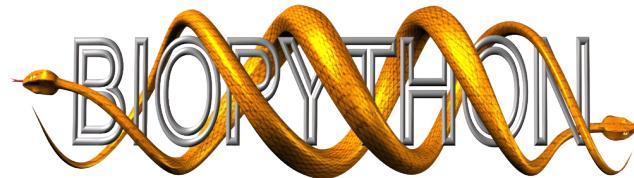
# Pba RBH Visualisation

- 114,087 RBH between *Pba* and target bacteria  
too much to go through by hand/eye
- Visualisation tools were essential

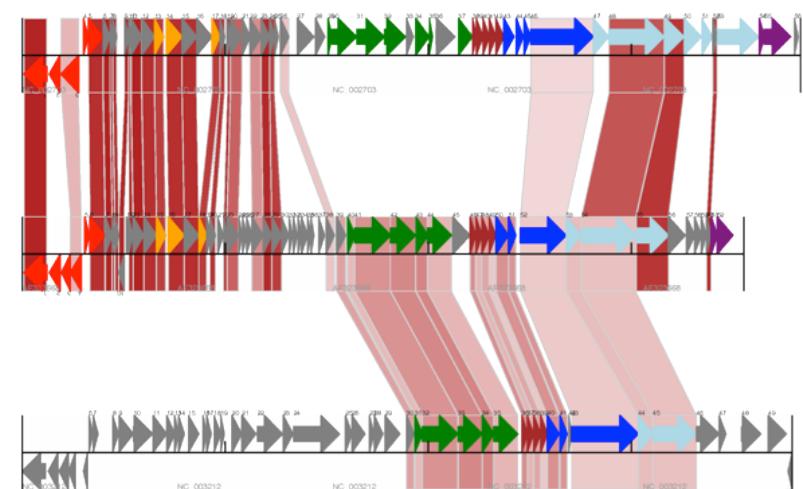
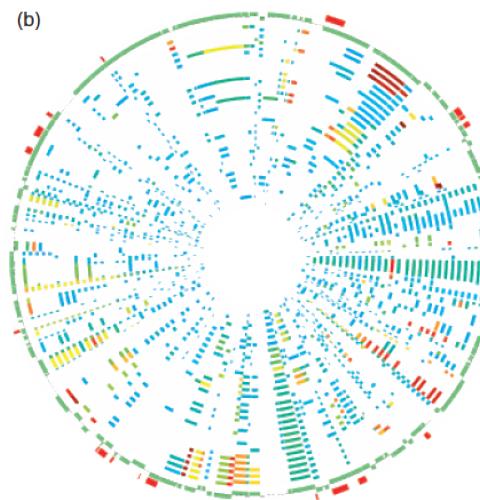
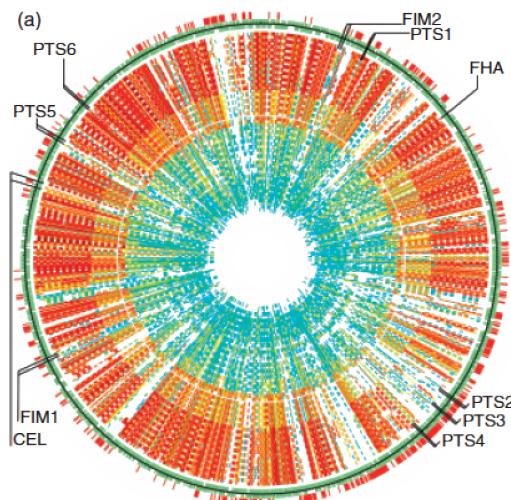




# GenomeDiagram



- Comparative genomics visualisation package
- Developed in 2003 for *Pba* sequencing, later incorporated into Biopython

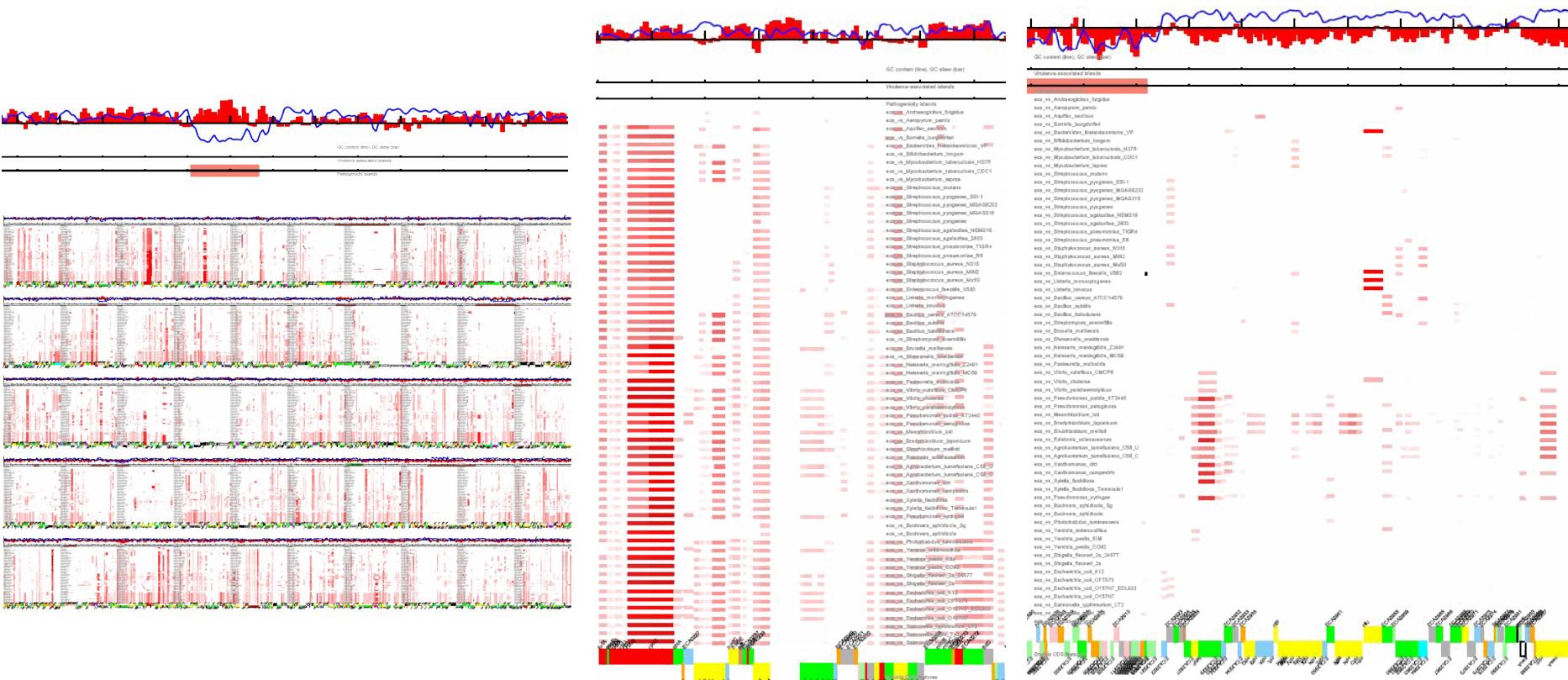


<http://www.biopython.org>

Pritchard et al. (2006) *Bioinformatics* doi:10.1093/bioinformatics/btk021.

# Visualising LGT in *Pba*

- Looked for CDS where RBH sequence similarity doesn't follow the taxonomic species tree, and GC% anomalies
- Colour sequence similarity: **more red = more similar**



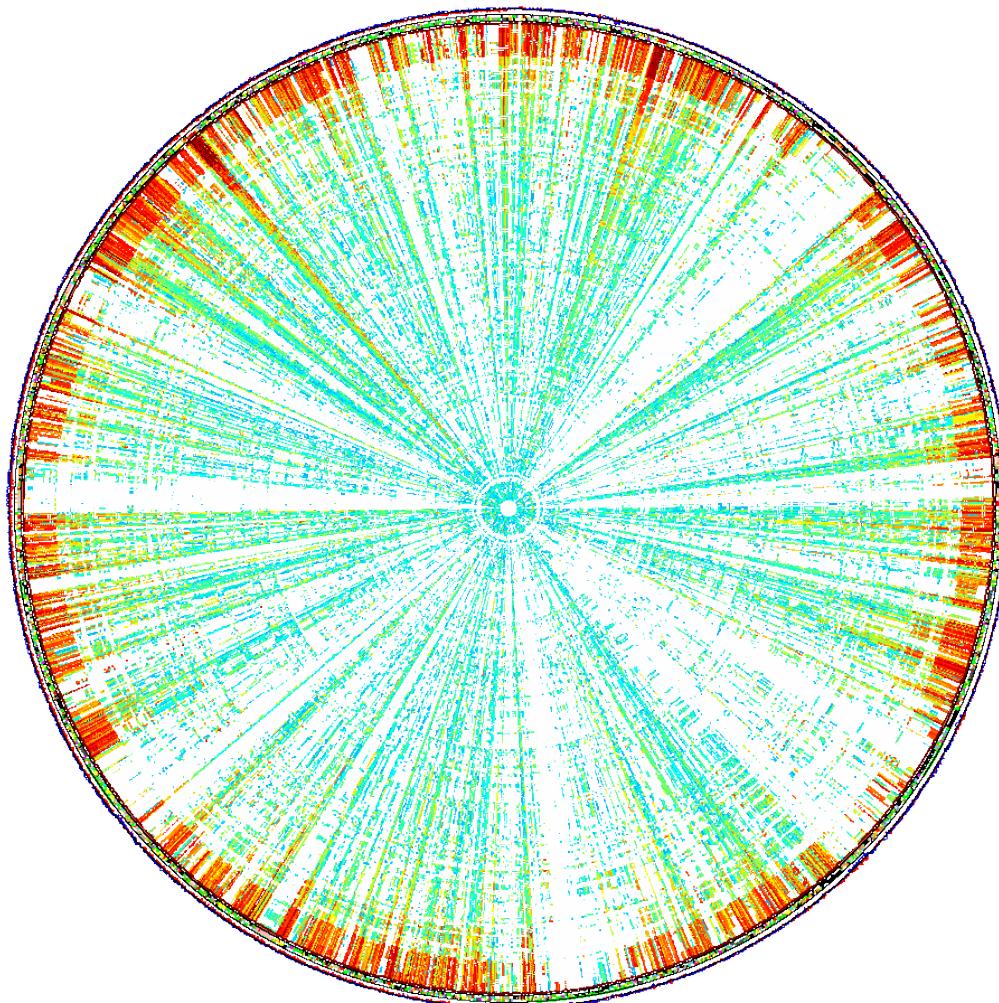


- Initially, 11 putative horizontally-acquired islands (HAIs) identified by atypical G+C, proximity of tRNA, etc.
- Six further putative HAIs identified on the basis of RBH (may be more ancient acquisitions)
- **HAI8:** Type III Secretion System, *Pseudomonas syringae*
- **HAI2:** Phytotoxin (*cfa*) synthesis, *P. syringae*
- **HAI15:** Adherence/agglutination, *Vibrio* spp.
- **HAI14:** Nitrogen fixation, *Bradyrhizobium* spp.



# *Pba* LGT: Plant-associated bacteria

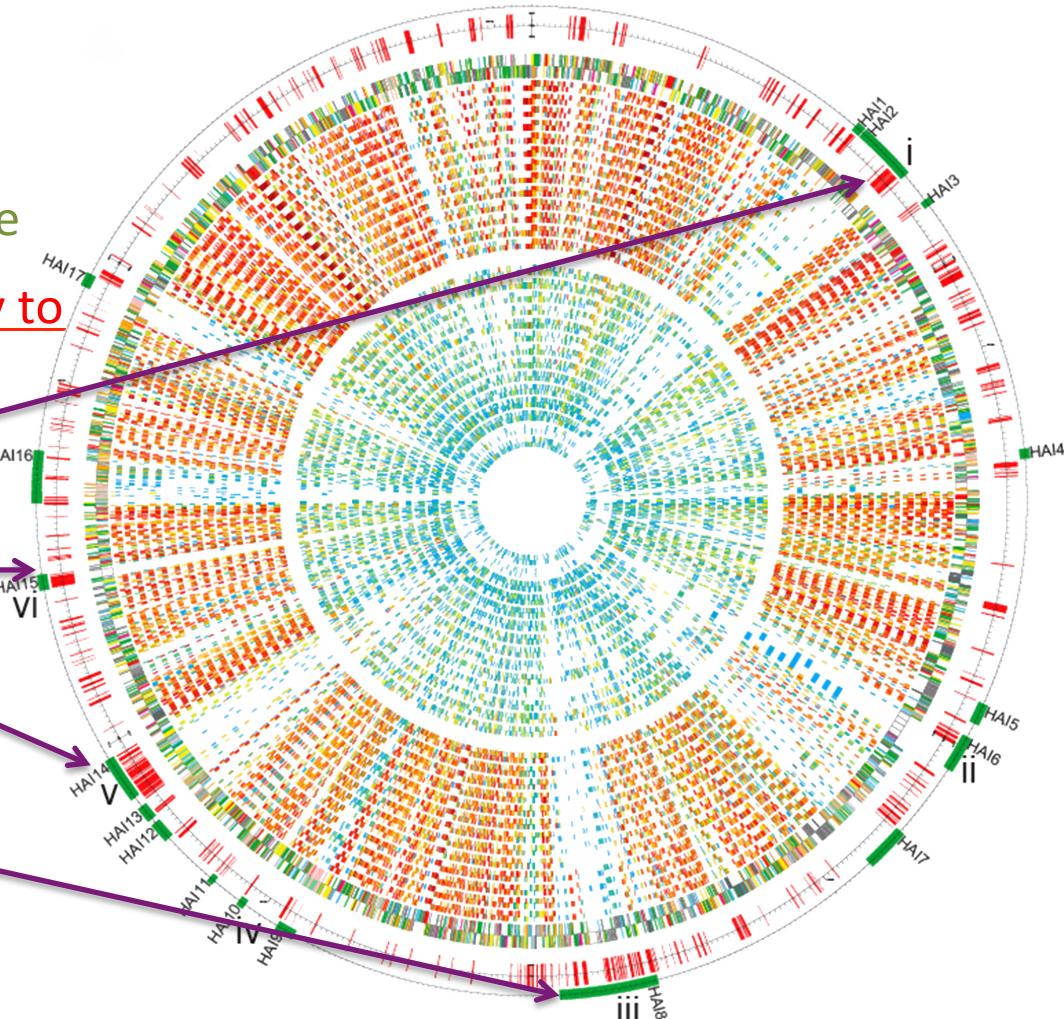
- 2006: Two years later, far more genomes (400) for comparison
- Larger selection of plant- and animal-associated bacteria



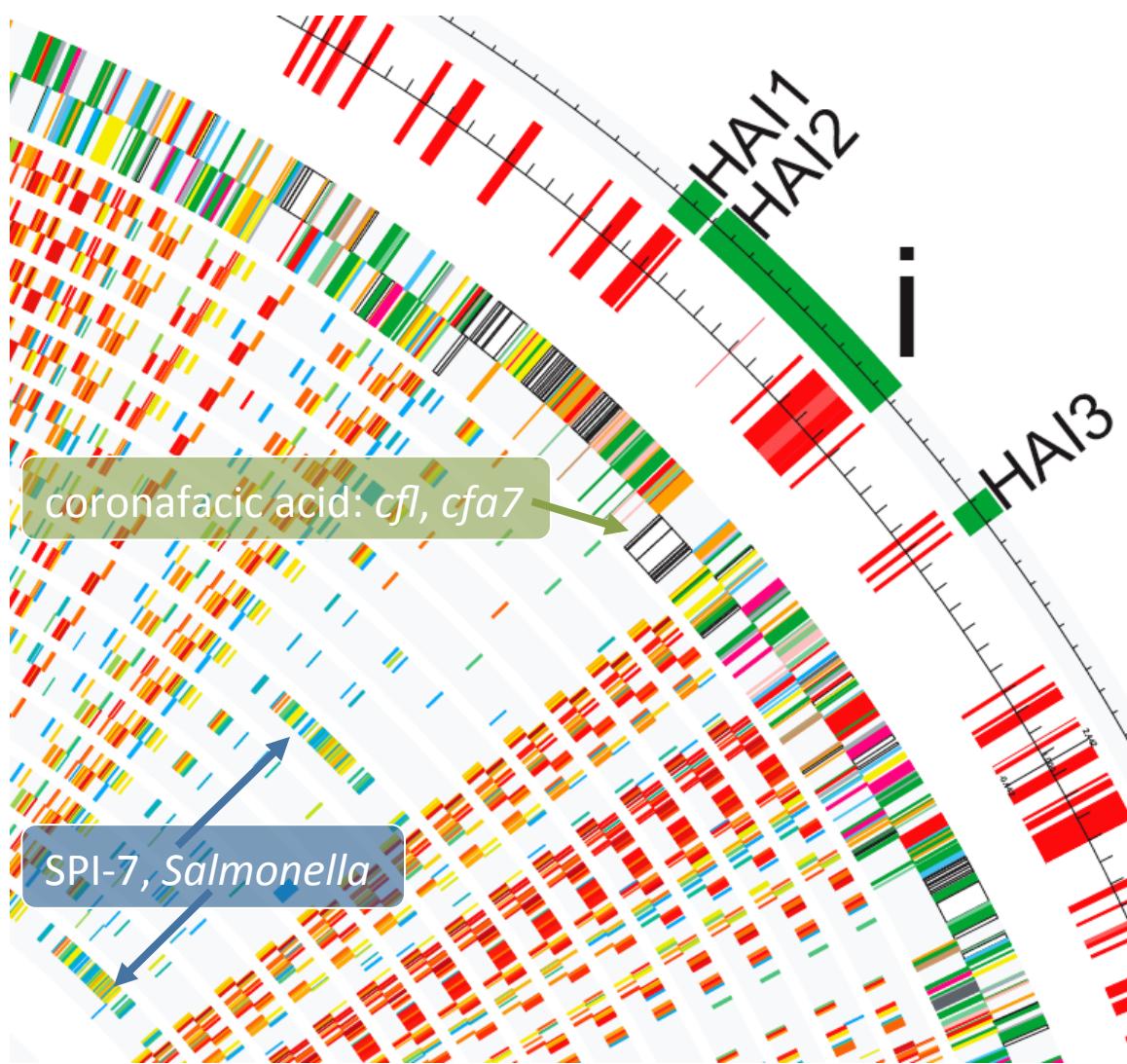


# Pba LGT: Plant-associated bacteria

- Comparison against plant- (13) and animal-associated (14) bacteria
- Plant-associated in centre
- Animal-associated on outside
- Red marks: greater similarity to plant-associated bacteria
- HAI2: Phytotoxin
- HAI15: Adherence
- HAI14: Nitrogen fixation
- HAI8: T3SS



# SPI-7 and Coronafacic Acid



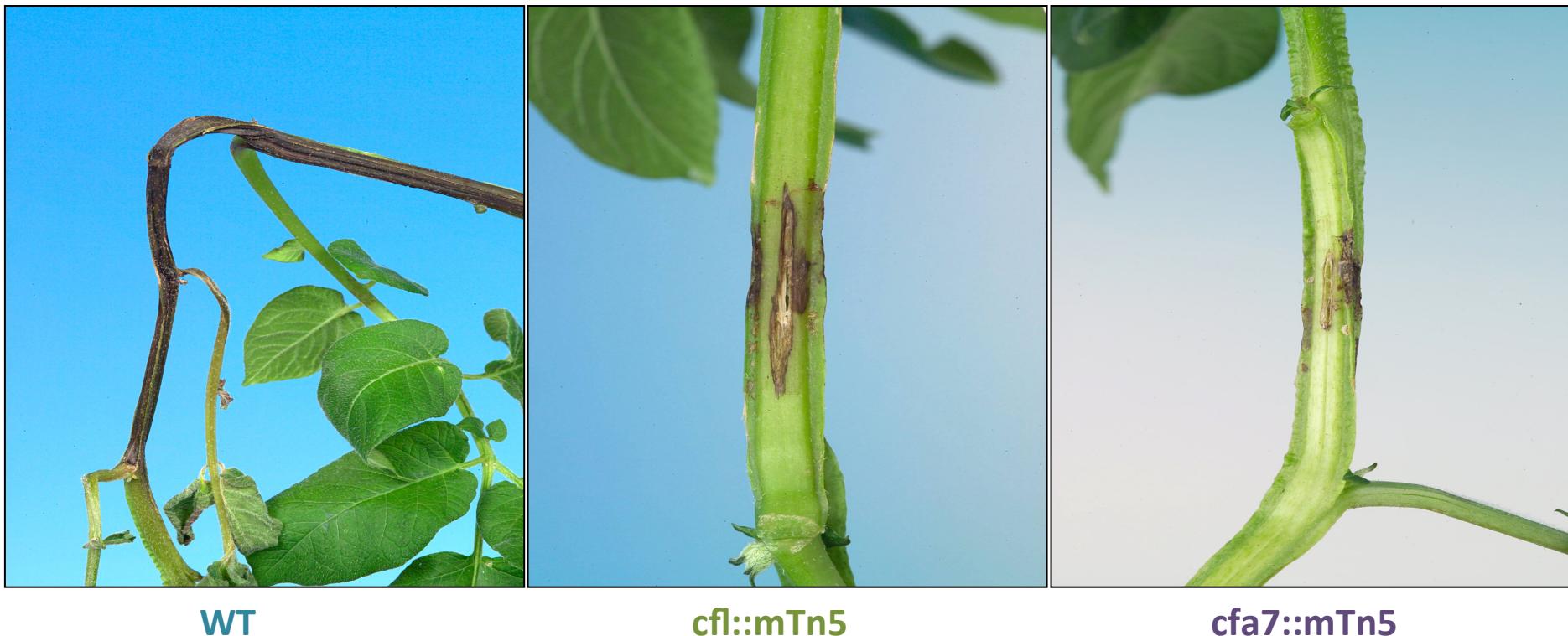
Coronatine (*P. syringae*) interferes with jasmonate responses in host, as a jasmonate mimic

Coronafacic acid –  
*Pseudomonas syringae* phytotoxin precursor (coronatine)  
- payload

SPI-7 -  
*Salmonella Typhi*  
Pathogenicity island  
- delivery system

# Coronafacic Acid Synthesis

- Important for virulence? Testable hypothesis.
- Present in *Pba*, some *Pcc*, virtually no *Dickeya*
- *cfl* and *cfa7* knockouts show reduced *Pba* virulence
  - 200cfu inoculation, measure 17dpi





# Conclusions I

- **Comparative genomics is powerful!**
- Good annotation is essential to understanding genomic data
- Useful, falsifiable hypotheses can be generated from inspection of genomic data
- **Visualisation is extremely useful for large genomic datasets**
- Determining orthology/functional equivalence on large scale data is not yet a solved problem

# 2009: *Phytophthora infestans*

Vol 461 | 17 September 2009 | doi:10.1038/nature08358

nature

## LETTERS

### Genome sequence and analysis of the Irish potato famine pathogen *Phytophthora infestans*

Brian J. Haas<sup>1\*</sup>, Sophien Kamoun<sup>2,3\*</sup>, Michael C. Zody<sup>1,4</sup>, Rays H. Y. Jiang<sup>1,5</sup>, Robert E. Handsaker<sup>1</sup>, Liliana M. Cano<sup>2</sup>, Manfred Grabherr<sup>1</sup>, Chinnappa D. Kodira<sup>1†</sup>, Sylvain Raffaele<sup>2</sup>, Trudy Torto-Alalibo<sup>1†</sup>, Tolga O. Bozkurt<sup>2</sup>, Audrey M. V. Ah-Fong<sup>6</sup>, Lucia Alvarado<sup>1</sup>, Vicki L. Anderson<sup>7</sup>, Miles R. Armstrong<sup>8</sup>, Anna Avrova<sup>8</sup>, Laura Baxter<sup>9</sup>, Jim Beynon<sup>9</sup>, Petra C. Boevink<sup>8</sup>, Stephanie R. Bollmann<sup>10</sup>, Jorunn I.B. Bos<sup>3</sup>, Vincent Bulone<sup>11</sup>, Guohong Cai<sup>12</sup>, Cahid Cakir<sup>3</sup>, James C. Carrington<sup>13</sup>, Megan Chawner<sup>14</sup>, Lucio Conti<sup>15</sup>, Stefano Costanzo<sup>16</sup>, Richard Ewan<sup>15</sup>, Noah Fahlgren<sup>13</sup>, Michael A. Fischbach<sup>17</sup>, Johanna Fugelstad<sup>11</sup>, Eleanor M. Gilroy<sup>8</sup>, Sante Gnerre<sup>1</sup>, Pamela J. Green<sup>18</sup>, Laura J. Grenville-Briggs<sup>7</sup>, John Griffith<sup>14</sup>, Niklaus J. Grünwald<sup>10</sup>, Karoly Horn<sup>14</sup>, Neil R. Horner<sup>7</sup>, Chia-Hui Hu<sup>19</sup>, Edgar Huitema<sup>3</sup>, Dong-Hoon Jeong<sup>18</sup>, Alexandra M. E. Jones<sup>2</sup>, Jonathan D. G. Jones<sup>2</sup>, Richard W. Jones<sup>20</sup>, Elinor K. Karlsson<sup>1</sup>, Sridhara G. Kunjeti<sup>21</sup>, Kurt Lamour<sup>22</sup>, Zhenyu Liu<sup>3</sup>, LiJun Ma<sup>1</sup>, Daniel MacLean<sup>2</sup>, Marcus C. Chibucos<sup>23</sup>, Hayes McDonald<sup>24</sup>, Jessica McWalters<sup>14</sup>, Harold J. G. Meijer<sup>5</sup>, William Morgan<sup>25</sup>, Paul F. Morris<sup>26</sup>, Carol A. Munro<sup>27</sup>, Keith O'Neill<sup>1†</sup>, Manuel Ospina-Giraldo<sup>14</sup>, Andrés Pinzón<sup>28</sup>, Leighton Pritchard<sup>8</sup>, Bernard Ramsahoye<sup>29</sup>, Qinghu Ren<sup>30</sup>, Silvia Restrepo<sup>28</sup>, Sourav Roy<sup>6</sup>, Ari Sadanandom<sup>15</sup>, Alon Savidor<sup>31</sup>, Sebastian Schornack<sup>2</sup>, David C. Schwartz<sup>32</sup>, Ulrike D. Schumann<sup>7</sup>, Ben Schwessinger<sup>2</sup>, Lauren Seyer<sup>14</sup>, Ted Sharpe<sup>1</sup>, Cristina Silvar<sup>2</sup>, Jing Song<sup>3</sup>, David J. Studholme<sup>2</sup>, Sean Sykes<sup>1</sup>, Marco Thines<sup>2,33</sup>, Peter J. I. van de Vondervoort<sup>5</sup>, Vipaporn Phuntumart<sup>26</sup>, Stephan Wawra<sup>7</sup>, Rob Weide<sup>5</sup>, Joe Win<sup>2</sup>, Carolyn Young<sup>3</sup>, Shiguo Zhou<sup>32</sup>, William Fry<sup>12</sup>, Blake C. Meyers<sup>18</sup>, Pieter van West<sup>7</sup>, Jean Ristaino<sup>19</sup>, Francine Govers<sup>5</sup>, Paul R. J. Birch<sup>34</sup>, Stephen C. Whisson<sup>8</sup>, Howard S. Judelson<sup>6</sup> & Chad Nusbaum<sup>1</sup>



Haas *et al.* (2009) *Nature* **461**: 393-398. doi:10.1038/nature08358.



# *Phytophthora infestans*

- The most destructive pathogen of potato (\$6.7bn/yr)
- Causes potato late blight
- Irish Potato Famine (1850s)
- Oomycete model organism, hemibiotrophic
- Adapts rapidly to overcome control measures and bred resistance





# *Phytophthora infestans*

- Major international collaboration, 34 institutions
- WGS at Broad Institute of MIT and Harvard, to 9X coverage
- Four library types, due to highly repetitive genome

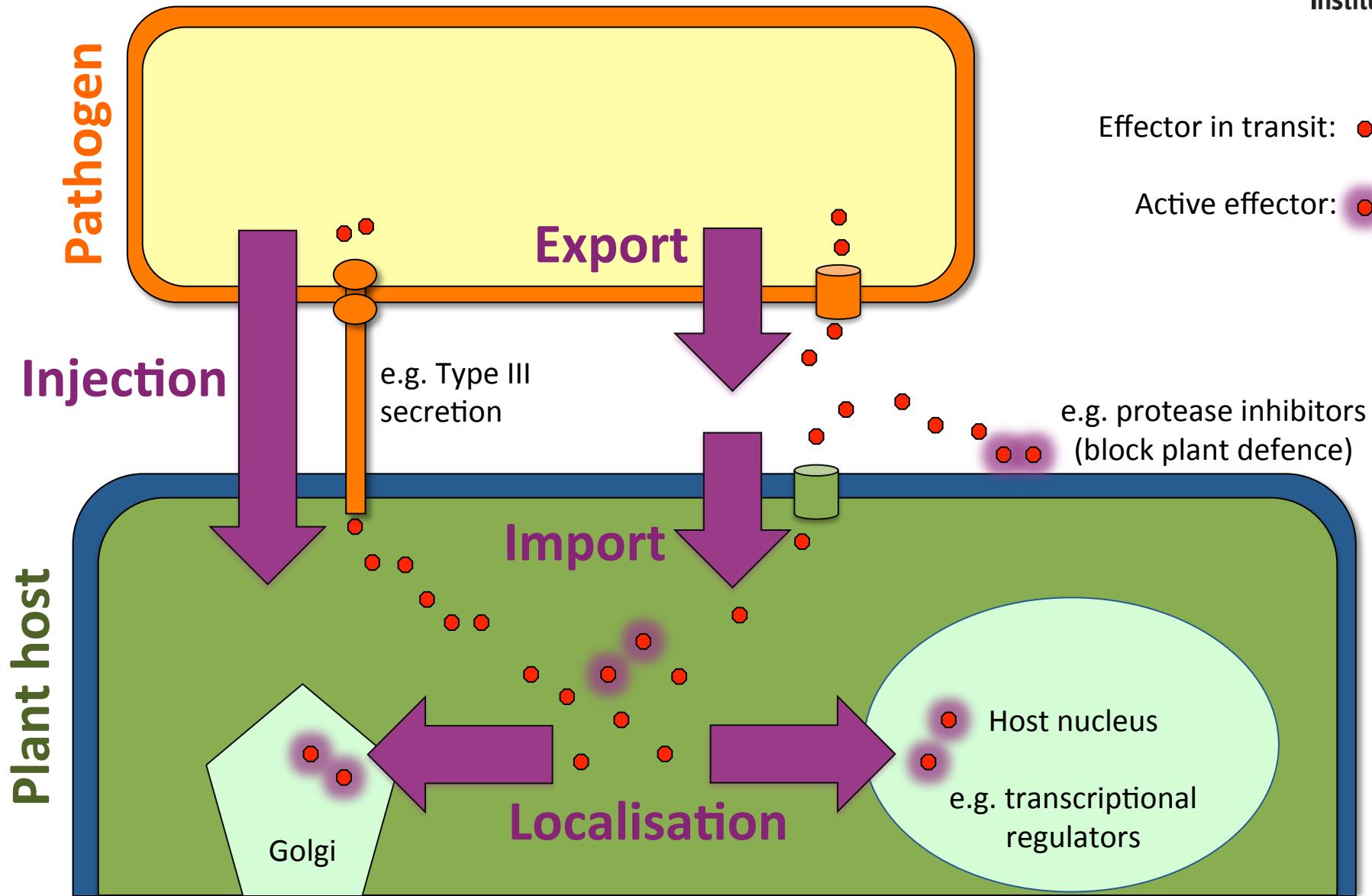
4kb plasmids, 10kb plasmids, 40kb fosmids, 80kb BACs

- Annotation mostly automated; some protein families received manual review by the community

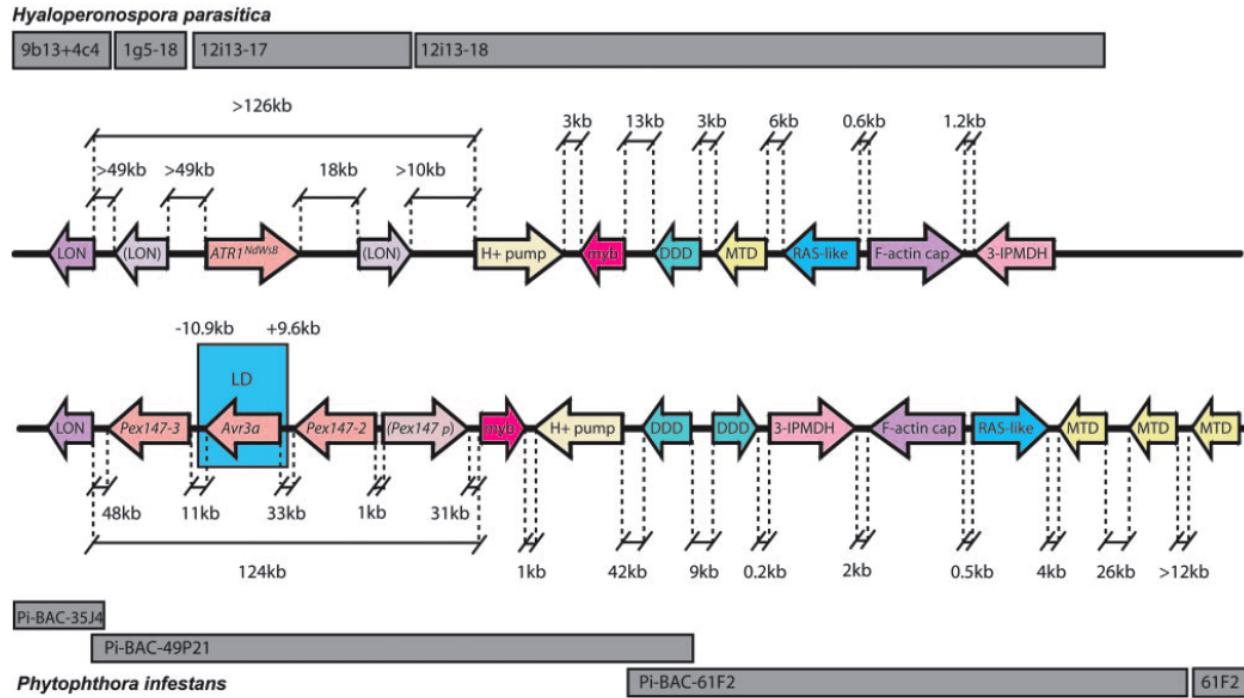
CRN, NPP1-like, elicitin-like, small secreted cysteine-rich, RxLR and transporter

- Published 2009, Nature (96 authors)  
*Haas et al. (2009) Nature* **461**: 393-398. [doi:10.1038/nature08358](https://doi.org/10.1038/nature08358).
- Public data: 18,288 contigs (4,921 supercontigs); 250Mbp (possibly triploid)  
[http://www.broadinstitute.org/annotation/genome/phytophthora\\_infestans](http://www.broadinstitute.org/annotation/genome/phytophthora_infestans)

# My Interest: Effectors



# *P. infestans* Avr3a locus



Contigged, annotated  
*P. infestans* and  
*H. peronospora* BAC  
fragments

Identified microsynteny  
surrounding the *Avr3a*  
and *ATR1* avirulence  
loci

*Avr3a*, *ATR1* belong to  
RXLR family

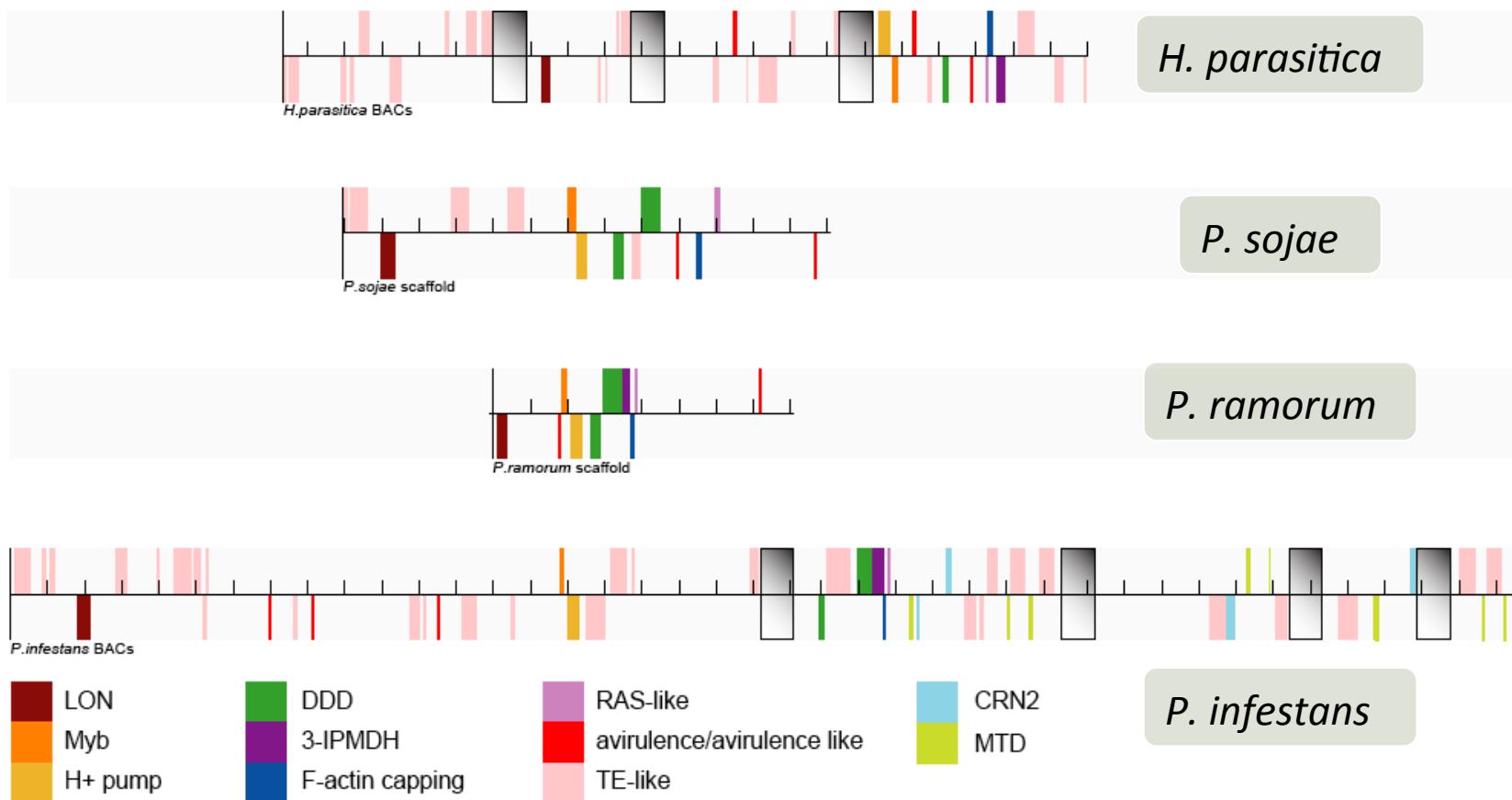
- RXLR N-terminal motif observed in a number of *P. infestans* effectors
- Also observed in *Plasmodium falciparum*, and thought associated with translocation
- **Necessary for translocation in *P. infestans***

Armstrong *et al.* (2005) Proc. Natl. Acad. Sci. USA [doi:10.1073/pnas.0500113102](https://doi.org/10.1073/pnas.0500113102)

Whisson *et al.* (2007) Nature [doi:10.1038/nature06203](https://doi.org/10.1038/nature06203)

# Genome expansion in oomycetes

- Early evidence from BAC assemblies of syntenous loci

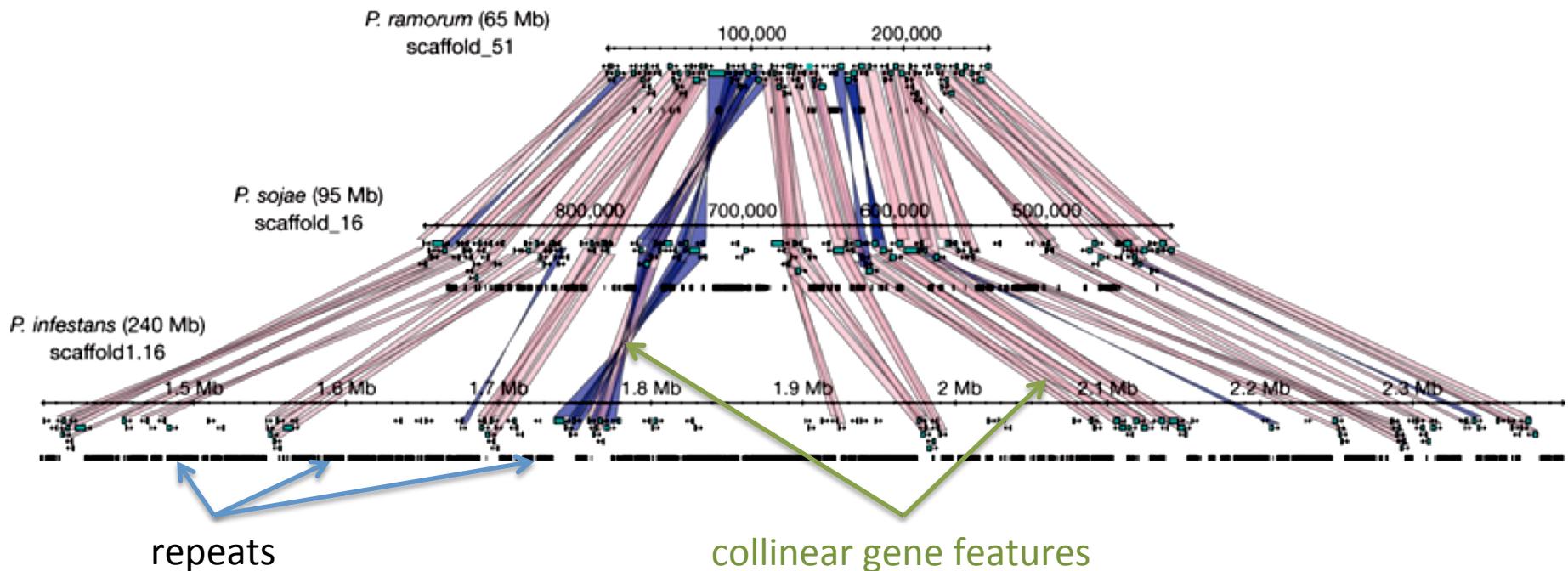


Armstrong *et al.* (2005) Proc. Natl. Acad. Sci. USA [doi:10.1073/pnas.0500113102](https://doi.org/10.1073/pnas.0500113102)

Whisson *et al.* (2007) Nature [doi:10.1038/nature06203](https://doi.org/10.1038/nature06203)

# Repeat-driven Genome Expansion

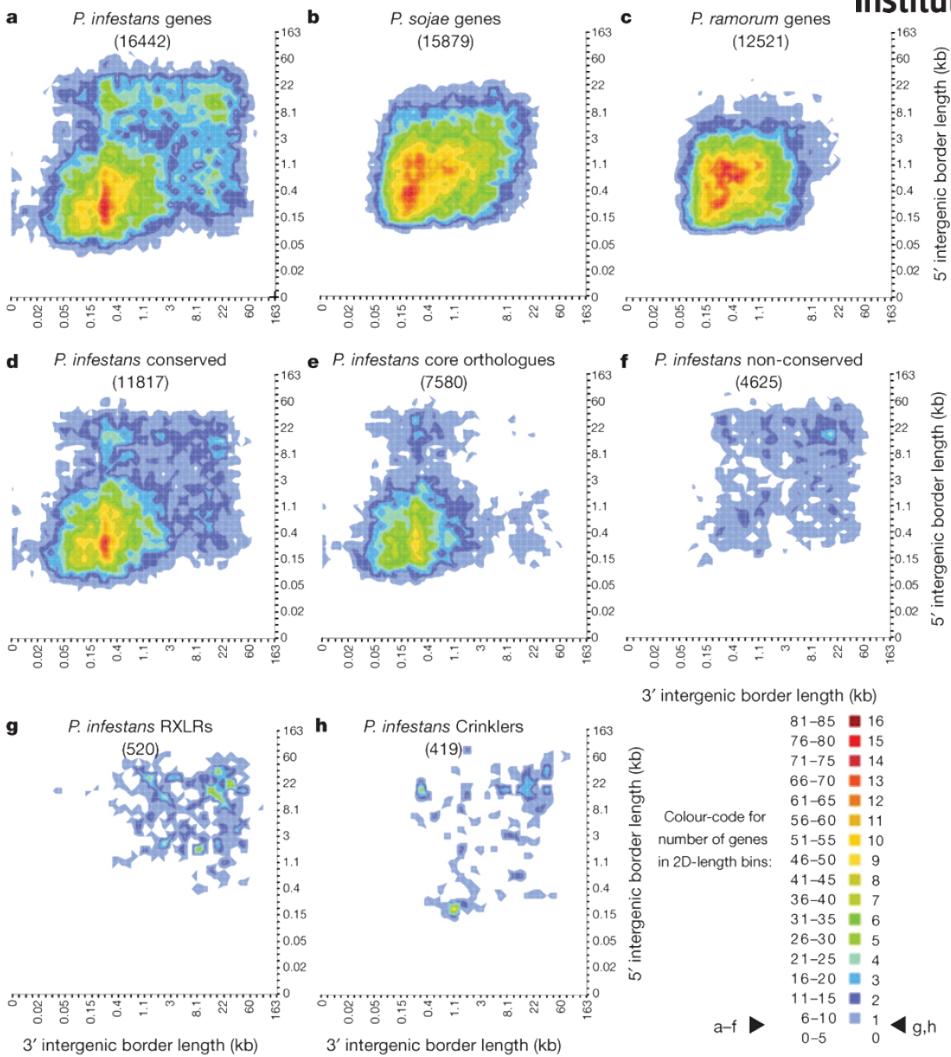
- *P. infestans* genome >70% repeat regions (mostly TEs)
  - “Our findings suggest a two-speed genome...”  
<https://www.broadinstitute.org/news/1328>
  - High repeat/TE density results in ‘experimentation’





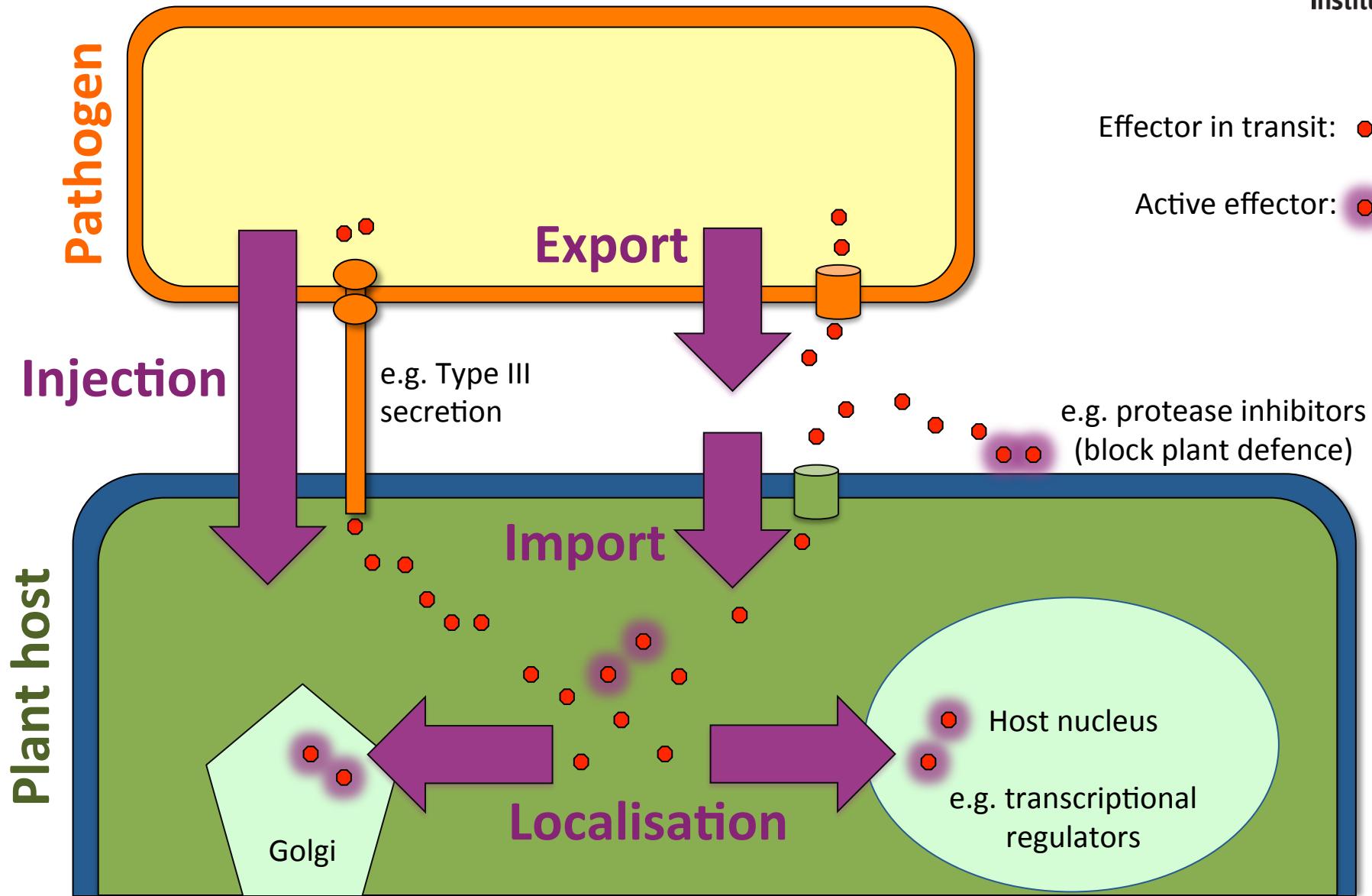
# *P. infestans* two-speed genome

- Diverse population of transposons drives genome expansion specifically in *P. infestans*
- Pathogen effectors disproportionately located in expansion regions
- Expansion drives effector diversity, drives adaptation to overcome novel resistance?



(see pi\_two\_speed.ipynb)

# Effector Properties: action and localisation





# What is an effector?

- Molecule produced by pathogen that (directly?) modifies host molecular/biochemical ‘behaviour’, e.g.
  - Inhibits enzyme action (*Cladosporium fulvum* AVR2, AVR4; *Phytophthora infestans* EPIC1, EPIC2B; *P. sojae* glucanase inhibitors)
  - Cleaves protein target (*Pseudomonas syringae* AvrRpt2)
  - (De-)phosphorylates protein target (*Pseudomonas syringae* AvrRPM1, AvrB)
  - Additional component in/retargeting host system, e.g. E3 ligase activity (*P. syringae* AvrPtoB; *P. infestans* Avr3a)
  - Regulatory control (*Xanthomonas campestris* AvrBs3, TAL effectors)



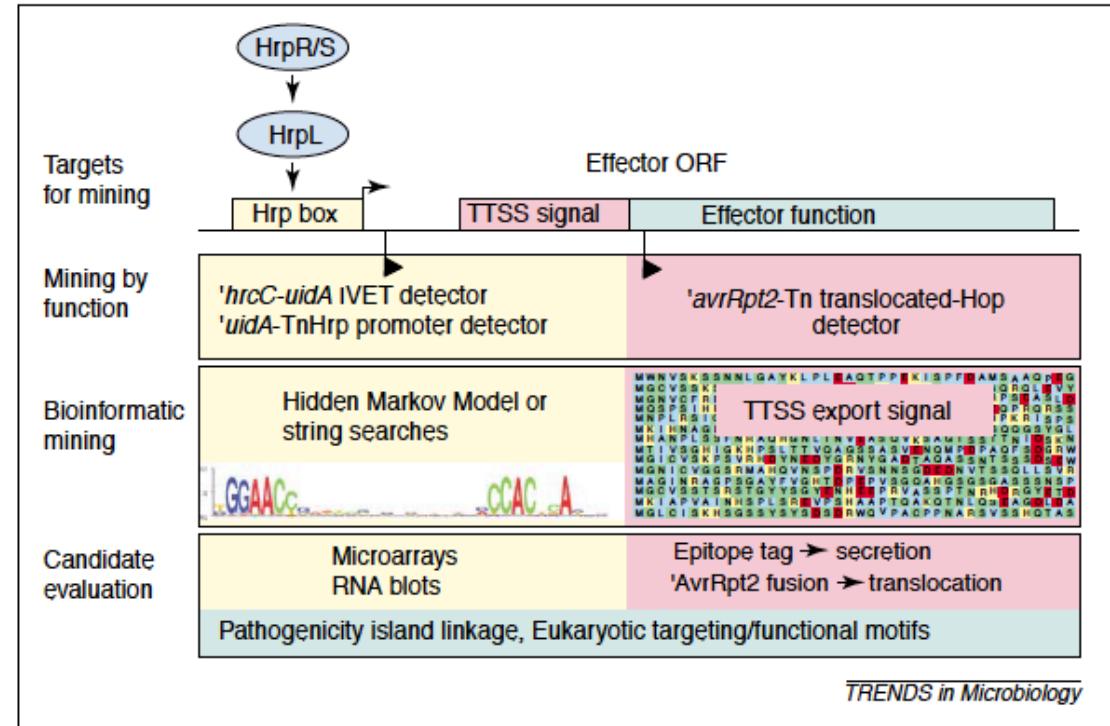
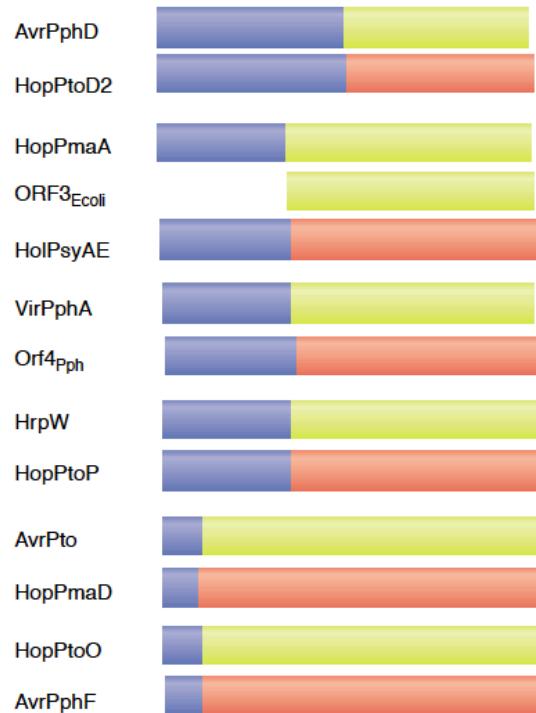
# What is an effector?

- No unifying biochemical mechanism; may act inside or outwith host cell, many functions and final locations
- **No single test** for ‘candidate effectors’
  - Tests are for protein family membership and/or evidence of ‘effector-like behaviour’
  - A general sequence classification problem (functional annotation)
  - Many possible bioinformatic/computational approaches
  - No big red button



# Characteristics of known effectors

- Modularity
  - Delivery: localisation/translocation domain(s)
  - Activity: functional/interaction domain(s)

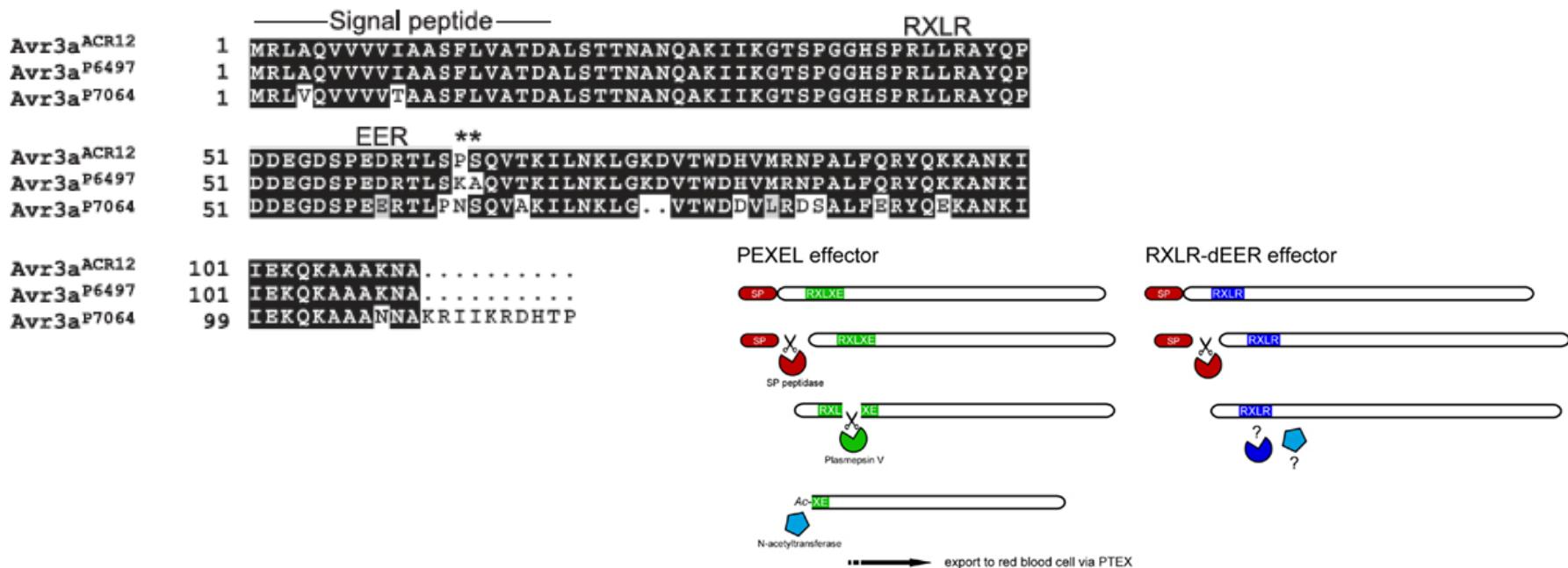


Greenberg & Vinatzer (2003) *Curr Opin Microbiol* doi:10.1016/S1369-5274(02)00004-8

Collmer *et al.* (2002) *Trends in Microbiology* doi:10.1016/S0966-842X(02)02451-4

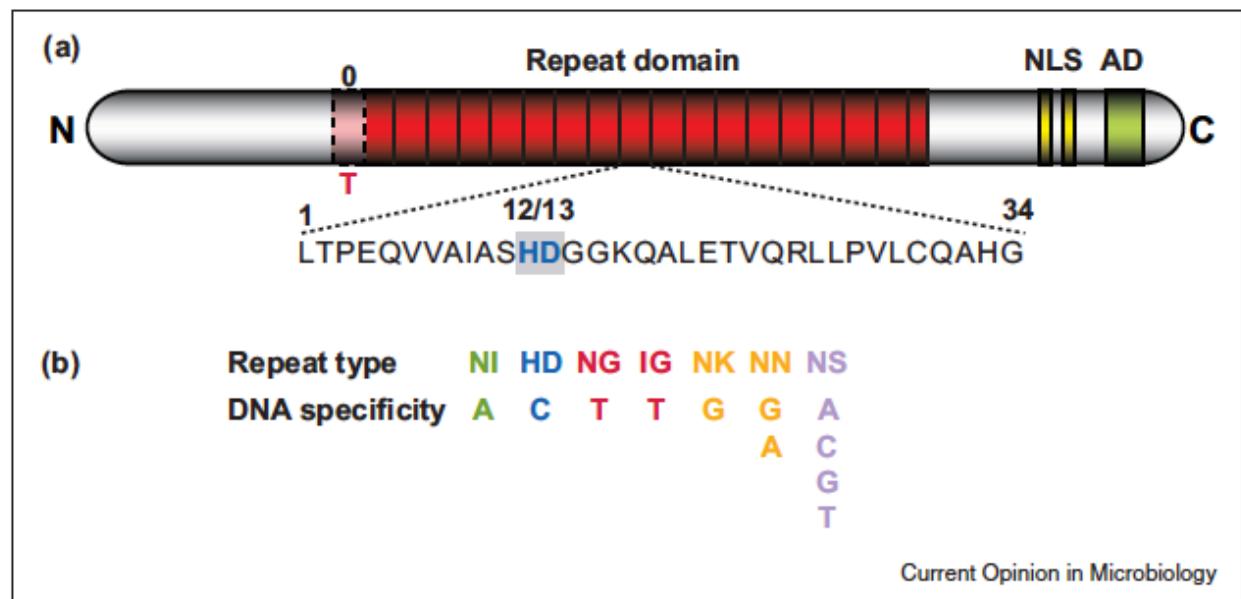
# Characteristics of known effectors

- Modularity
  - Delivery: localisation/translocation domain(s)
  - Activity: functional/interaction domain(s)



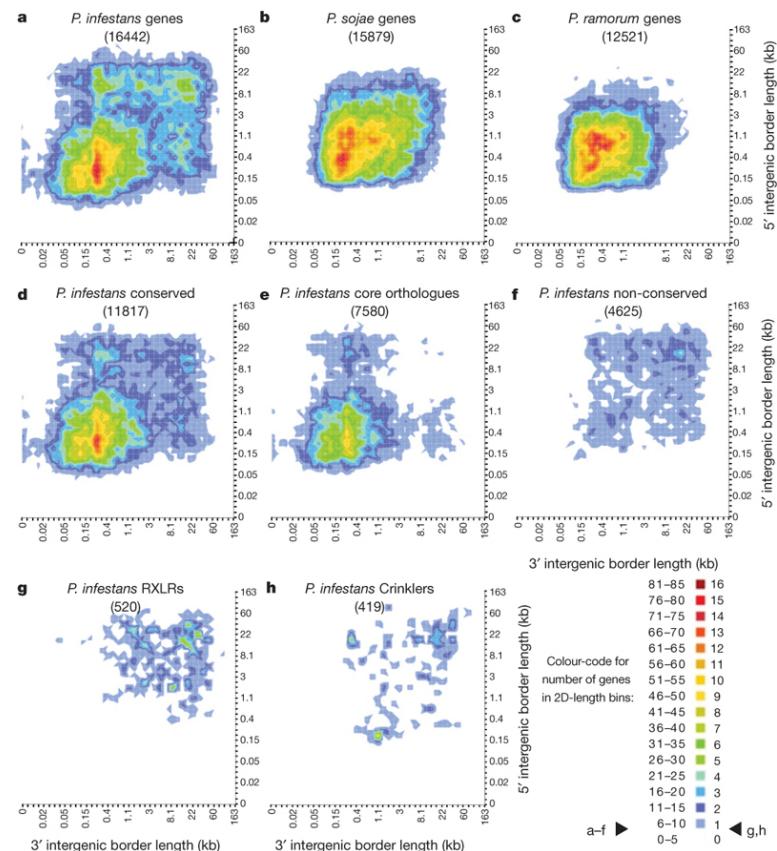
# Characteristics of known effectors

- Sequence motifs
  - Localisation/translocation domain(s) typically common to effector class (e.g. RxLR, T3E, CHxC)
  - Functional domain(s) may be common to effector class (e.g. TAL), or divergent (e.g. RxLR, T3E in general)



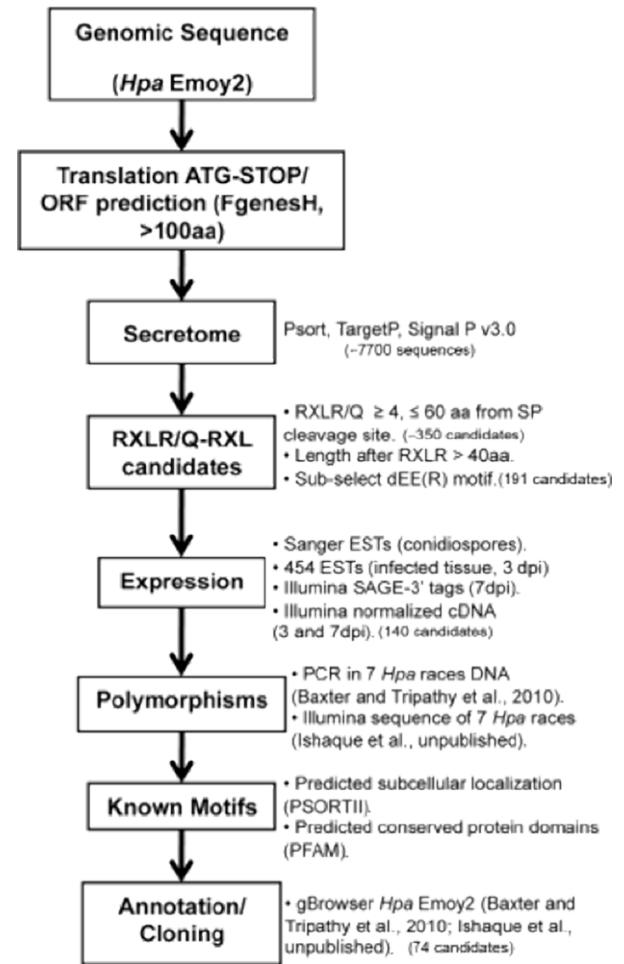
# Characteristics of known effectors

- “Arms Races” occur:
  - Host defences track effector evolution
  - Effectors evade host defences
- Divergence of effectors under selection pressure
  - Diversifying selection; divergence may result from evasion of detection, rather than change of biochemical ‘function’
- Effectors may be found preferentially in characteristic locations
  - P. infestans* ‘gene sparse’ regions



# Characteristics of known effectors

- Application of ‘filters’: reduce the number of sequences to check (similar to stratification)
  - Presence/absence filters:
    - SignalP (export signal)
    - RxLR/T3SS (translocation signal)
    - Expression (used by pathogen)
    - Positive selection (suggests arms race)
    - Location in repeat-rich regions
    - etc...
- Workflows (e.g. Galaxy, Taverna) can be useful here



Fabro et al. (2011) *PLoS Pathog.* doi:10.1371/journal.ppat.1002348.

Cock et al. (2013) *PeerJ* doi:10.7717/peerj.167/.

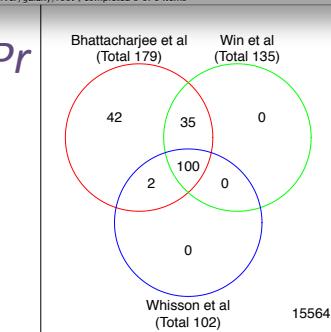
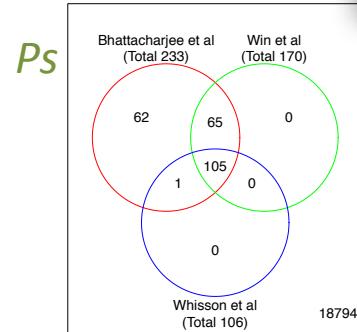
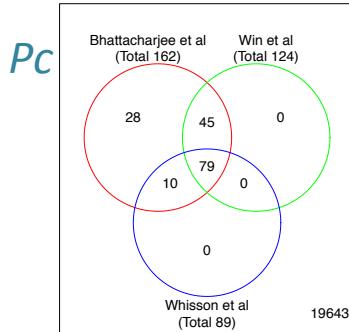
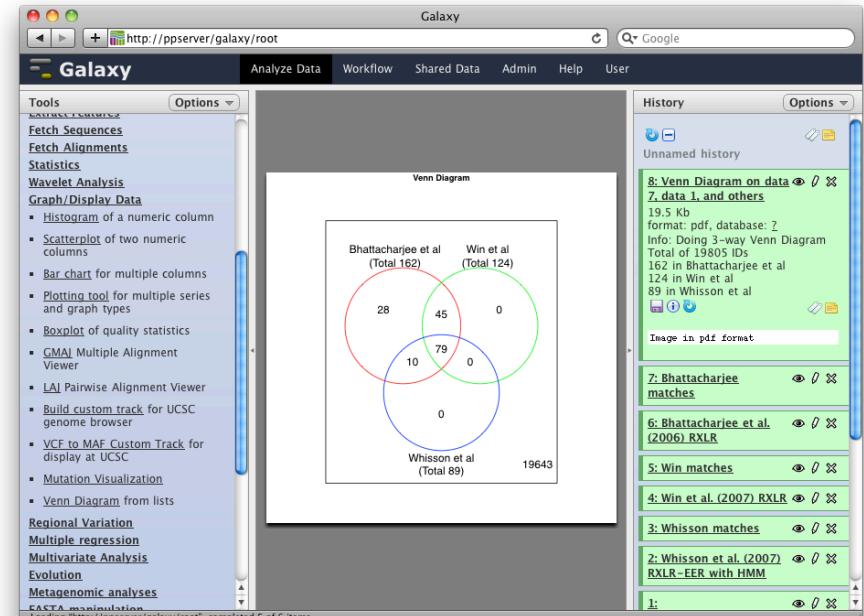
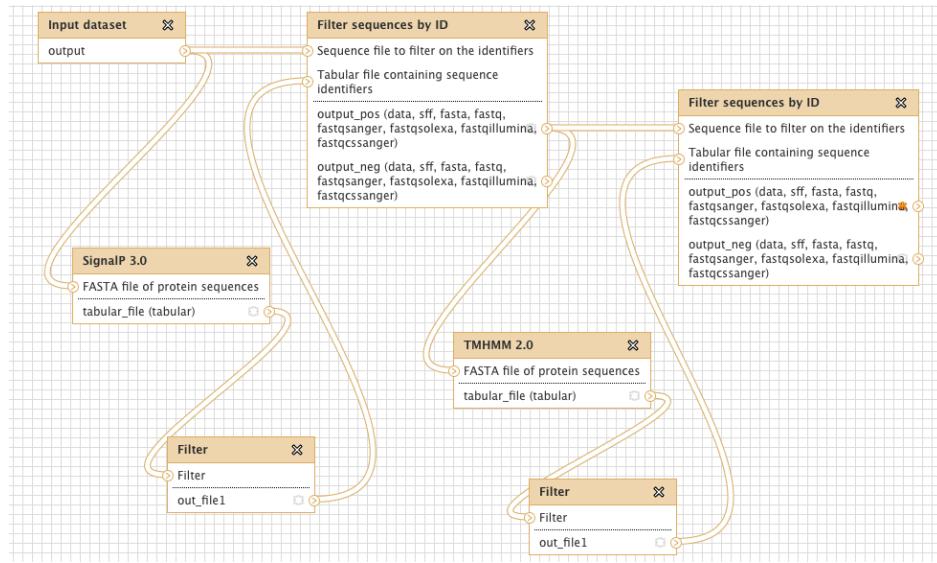
Cock & Pritchard (2014) *Methods Mol. Biol.* doi:10.1007/978-1-62703-986-4\_1

Pritchard & Broadhurst (2014) *Methods Mol. Biol.* doi:10.1007/978-1-62703-986-4\_4

# Galaxy Workflow – Effector finding

- Public Galaxy workflow for RxLR-finding

[http://toolshed.g2.bx.psu.edu/repository?repository\\_id=0984925f74bd671c](http://toolshed.g2.bx.psu.edu/repository?repository_id=0984925f74bd671c)

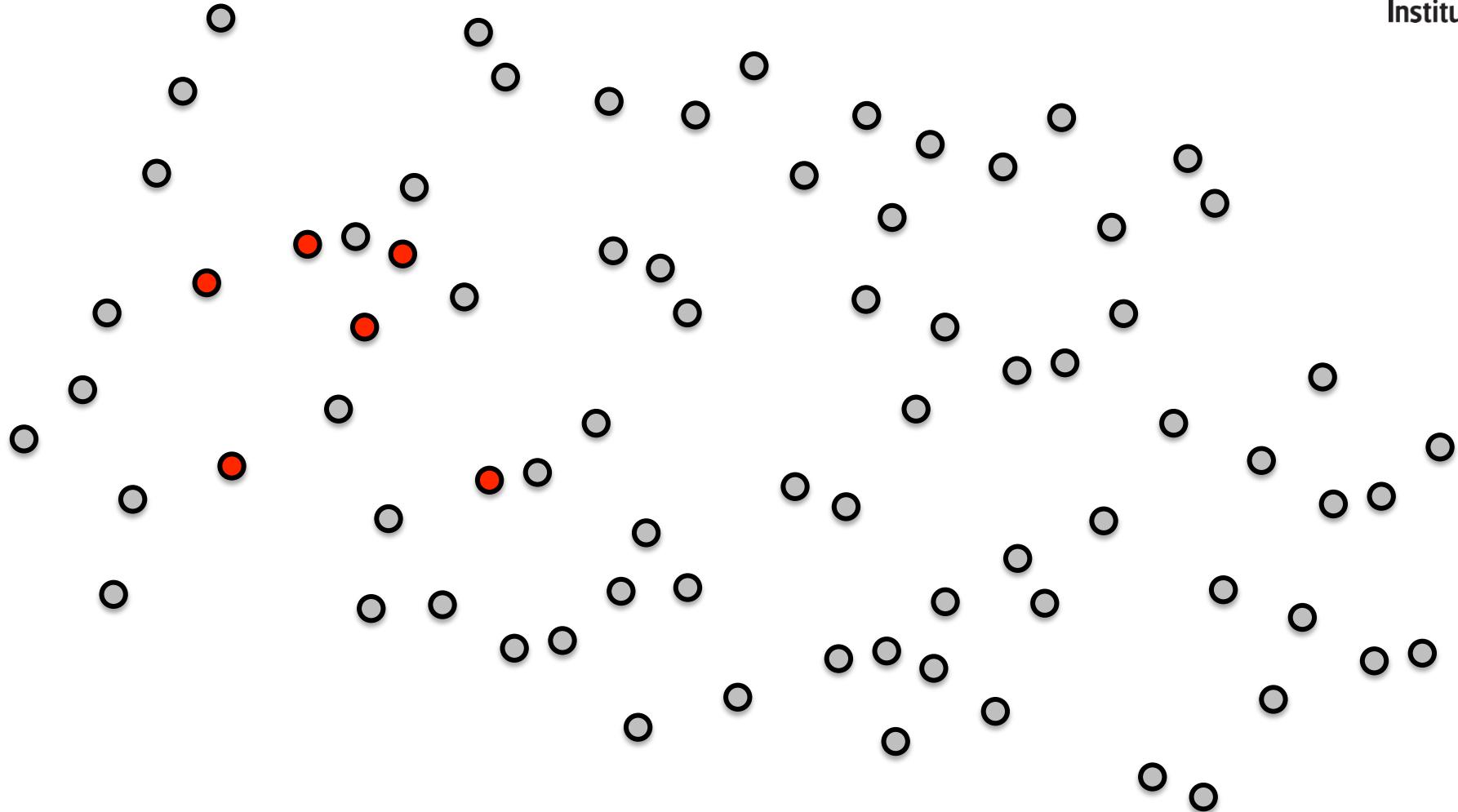


Cock et al. (2013) PeerJ doi:10.7717/peerj.167/

Cock & Pritchard (2014) Methods Mol. Biol. doi:10.1007/978-1-62703-986-4\_1

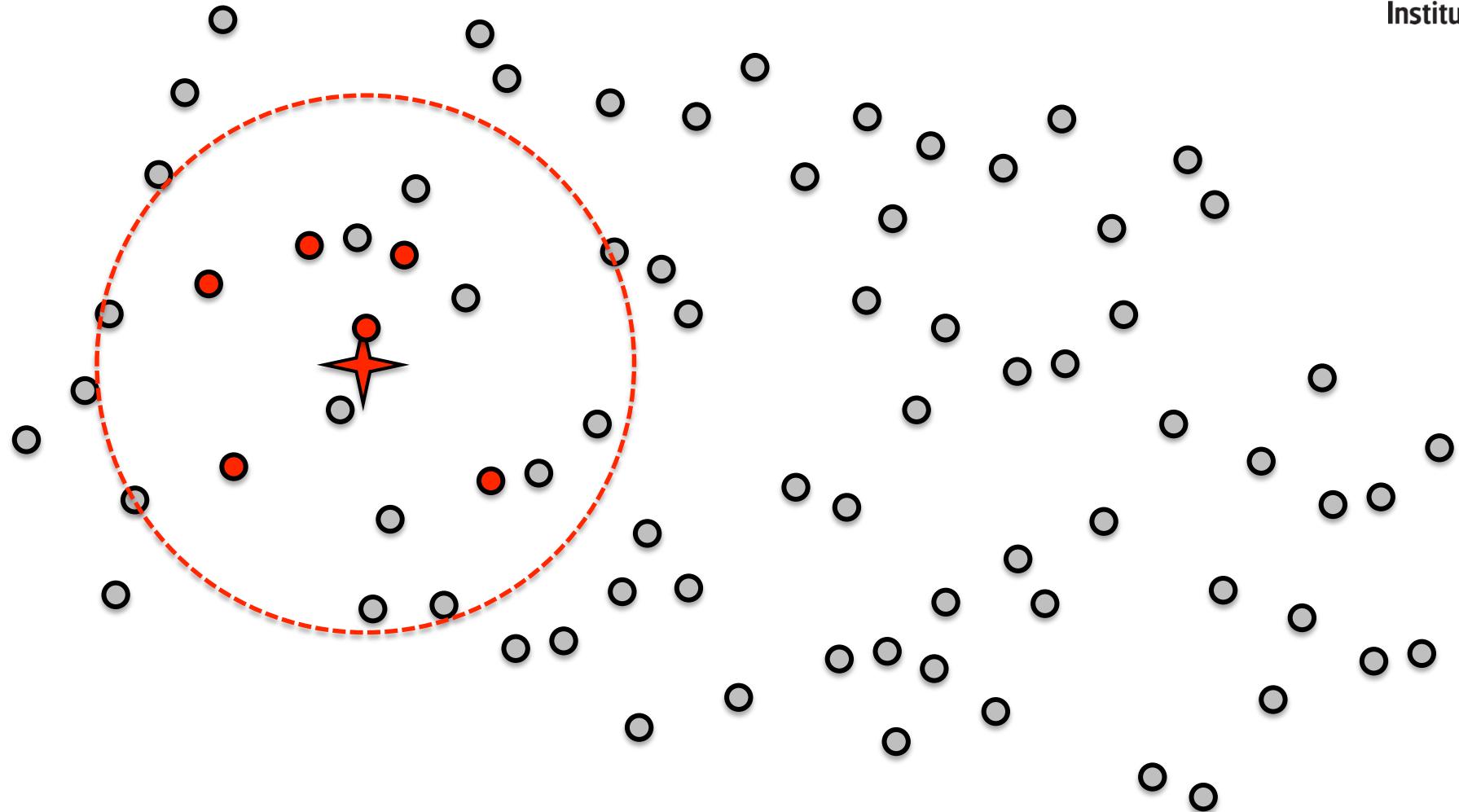


# Sequence space



Known exemplars: red

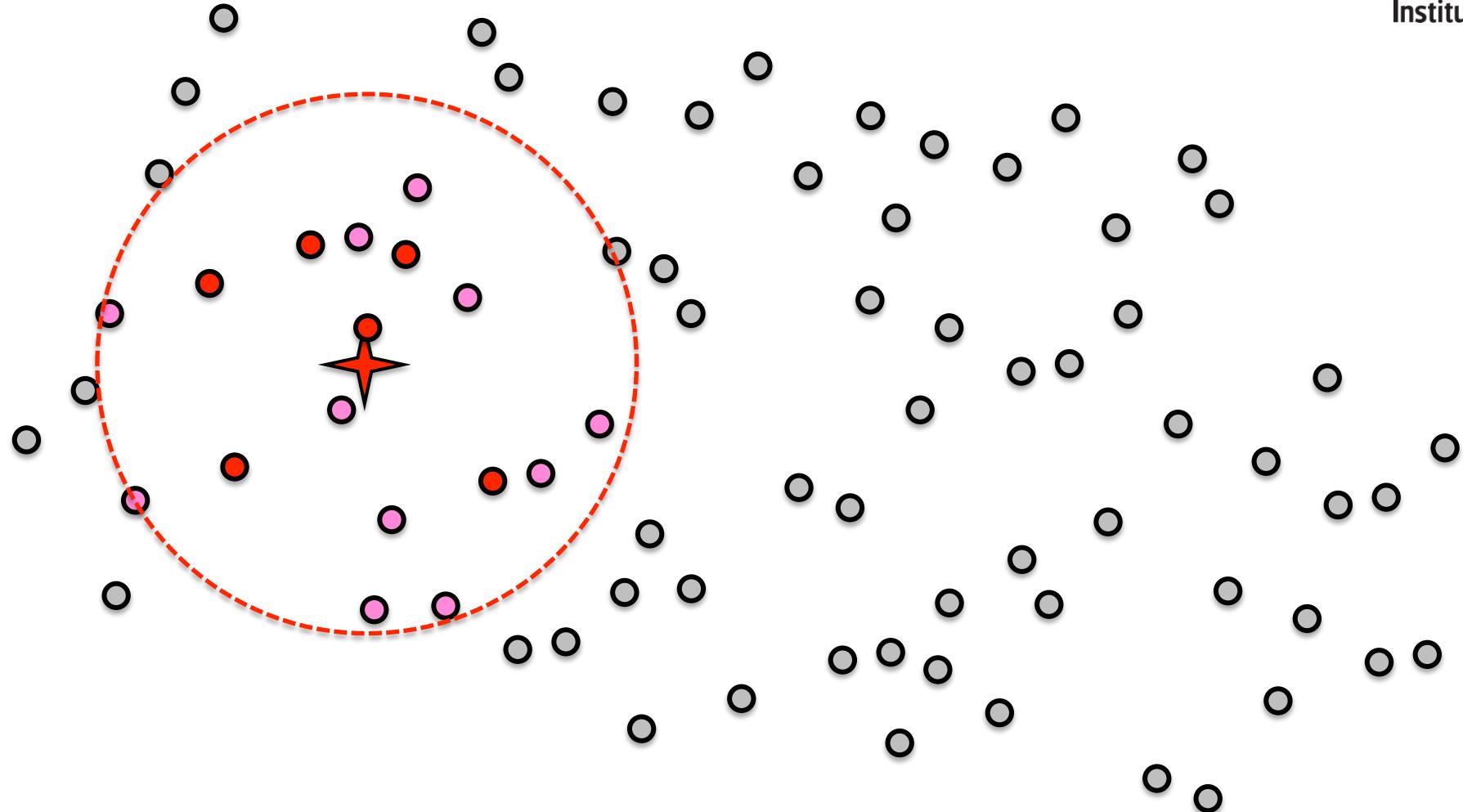
# Sequence space



Define a centre, and a distance that includes the examples



# Sequence space



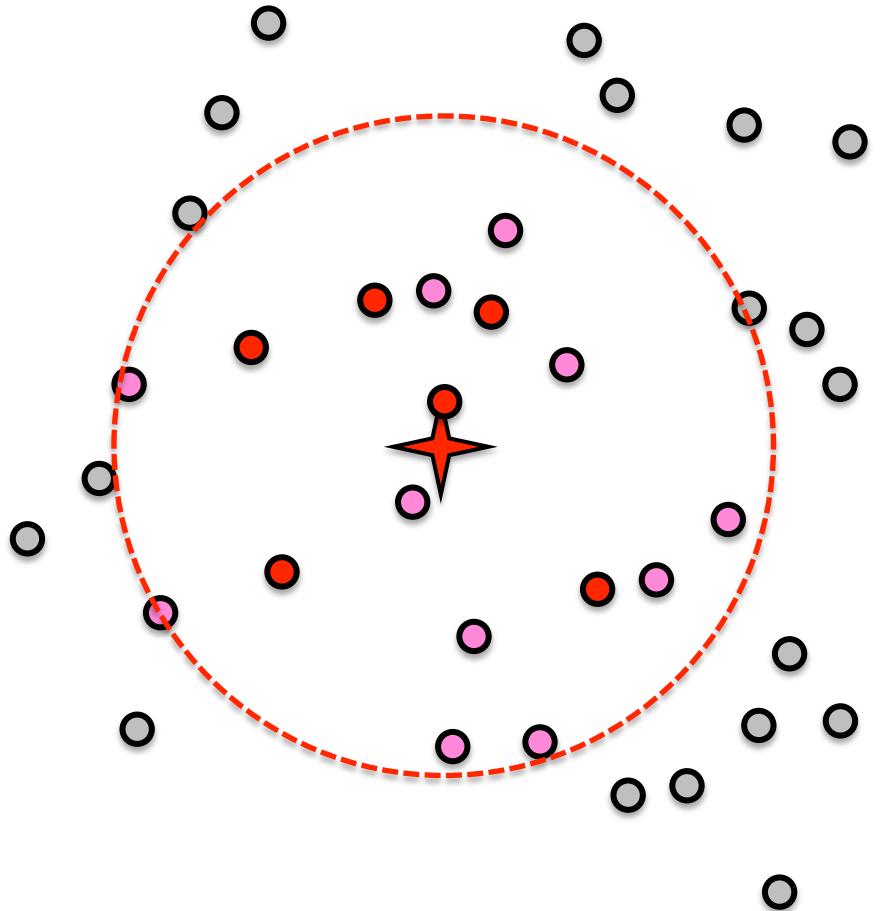
Classify ‘similar’ sequences

# Finding effectors

- Finding effectors is simple(-ish):
  1. Have one or more examples of your effector (class)
  2. Define some kind of appropriate threshold of similarity
  3. Check all the gene/gene product sequences in the genome against that threshold



# It's not *that* simple



- How do we define 'distance'?
- How large a 'distance' do we take?
- How do we know we've chosen a sensible 'distance'?



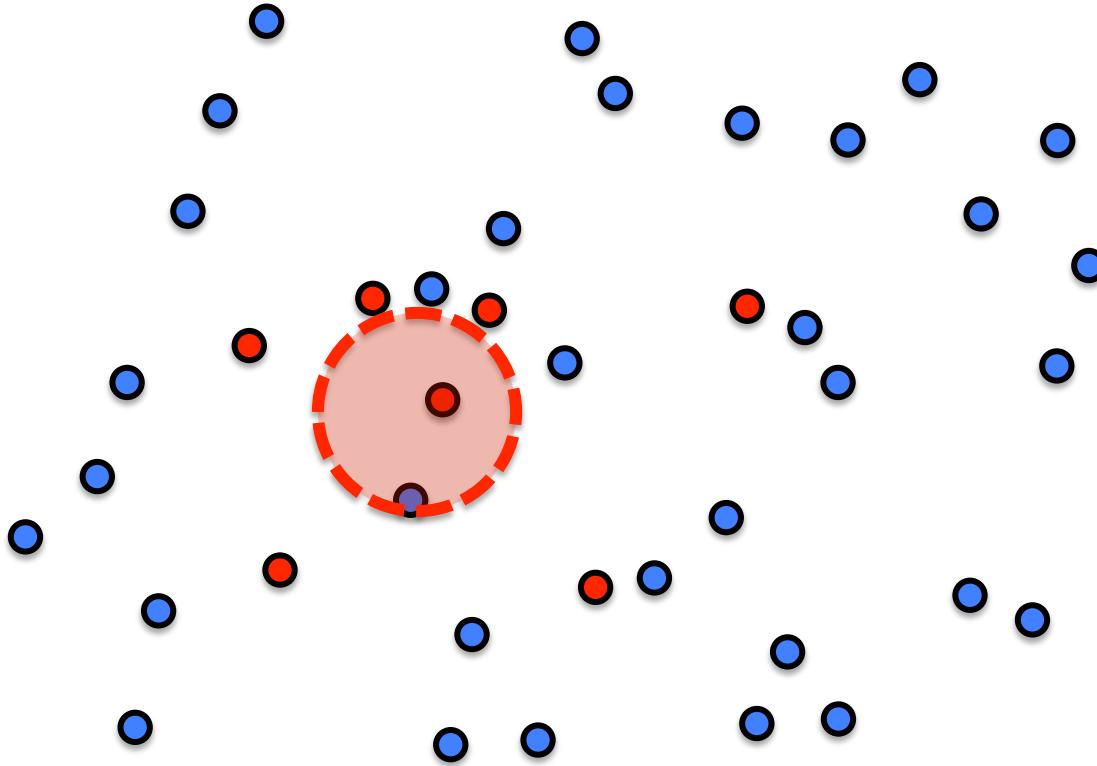
# How do we define 'distance'

- Sequence identity
- Derived score (based on sequence identity/alignment)
  - Bit score in BLAST
  - E-value in BLAST
- Derived score (based on other measures)
  - Bit score in HMMer
- Clustering (not strictly a distance)
  - Sequence identity (e.g. CD-HIT)
  - MCL

(we're really assessing criteria for class membership)



# How large a distance do we take?



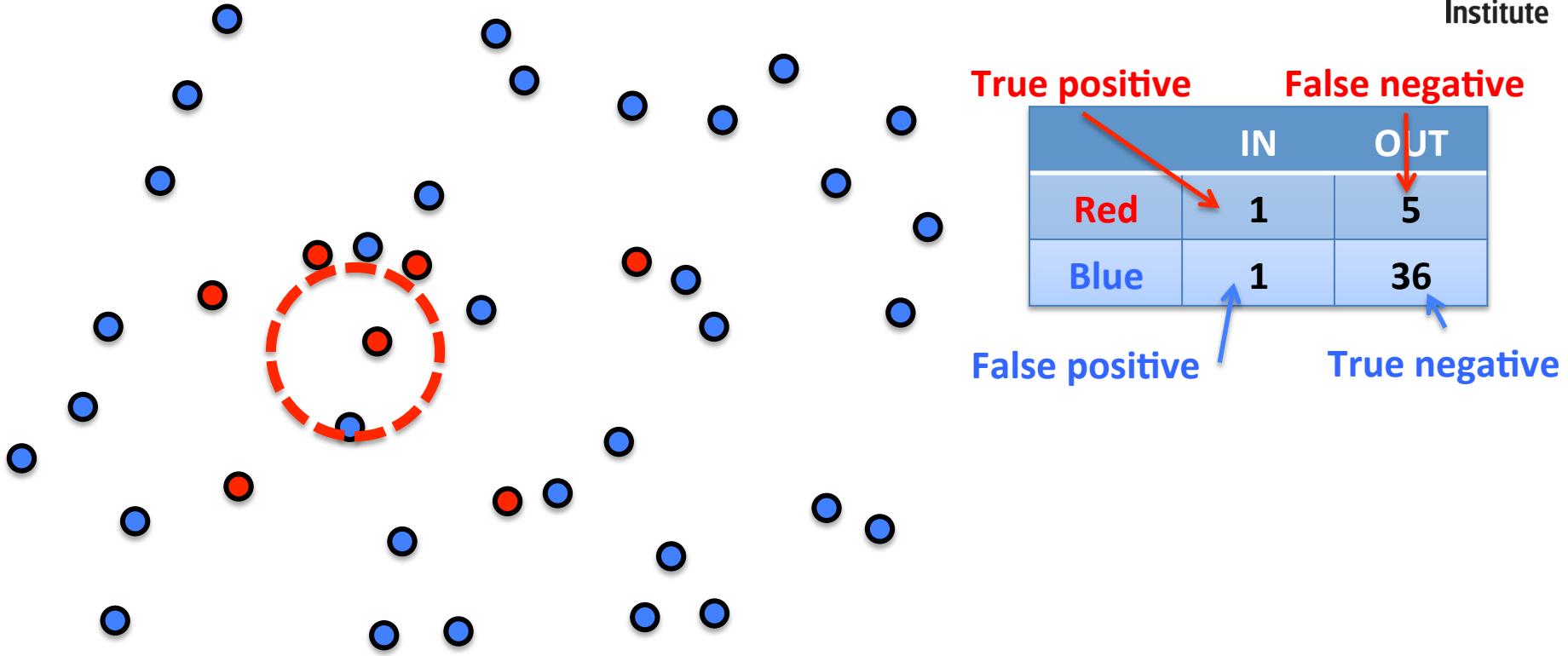
Confusion matrix:

	IN	OUT
Red	1	5
Blue	1	36

- Our distance/boundary classifies sequences as ‘in’ or ‘out’
  - ‘red’ or ‘blue’
- Changing distance/bound results in various degrees of success...



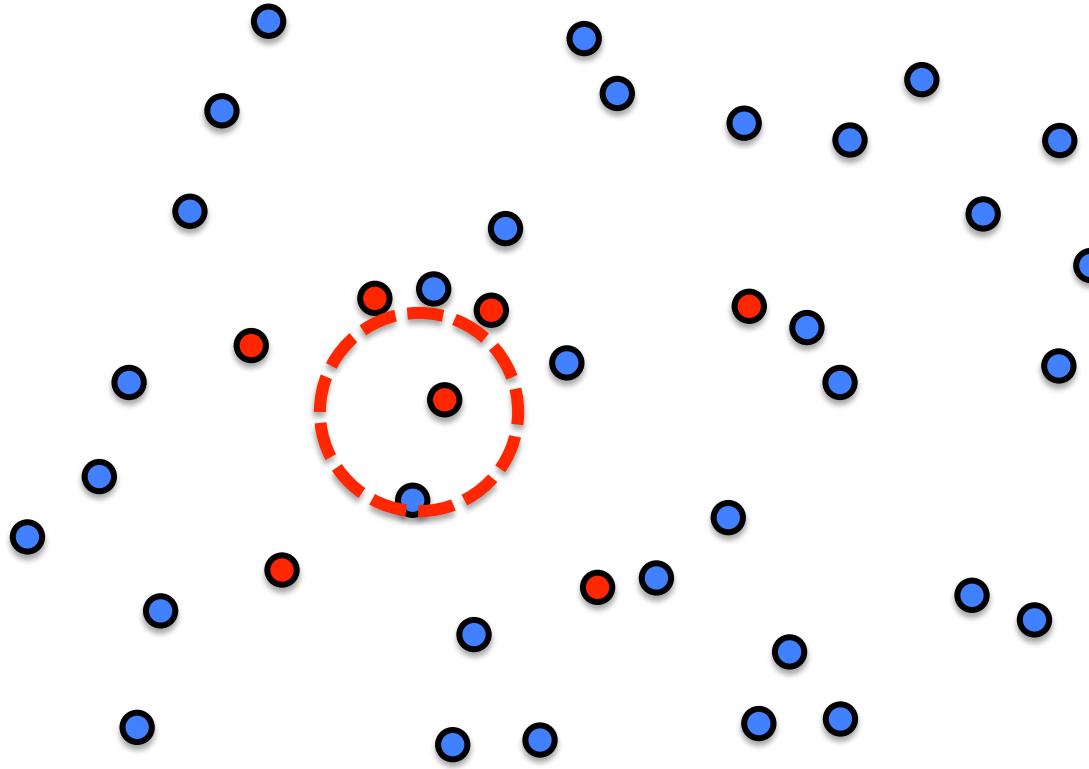
# How large a distance do we take?



- Our distance/boundary classifies sequences as 'in' or 'out'
  - 'red' or 'blue'
- Changing distance/bound results in various degrees of success...



# How large a distance do we take?



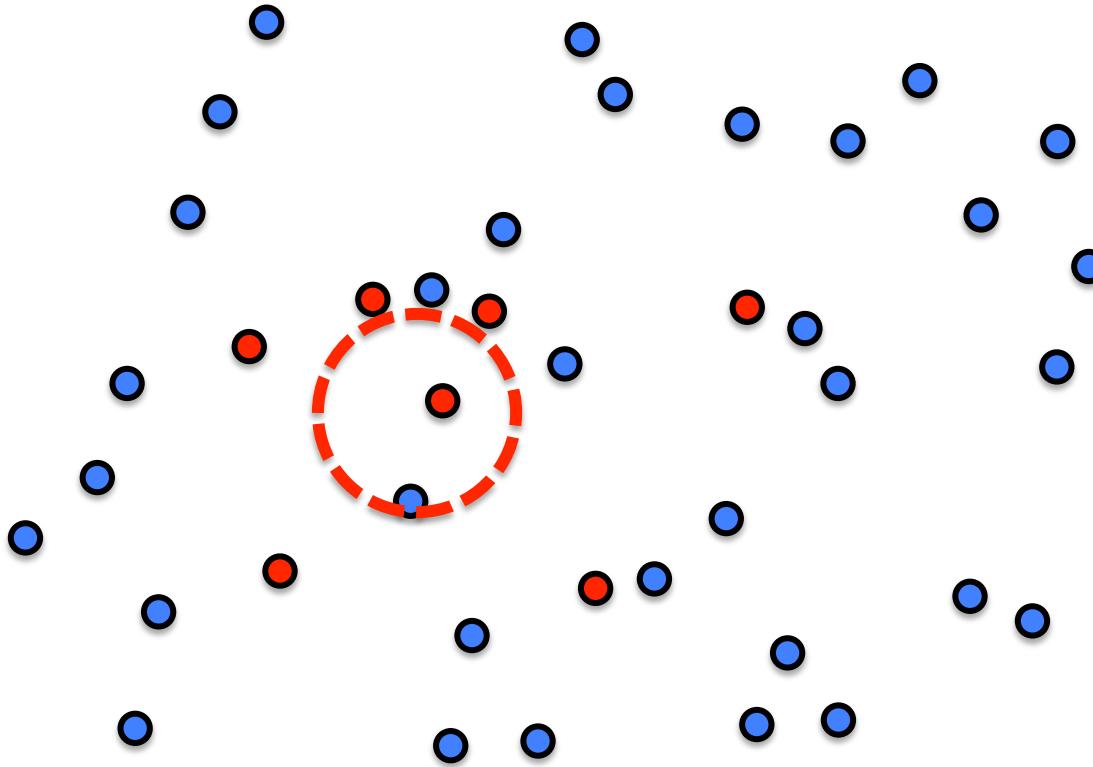
	IN	OUT
Red	1	5
Blue	1	36

False positive rate	$FP/(FP+TN)$
False negative rate	$FN/(TP+FN)$
Sensitivity	$TP/(TP+FN)$
Specificity	$TN/(FP+TN)$
False discovery rate	$FP/(FP+TP)$

- Our distance/boundary classifies sequences as ‘in’ or ‘out’
  - ‘red’ or ‘blue’
- Changing distance/bound results in various degrees of success...



# How large a distance do we take?



	IN	OUT
Red	1	5
Blue	1	36

False positive rate  $1/37 = 0.03$

False negative rate  $5/6 = 0.83$

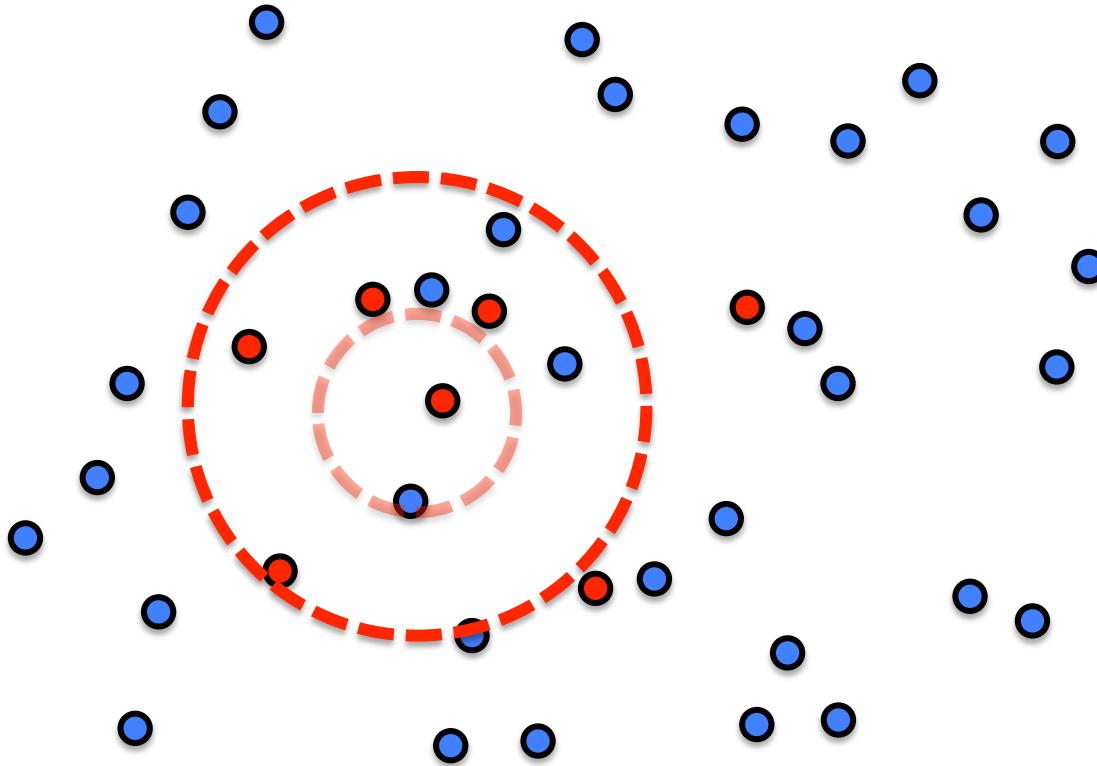
Sensitivity  $1/6 = 0.17$

Specificity  $36/37 = 0.97$

- Our distance/boundary classifies sequences as ‘in’ or ‘out’
  - ‘red’ or ‘blue’
- Changing distance/bound results in various degrees of success...



# How large a distance do we take?

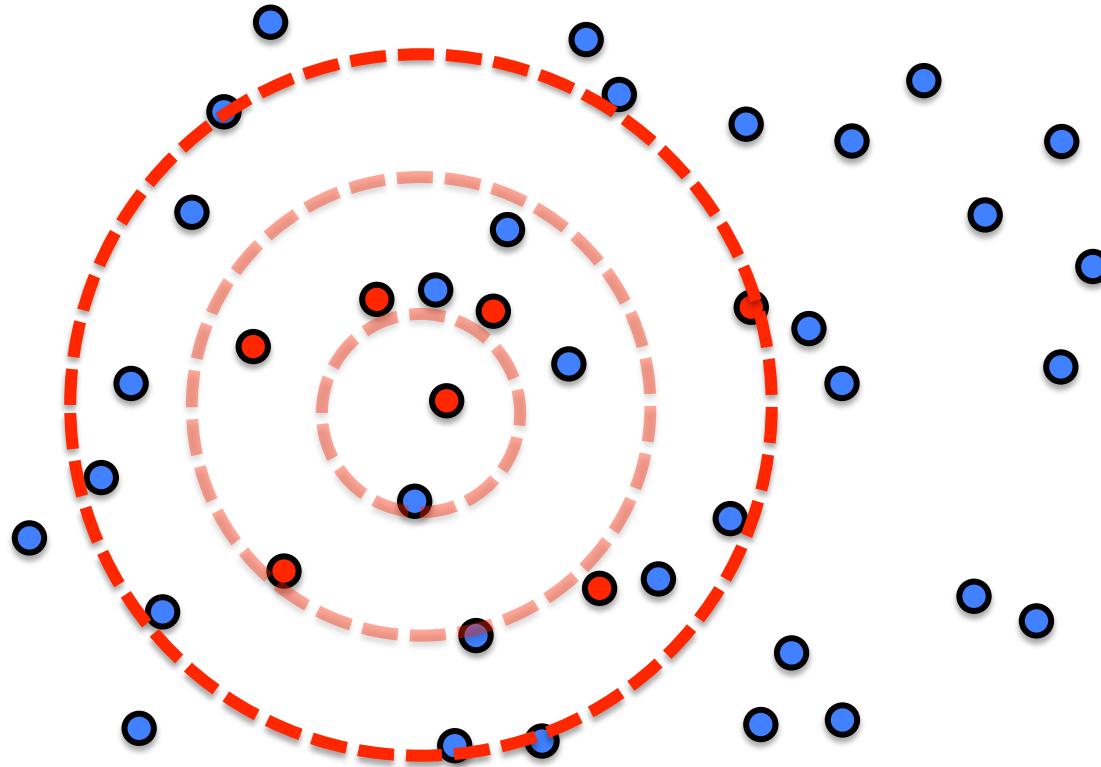


	IN	OUT
Red	5	2
Blue	4	33

False positive rate	0.11
False negative rate	0.29
Sensitivity	0.81
Specificity	0.89

- Our distance/boundary classifies sequences as ‘in’ or ‘out’
  - ‘red’ or ‘blue’
- Changing distance/bound results in various degrees of success...

# How large a distance do we take?



	IN	OUT
Red	7	0
Blue	14	23

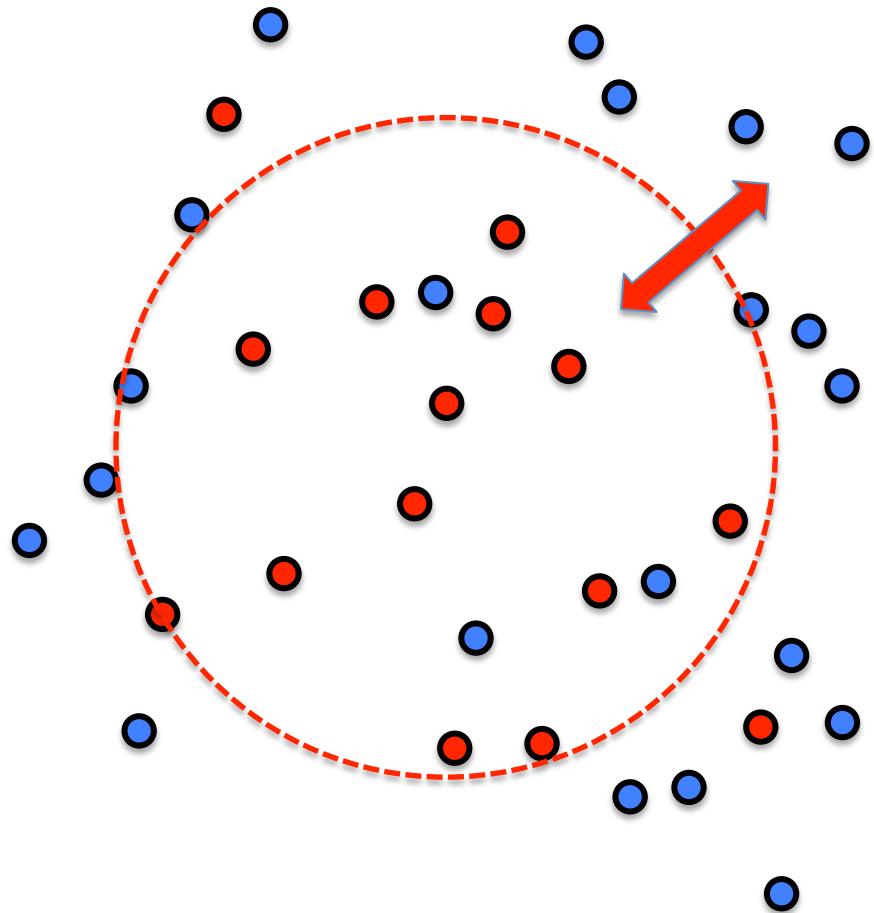
---

False positive rate	0.38
False negative rate	0
Sensitivity	1
Specificity	0.62

- Our distance/boundary classifies sequences as ‘in’ or ‘out’
  - ‘red’ or ‘blue’
- Changing distance/bound results in various degrees of success...



# How large a distance do we allow?



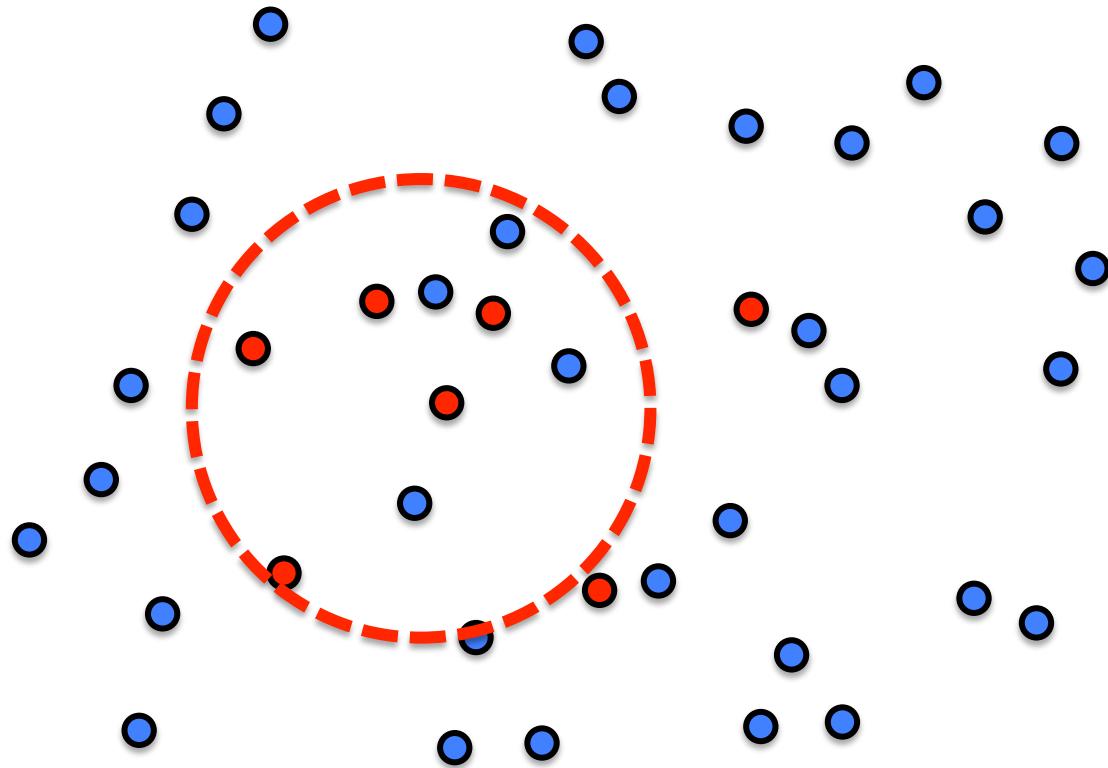
- Assign known 'positive' and 'negative' examples
- Vary distances and measure performance (F-measure, AUC, ...)
- Choose distance that gives the best performance



# Cross-validation

- Estimation of classifier performance depends on
  - distance measure
  - composition of training set ('positives' and 'negatives')
- Cross-validation gives objective estimate of performance
- Many strategies available, including:
  - leave-one-out (LOO)
  - $k$ -fold crossvalidation
  - repeated (random) subsampling
- Essentially: always keep a **hold-out set** (not used to train)

# After Cross-validation

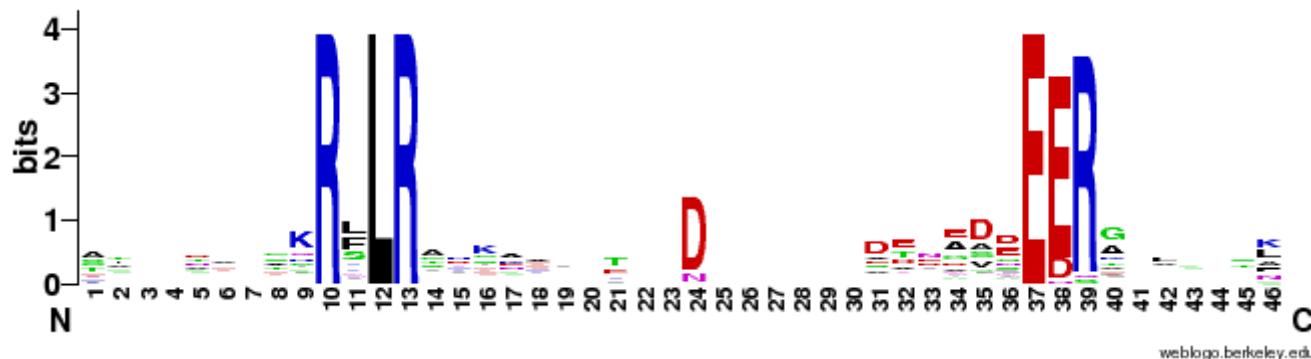


False positive rate	0.11
False negative rate	0.29
Sensitivity	0.81
Specificity	0.89
Precision	0.56

- Used crossvalidation to find ‘best’ method & parameters
- This gives estimated performance metrics on unseen data
- Apply ‘best’ method to complete dataset for prediction

# Building an RxLR training set

- Starting point: 49 candidate sequences (reference set)
- Known:
  - Contain (putatively) RxLR-EER motif
  - All but one transcribed (i.e. not bad gene calls)
- Assumed:
  - Presence of signal peptide and RxLR-EER categorises effectors



# Building a classifier

- We have a recognisable motif, with substantial local variation and indels
  - Therefore chose profile HMM
  - Use HMMer software
- Profile HMMs sensitive to quality of alignment**
- Therefore treat alignment as a parameter of the HMM (much difference between alignments – several tested!)

```

Contig7763_1
RD20□ (rpvb_3739.y1.
E7.6019.C1
rpvb_13318.y1.abd
E7.5364.C1
PG001C8
E7.9705.C1
PIBAC-49P21|rpvb_119
PK001F10.XT7
rpch_8164.y1.abd
Contig22627_1
E7.5930.C1
MY-20-B-07
E7.6418.C1
ipi01_a
Contig1828_1
rpch_15494.y1.abd
PiPEX34A111
rpvb_12884.y1.abd
ipi01_b
Contig4167_1
rpccy_0850.y1.abd
E7.1257.C1
Contig5688_1
E7.9043.C1
rpccm_4494.y1.abd□
Contig5386_1
E7.4706.C1
E7.6301.C1
PV004C8.XT7
E7.6199.C1
Contig5191_1
MY-15-C-09|avr3a
E7.4461.C1
E7.8373.C1
E7.8848.C1
rpch_12409.y1.abd
E7.6112.c1
rpvb_11175.y1.abd
1 ASDSVVGVRSLRSHKVNKEEKEDL---VDSREDVEERK-G---- 37
1 AVRASFNTKRALRSHTKATDH-----GEERA-YKPSL 31
1 --VDYATTERLLRAHSSDKE-----EQKEEEER-AISINF 32
1 TDDGPKHTKRFRLRGESSKIV--NL--KQEEGVFEERKGVSQKL 39
1 --EATDIMRQLRVLGKAV-----ASLFENHQSTREL 30
1 DTVAKDNKKRLLRAYKDA-----EDDSEDIKNVKPTA 32
1 ADQSIVIEQPRFLRDGKIA-----EGDNEERVNAKEA 32
1 TVINTNHDRVILRNR-----VTNAAEERAGWGVSK 31
1 -----SGRVLRADATYN-----QATNAVDEERAS---- 24
1 SIV-NTEGGRILLRGVKKR-----TAEREVQEERMS---- 29
1 RYPLPRNNGRILRRRAKQHA-----TNEVGVEEERFYT---K 33
1 SSEGTHETARLLRNLAQP-----PVETGNQEERTIN-FAS 35
1 VHLAGEREKRLLRFDNSD-----YRDDDEEERA----N 30
1 --NNNQEFAIRLNTEERSIAA-----ILAEGEE|RRAA---- 32
1 STEYNADEEKSRLRGDYNNEVTK-----EPNTSDEERAFS-ISK 37
1 --MTYTVSRKRLLRVAGED-----DDATDDEE|RG-LGS-- 32
1 --DAHIHAGRVLRDRSRV-----DEERG-LPT-- 24
1 -----CAKSRHLRNGKDALWNYDTGGINSIVADDEERV-VS-- 37
1 -AEEENSIVRSLRAVE-----TSEDEEERD-LLGLF 29
1 -VKLIVNKNRRLRGAQ-----NQHEEEERD-VGISN 28
1 -----LSQRHLRSHDTPV-----LV-DDYNADEERG-LDNA-- 29
1 ---VSIVSIRKLRSH-----EERG-LP-- 18
1 -IAGDFVSKRKLRADKA-----A-DQVSTEERG-VTGF- 31
1 -IDSSNTAKRSRSLRQVAKASQESVDYLSSK-SKYVAAEERG-WLTK- 42
1 -LNVDAVDRFLRERKE-----SMEERG-YHLT- 26
1 -QFVVAMGRRSLRTSGE-----ANEERT-RLNT- 26
1 -EQN-VLGKRSRSLRTDHMRVSTV-----EDQEGDEERI-FRRF- 34
1 -EPNKHVATRSLRHP-----I-----DDSSDGEERL-LNGM- 30
1 -VQ-DDNYDRQLRGFYATENTD-PVNNQDTAHEDEERG-NVAT- 40
1 --I-TSQVQRILLRTHHATIK-----VNADSEERF-LTEP- 30
1 -----DDITSRNLR-----SCEERA-YAFV- 20
1 -NSAKDAPAKRFLRKTSK-----TDEDNE|RRA-F---- 27
1 --IHDSAGRRLLRKNEE-----NEETSEERA-PNF- 27
1 -IPTHVHDQRLRRVRVN-----DEGELTEERT--GG-- 28
1 -TTIATSSKRFRLRYDAEVR---DTVRG-DNDVDREERG-SSP-- 37
1 -AQLVRSGKRFRLRPTAEQL---DDEL---GQVDEERR-MEI-- 34
1 -TKNGVLAKRFRLRAQGP-----PDEERGRLKD-- 26
1 -TRWNVADKRLLRANDG-----TNAAAEERG-MAD-- 28
1 -YSPTTESYRFLRARKN-----EDASDEERG-ISN-- 28

```

# Building a classifier

- Parameters modified for HMM
  - Alignment package (no alignment, anchored, Clustal, DiAlign, T-Coffee) on default settings
  - Full-length and truncated (no signal peptide) alignments to test for influence of signal peptide region on classifier
    - Plus one alignment of RxLR-EER plus flanking region only ('cropped')
- HMM built for each of eleven alignments
  - Default HMMer parameters
- Once built, the HMM was the classifier.  
Stratification by presence of signal peptide (SignalP)

**hmmbuild --amino <output> <alignment>**

# Testing the classifiers

- Eleven classifiers to test
- **Step 1: Consistency test**
  - Does the classifier correctly call as positive the sequences used to train it?
  - Estimates recovery of the information in the training set
- **Step 2: Recovery of full sequences**
  - Estimates performance of classifier on complete sequence data

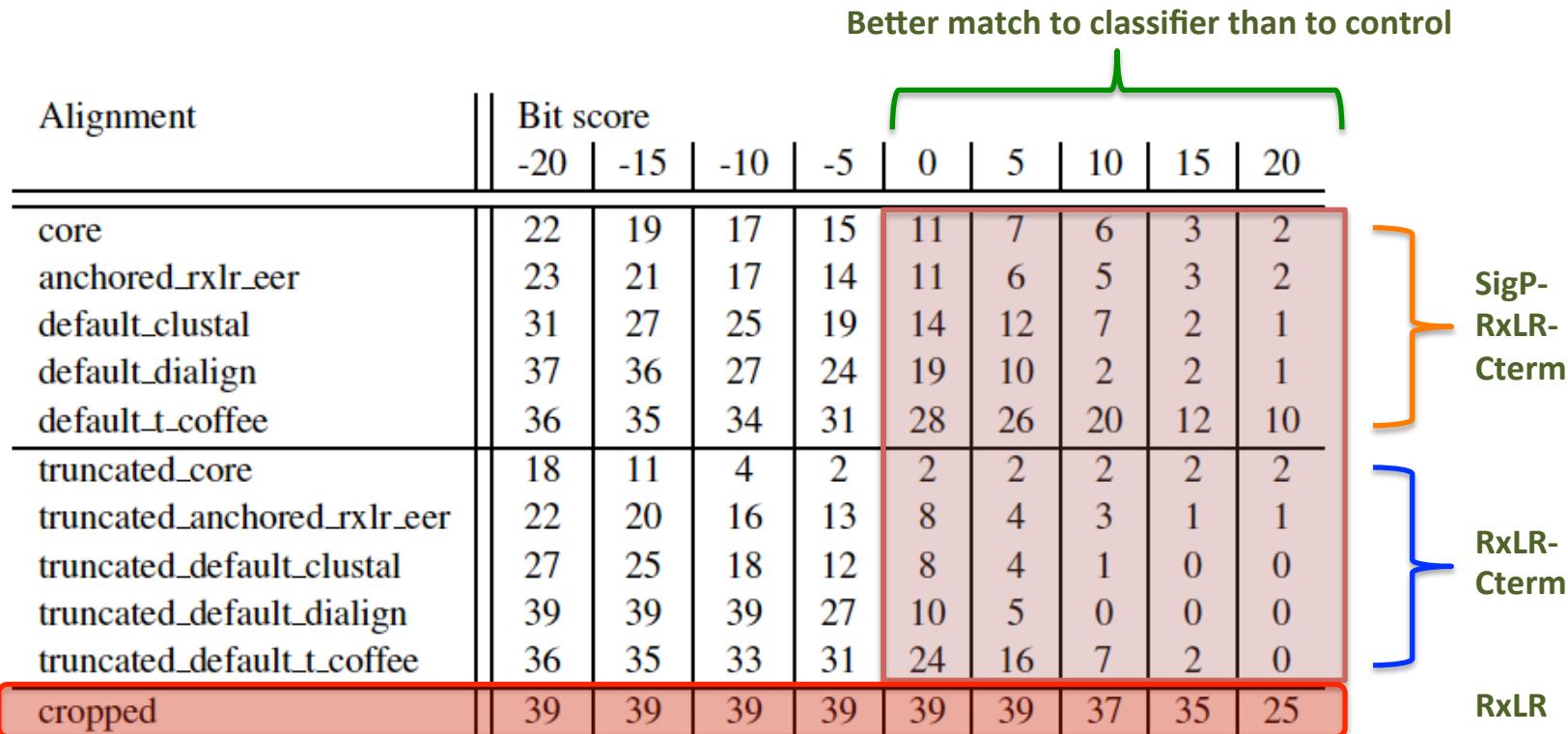
Alignment	Consistency	Reference	Minimum bit score
core_rxlr_aa	38	38	10.0
core_aa_anchored_rxlr_eer	36	36	18.6
core_default_clustal	33	38	12.9
core_default_dialign	2	38	6.8
core_default_t_coffee	14	38	5.8
truncated_core_aa	33	33	3.6
truncated_core_aa_anchored_rxlr_eer	34	34	13.8
truncated_default_clustal	31	35	4.4
truncated_default_dialign	5	34	1.9
truncated_default_t_coffee	17	37	2.0
cropped	36	39	5.6

} SigP-RxLR-Cterm  
} RxLR-Cterm  
} RxLR

# Testing the classifiers

- **Step 3: Leave-One-Out Crossvalidation**

- But only have positive examples!
- Removes possibility that classifier matches on basis of having ‘seen’ a sequence before





# Testing the classifiers

- **Step 4: Tests on negative samples**
  - Completely shuffled sequences
  - Shuffled downstream of the signal peptide only
  - Replace RxLR-EER with AAAA-AAA

**No classifier identifies a false positive**  
(no classifier matches on sequence composition alone)



# Testing the classifiers

- **Step 4: Tests on negative samples**
  - Completely shuffled sequences
  - Shuffled downstream of the signal peptide only
  - Replace RxLR-EER with AAAA-AAA

(some recognition on basis of signal peptide)

Alignment	Recovered	Reference	Rate	
core_rxlr_aa	3	38	8%	SigP- RxLR- Cterm
core_aa_anchored_rxlr_eer	0	36	0%	
core_default_clustal	9	38	24%	
core_default_dialign	2	38	5%	
core_default_t_coffee	0	38	0%	
truncated_core_aa	0	33	0%	RxLR- Cterm
truncated_core_aa_anchored_rxlr_eer	0	34	0%	
truncated_default_clustal	0	35	0%	
truncated_default_dialign	0	34	0%	
truncated_default_t_coffee	0	37	0%	
cropped	0	39	0%	RxLR



# Testing the classifiers

- **Step 4: Tests on negative samples**

- Completely shuffled sequences
- Shuffled downstream of the signal peptide only
- Replace RxLR-EER with AAAA-AAA

(some recognition on sequence other than motif)

Alignment	Recovered	Reference	Rate	
core_rxlr_aa	39	38	102%	SigP- RxLR- Cterm
core_aa_anchored_rxlr_eer	27	36	75%	
core_default_clustal	35	38	92%	
core_default_dialign	38	38	100%	
core_default_t_coffee	32	38	84%	
truncated_core_aa	20	33	61%	RxLR- Cterm
truncated_core_aa_anchored_rxlr_eer	16	35	46%	
truncated_default_clustal	33	35	94%	
truncated_default_dialign	12	34	35%	
truncated_default_t_coffee	16	37	43%	
cropped	0	39	0%	RxLR

# Choosing a classifier

- The ‘cropped’ classifier has:
  - 100% recovery of positive training sequences
  - 0% recovery of negative test sequences
- Some variation in classifier performance on whole genome:  
**284 candidates from ‘cropped’** (~80-200 common to other models)

Method	total	HMM					truncated					cropped
		core	anchored	full-length	clustal	dialign	t-coffee	core	anchored	clustal	dialign	
core_rxlr_aa	1999	X										
core_aa_anchored_rxlr_eer	1994	1890	X									
core_default_clustal	541	478	465	X								
core_default_dialign	196	89	90	107	X							
core_default_t_coffee	292	106	111	141	166	X						
truncated_core_aa	929	923	925	402	60	63	X					
truncated_core_aa_anchored_rxlr_eer	195	170	184	149	75	87	163	X				
truncated_default_clustal	81	65	70	64	68	77	62	69	X			
truncated_default_dialign	131	44	50	52	63	75	35	44	39	X		
truncated_default_t_coffee	143	70	87	79	103	118	60	77	69	53	X	
cropped	284	84	98	114	133	191	60	86	71	80	120	X

Whisson *et al.* (2007) *Nature* [doi:10.1038/nature06203](https://doi.org/10.1038/nature06203)

Haas *et al.* (2009) *Nature* [doi:10.1038/nature08358](https://doi.org/10.1038/nature08358)

# A Trip To The Doctor

- Routine medical checkup
- Test for disease X (horrible, unpleasant, potentially suppurating)
- Test has *sensitivity* (i.e. predicts disease where there is disease) of **95%**
- Test has *false positive rate* (i.e. predicts disease where there is no disease) of **1%**
- Your test is **positive** (bad luck, there)
- What is the probability that you have disease X?

0.01 0.05 0.50 0.95 0.99



# A trip to the doctor, part II

- Test for disease X (horrible, unpleasant, potentially suppurating)
- Test has *sensitivity* (i.e. predicts disease where there is disease) of **95%**
- Test has *false positive rate* (i.e. predicts disease where there is no disease) of **1%**
- Your test is positive
- To calculate the probability that the test correctly determines whether you have the disease, you need to know the **baseline occurrence**.

**Baseline occurrence: 1%  $\Rightarrow P(\text{disease} | +\text{ve})=0.490$**

**Baseline occurrence: 80%  $\Rightarrow P(\text{disease} | +\text{ve})=0.997$**

## What is the baseline occurrence for effectors?

- Usually rely on predictions for expected baseline
- Bacterial genomes:  $\approx 4500$  genes
  - Type III effectors: 1-10% (Arnold *et al.* 2009); 1-2% (Collmer *et al.* 2002); 1% (Boch and Bonas, 2010)
- Oomycete/fungal genomes:  $\approx 20000$  genes
  - RxLRs: 120-460 (1-2%; Whisson *et al.* 2007);  $\leq 563$  ( $\leq 2\%$  Haas *et al.* 2009);
  - CRNs: 19-196 ( $\leq 1\%$ ; Haas *et al.* 2009)
  - CHxC:  $\approx 30$  (<1%; Kemen *et al.* 2011)
- **We need to take care over result interpretation:**
  - Prediction method with 5% false negative rate and 1% false positive rate, with 1% baseline, predicting 500 effectors:
    - $P(\text{effector} | \text{positive test}) \approx 0.5$



# A lesson from the literature?

- “The resulting computational model revealed a strong type III secretion signal in the N-terminus that can be used to detect effectors with sensitivity of 71% and [specificity] of 85%.”
  - Sensitivity [ $P(+ve|T3E)$ ] = 0.71; FPR [1-Specificity;  $P(+ve|not T3E)$ ] = 0.15
  - Base rate [ $P(T3E)$ ] ≈ 3%; Genes = 4500

$$P(T3E|+ve) = \frac{P(+ve|T3E)P(T3E)}{P(+ve|T3E)P(T3E) + P(+ve|\overline{T3E})P(\overline{T3E})}$$

- We expect  $P(T3E|+ve) \approx 0.13$
- (and a significant number, up to 15% of the genome, of false positives... see `pred_acc.ipynb`)

# A lesson from the literature?

- “The surprisingly high number of (false) positives in genomes without TTSS exceeds the expected false positive rate (Table 1)”

TTSS status	G+C content	Number of proteins	Positives	Z-Score
-	56.3%	1841	3.7%	0.6
-	43.1%	1963	1.6%	-2.3
-	51.5%	5349	2.3%	8.3
-	61.0%	1818	2.9%	6.8
-	62.0%	4129	2.9%	5.7
+	57.9%	5169	3.8%	7.2
+	59.2%	3088	3.8%	0.7
+	58.3%	5607	3.7%	0.6
-	42.8%	2120	10.8%	17.0
-	44.8%	2385	12.2%	20.0
-	40.1%	3541	3.3%	0.6
-	45.8%	4008	3.5%	0.6
+	51.4%	4510	2.8%	0.6
+	52.0%	4614	2.7%	0.6
+	52.1%	4659	2.8%	0.6
+	52.1%	4617	2.8%	0.6
+	52.2%	4205	2.8%	0.6
+	52.2%	3964	2.6%	0.6
+	52.1%	4779	2.5%	0.6
+	52.2%	4805	2.6%	0.6
+	52.2%	4091	2.7%	0.6
+	52.1%	5601	3.1%	0.6
+	52.2%	4627	2.8%	0.6
+	51.9%	4753	2.9%	1.9
+	52.1%	4312	2.7%	0.8
+	52.2%	4523	2.8%	1.8
-	53.6%	3645	3.0%	5.0
+	46.2%	4489	4.7%	7.6
-	46.3%	4394	4.7%	7.2
+	46.2%	4687	4.8%	7.7
-	46.3%	4440	4.8%	7.3
-	45.1%	3754	4.9%	5.0

196 Type III effectors?

TTSS status	G+C content	Number of proteins	Positives	Z-Score
-	58.4%	4777	3.3%	-1.0
-	66.9%	2157	3.6%	1.1
-	62.4%	4587	7.7%	18.9
-	65.4%	4506	7.2%	20.7
-	59.2%	1631	6.1%	6.9
-	60.2%	1998	6.2%	8.8
-	60.1%	1728	6.4%	6.2
-	59.9%	2416	5.3%	7.3
-	72.5%	3079	5.0%	12.6
-	72.4%	2940	4.0%	8.0
-	53.5%	2272	5.8%	8.5
-	63.1%	2950	9.3%	19.4
-	53.8%	3056	5.9%	9.5
-	54.1%	3076	6.4%	11.9
-	61.4%	2119	10.3%	21.1

218 Type III effectors, no T3SS?

$$0.038 \times 5169 \times 0.13 \approx 26$$

[No. +ve x  $P(T3E | +ve)$ ]

1-2% (Collmer *et al.* 2002); 1% (Boch and Bonas, 2010)



# Conclusions II

- Genomic context is almost as interesting and useful as genome features
- Functional prediction is largely sequence classification
- The statistics of large datasets can be unintuitive:  
**always try to think statistically**
- Workflows (e.g. Galaxy) are great tools

# 2013: *Dickeya* spp.



## Draft Genome Sequences of 17 Isolates of the Plant Pathogenic Bacterium *Dickeya*

Leighton Pritchard,<sup>a</sup> Sonia Humphris,<sup>b</sup> Gerry S. Saddler,<sup>c</sup> John G. Elphinstone,<sup>d</sup> Minna Pirhonen,<sup>e</sup> Ian K. Toth<sup>b</sup>

## Draft Genome Sequences of Four *Dickeya dianthicola* and Four *Dickeya solani* Strains

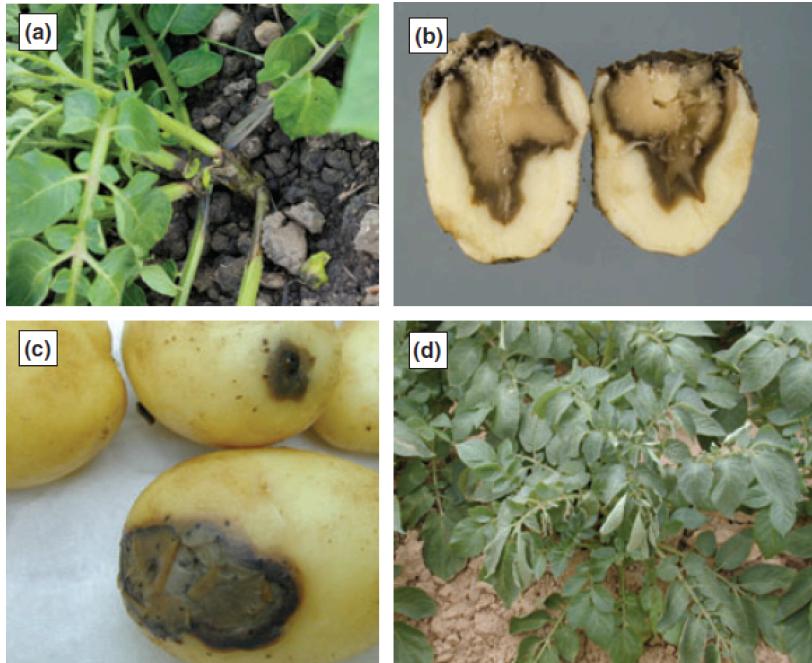
Leighton Pritchard,<sup>a</sup> Sonia Humphris,<sup>b</sup> Steve Baeten,<sup>c</sup> Martine Maes,<sup>c</sup> Johan Van Vaerenbergh,<sup>c</sup> John Elphinstone,<sup>e</sup> Gerry Saddler,<sup>d</sup> Ian Toth<sup>b</sup>

Pritchard *et al.* (2013) *Genome Announc.* [doi:10.1128/genomeA.00978-13](https://doi.org/10.1128/genomeA.00978-13).  
Pritchard *et al.* (2013) *Genome Announc.* [doi:10.1128/genomeA.00087-12](https://doi.org/10.1128/genomeA.00087-12).



# Dickeya spp.

- Soft rot pathogen of crop and ornamental plants
- *Dickeya* species have distinct host ranges
- Expanding and diversifying across Europe
- ***D. solani* a particular emerging threat**
- Promising SynthBio/BioFuels target



Toth et al. (2011) *Plant Path.* **60**: 385-399  
[doi:10.1111/j.1365-3059.2011.02427.x](https://doi.org/10.1111/j.1365-3059.2011.02427.x)



Host: chrysanthemum (*Chrysanthemum* spp.)  
Landesanst. f. Pflanzenbau und Pflanzenschutz, Mainz Archive



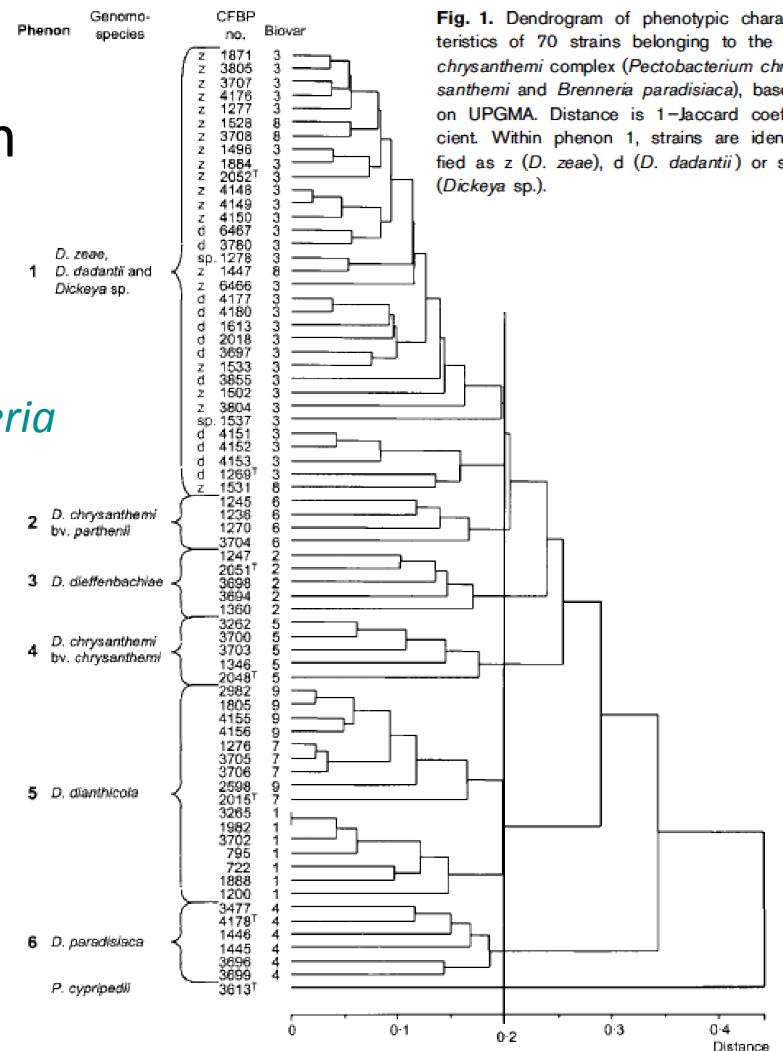
# Dickeya spp.

- 25 new bacterial genomes, spanning six species of genus
- Several short read methods, several sequencing services 454 (Liverpool), Illumina (Glasgow, Belgium), SE and PE
- 6-84X coverage, 170k-4m reads
- Genomes in draft form: 27-273 contigs
- Annotation mostly automated
- Published 2013, *Genome Announcements* (6, 8 authors)  
Pritchard *et al.* (2013) *Genome Announc.* [doi:10.1128/genomeA.00978-13](https://doi.org/10.1128/genomeA.00978-13).  
Pritchard *et al.* (2013) *Genome Announc.* [doi:10.1128/genomeA.00087-12](https://doi.org/10.1128/genomeA.00087-12).

# Dickeya classification is in flux

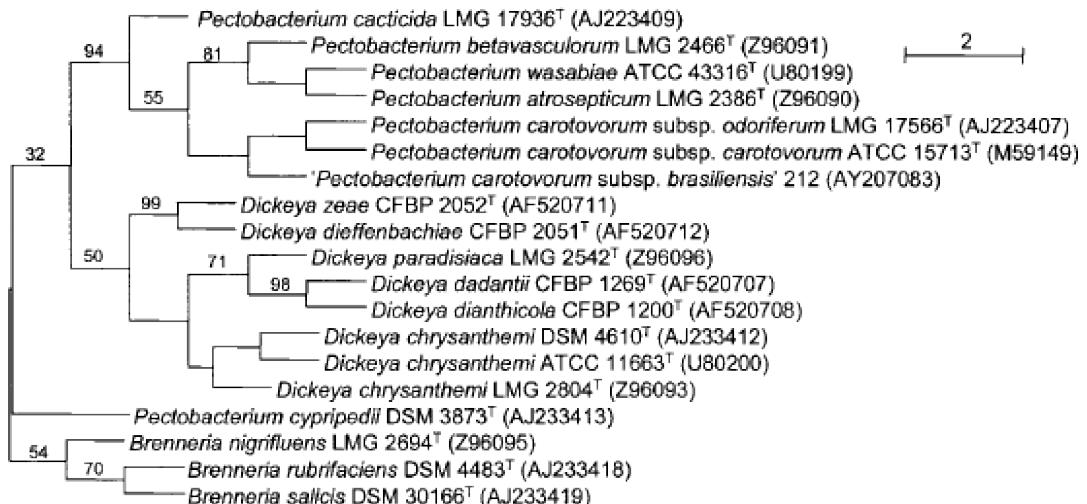
- a.k.a. *Pectobacterium chrysanthemi* or *Erwinia chrysanthemi*
- *Erwinia chrysanthemi* complex (from 1980)
  - Pathovars *chrysanthemi*, *dianthicola*, *dieffenbachiae*, *parthenii*, *zeae*, *paradisiaca*
  - (*E.chrysanthemi* pv *paradisiaca* later *Brenneria paradisiaca*)
- Classification by phenotype

Characteristic	Phenon 1	Phenon 2	Phenon 4	Phenon 3	Phenon 5	Phenon 6	<i>P. cypripedii</i>
(−)-D-Arabinose	+	−	−	+	−	+	+
(−)-D-Tartrate	−	d (25)	−	−	+	+	+
Inulin	−	−	+	−	d (88)	−	−
Lactose	+	d (75)	d (20)	−	−	d (17)	+
Growth at 39 °C	+	+	+	+	−	d (83)	−
cis-Aconitate	+	−	d (20)	d (80)	−	−	+
(+)-D-Melibiose, (+)-D-raffinose	+	+	+	−	d (44)	d (83)	−
5-Keto-D-gluconate	−	−	−	d (20)	−	+	+
Mannitol	+	+	+	+	+	−	+
Lecithin	+	+	+	+	+	−	−
ADH Moeller	d (15)	−	+	−	d (69)	−	−
meso-Tartate	+	d (75)	−	+	+	+	+
myo-Inositol	+	+	d (80)	+	+	−	+
Casein	+	d (75)	+	d (80)	d (75)	−	−
Novel species	<i>D. dadantii</i> + <i>D. zeae</i>	<i>D. chrysanthemi</i> bv. <i>parthenii</i>	<i>D. chrysanthemi</i> bv. <i>chrysanthemi</i>	<i>D. chrysanthemi</i> bv. <i>dieffenbachiae</i>	<i>D. chrysanthemi</i> bv. <i>dianthicola</i>	<i>D. chrysanthemi</i> bv. <i>paradisiaca</i>	—



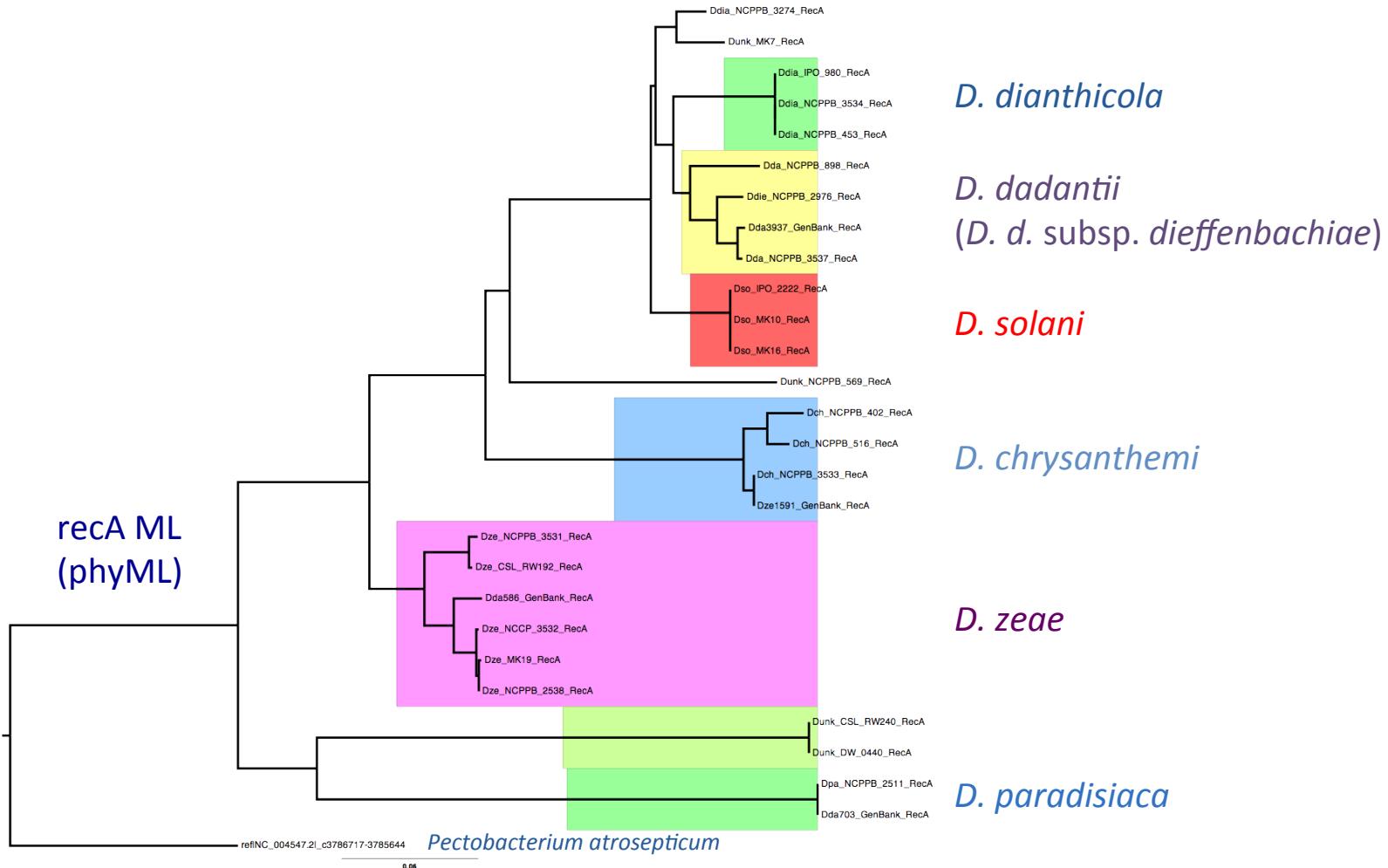
# Current classification of *Dickeya*

- Complex transferred to *Dickeya* gen. nov. with six assigned species:
  - Dickeya chrysanthemi*
  - Dickeya paradisiaca*
  - Dickeya dadantii*
  - Dickeya dianthicola*
  - Dickeya dieffenbachiae*
  - Dickeya zeae*



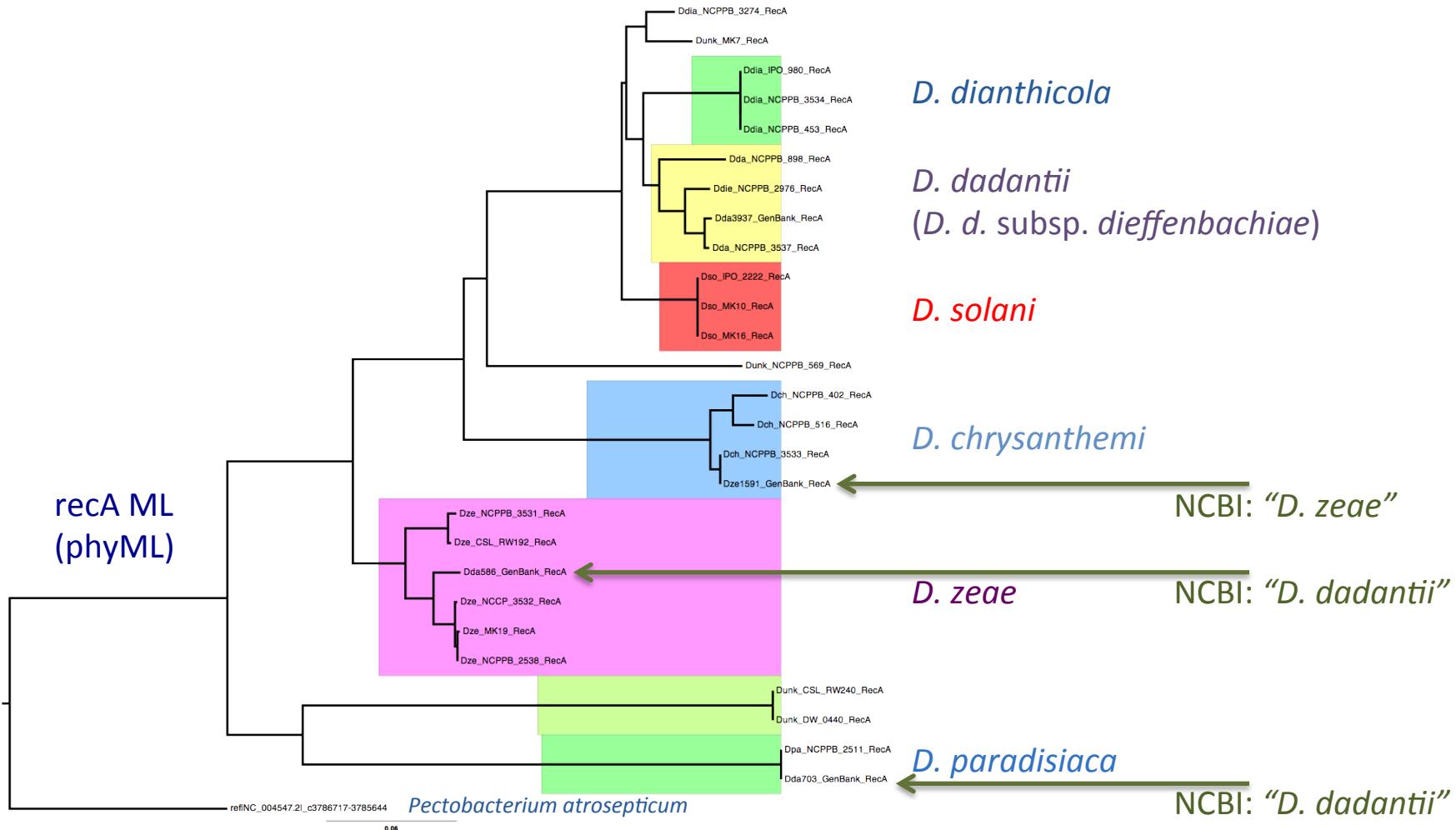
**Fig. 2.** Unrooted tree, the result of a phylogenetic analysis (parsimony shown) of 16S rRNA gene sequences of type strains of *Dickeya*, *Pectobacterium* and *Brenneria* species. Bootstrap values (expressed as percentages of 1000 replications, BIONJ+Kimura two-parameter) are indicated only for branches also retrieved by BIONJ and maximum-likelihood analyses ( $P<0.01$ ).

# Classification of sequenced *Dickeya*



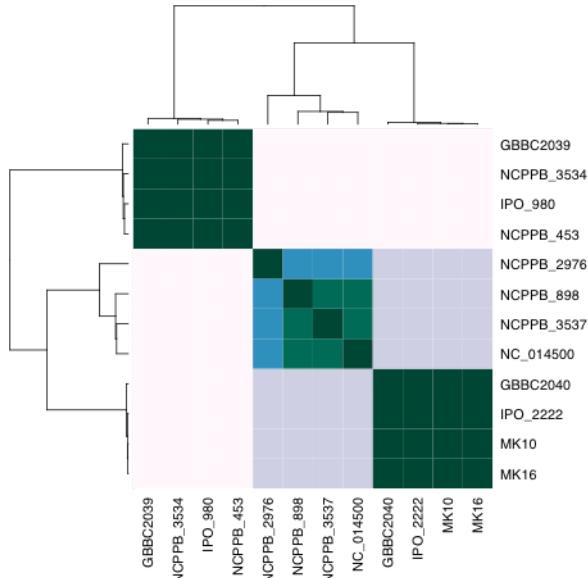
# GenBank *Dickeya* genomes misclassified

- 3/4 *Dickeya* genomes in GenBank are misassigned species



# *D. solani* recently described

- There was some question over whether *D. solani* is a species, or subspecies of *D. dadantii*.
- Species boundaries: 70% DNA-DNA hybridisation (DDH)
  - equivalent to 95% Average Nucleotide Identity (ANI), determined by whole-genome alignment



	<i>D. dianthicola</i>				<i>D. dadantii</i>				<i>D. solani</i>			
	GBBC2039	NCPPB_453	NCPPB_3534	IPO_980	NCPPB_2976	NCPPB_3537	NCPPB_898	NC_014500	GBBC2040	IPO_2222	MK10	MK16
<i>D. dianthicola</i>	1	0.9936	0.9927	0.993	0.9222	0.9219	0.9221	0.9223	0.9227	0.9229	0.923	0.9231
GBBC2039	0.9936	1	0.9947	0.9952	0.9221	0.9216	0.922	0.9219	0.9226	0.9226	0.9229	0.9231
NCPPB_3534	0.9927	0.9947	1	0.9937	0.922	0.9219	0.922	0.9221	0.9222	0.9223	0.9225	0.9227
IPO_980	0.993	0.9952	0.9937	1	0.9217	0.9214	0.9218	0.9217	0.9223	0.9224	0.9227	0.9227
<i>D. dadantii</i>	0.9222	0.9221	0.922	0.9217	1	0.9684	0.9669	0.968	0.9411	0.9413	0.9415	0.9416
NCPPB_2976	0.9219	0.9216	0.9219	0.9214	0.9684	1	0.985	0.985	0.9406	0.941	0.9412	0.9414
NCPPB_3537	0.9221	0.922	0.922	0.9218	0.9669	0.985	1	0.9838	0.94	0.9403	0.9405	0.9407
NCPPB_898	0.9223	0.9219	0.9221	0.9217	0.968	0.985	0.9838	1	0.9408	0.9413	0.9415	0.9417
NC_014500												
<i>D. solani</i>	0.9227	0.9226	0.9222	0.9223	0.9411	0.9406	0.94	0.9408	1	0.9987	0.9983	0.9983
GBBC2040	0.9229	0.9226	0.9223	0.9224	0.9413	0.941	0.9403	0.9413	0.9987	1	0.9994	0.9995
IPO_2222												
MK10	0.923	0.9229	0.9225	0.9227	0.9415	0.9412	0.9405	0.9415	0.9983	0.9994	1	0.9997
MK16	0.9231	0.9231	0.9227	0.9227	0.9416	0.9414	0.9407	0.9417	0.9983	0.9995	0.9997	1

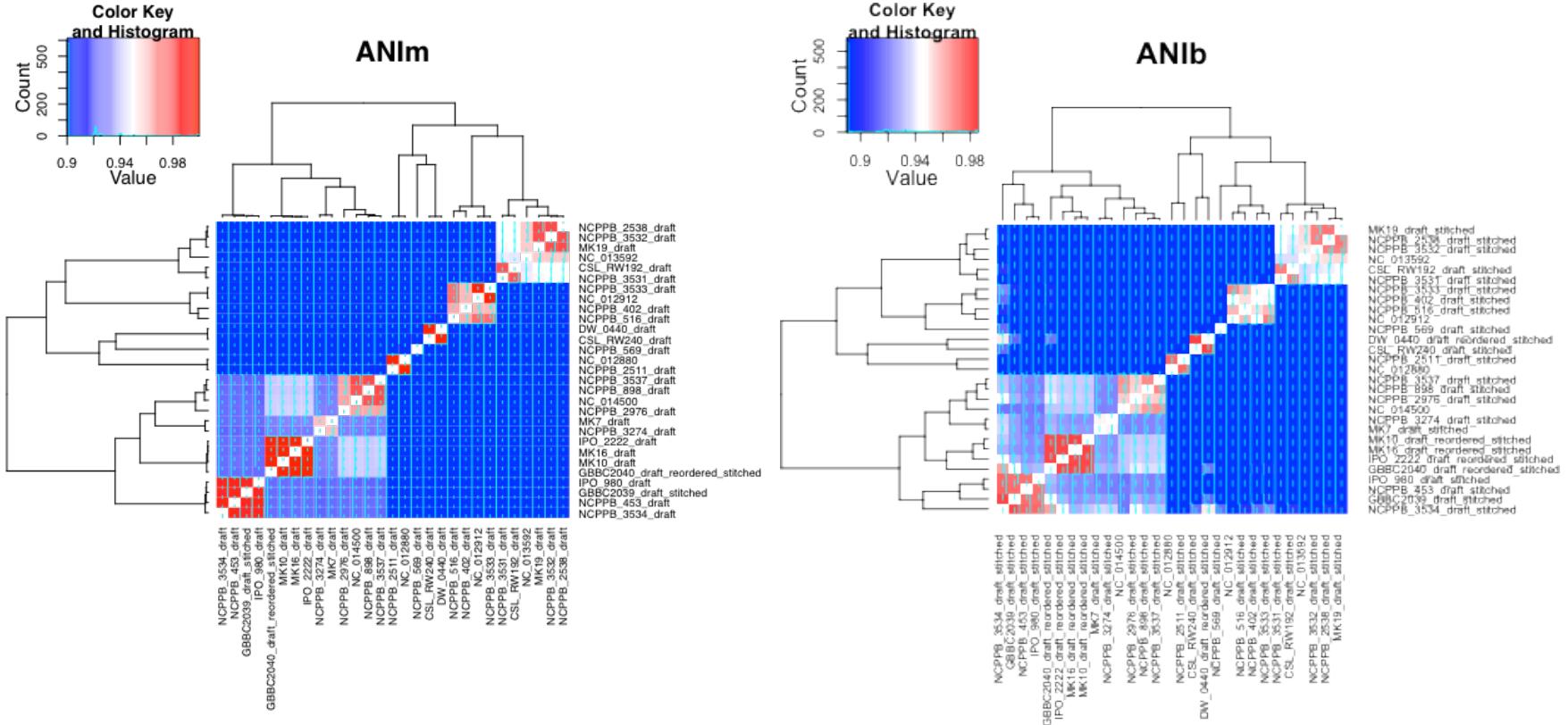
International Journal of Systematic and Evolutionary Microbiology (2014), 64, 768–774

DOI 10.1099/ijss.0.052944-0

*Dickeya solani* sp. nov., a pectinolytic plant-pathogenic bacterium isolated from potato (*Solanum tuberosum*)

Jan M. van der Wolf,<sup>1</sup> Els H. Nijhuis,<sup>1</sup> Małgorzata J. Kowalewska,<sup>2</sup> Gerry S. Saddler,<sup>2</sup> Neil Parkinson,<sup>3</sup> John G. Elphinstone,<sup>3</sup> Leighton Pritchard,<sup>4</sup> Ian K. Toth,<sup>4</sup> Ewa Lojkowska,<sup>5</sup> Marta Potrykus,<sup>5</sup> Małgorzata Waleron,<sup>5</sup> Paul de Vos,<sup>6</sup> Ilse Cleenwerck,<sup>6</sup> Minna Pirhonen,<sup>7</sup> Linda Garlant,<sup>7</sup> Valérie Hélias,<sup>8</sup> Joël F. Pothier,<sup>9,10</sup> Valentin Pflüger,<sup>11</sup> Brion Duffy,<sup>10</sup> Leah Tsror<sup>12</sup> and Shula Manulis<sup>12</sup>

# Nine *Dickeya* species?



- Availability of a genome for your pathogen may reveal some taxonomic surprises...  
(*Rhodococcus*...)



# Dickeya Genomic Diagnostics

- To legislate with respect to, and quarantine materials infected with, a pathogen we must be able to identify it accurately and precisely.

 The Scottish Government

science and Advice for Scottish Agriculture (SASA)

SASA HOME ABOUT US TOPICS DIAGNOSTICS RAD RESOURCES TRAINING STAFF

Home > Diagnostics > Bacteriology > Dickeya

## Dickeya

Potato Industry | Potatoes | Bacteria | Bacteriology



**Top wilt/ blackleg/ soft rot – *Dickeya* spp.**

The genus *Dickeya* was previously known as *Erwinia chrysanthemi*. There are two significant potato pathogens within this genus currently affecting Europe: *Dickeya dianthicola* and *Dickeya solani*. *D. dianthicola* has never been found on Scottish potatoes. *D. solani* emerged in Europe around 2005-2006. It is highly aggressive on potato, causing rapid wilting and blackleg-like symptoms across wide environmental conditions.

Preliminary research indicates that *D. solani* may be significantly more aggressive than *D. dianthicola* and *Pectobacterium* spp. and appears to be able to easily induce blackleg-like symptoms and to colonise sprouting potato tubers, even when temperatures are low. Aggressiveness of the new *Dickeya* pathogen appears to further increase at higher temperatures so there are implications for increased importance of this pathogen in response to global warming. As yet, there is little substantiated practical information on the biology of this strain in relation to its host range, its ability to survive, establish and spread in the environment or its behaviour on stored potato tubers.

*'D. solani'* has only ever been found on small number of Scottish ware crops (15 in total, from 2009 & 2010), all of which were grown from non-Scottish seed.

For further information see the [Dickeya solani - a threat to our potatoes](#) page on the Scottish Government website.

Information on a range of measures to protect Scotland's ware growers against Dickeya, and a range of threats, can be found in the [Defending your potato crop against disease leaflet](#).

View recent research from the 2010-2011 Scottish Potato Tuber Survey and 2011 inspection of growing crops of seed and ware potatoes.

Pectobacterium up Entomology

Printer-friendly version

**QUICK LINKS**

- Community Outreach
- Freedom of Information
- Job Vacancies
- Quality Assurance
- SASA in the media

SEARCH  Search

TEXT SIZE + - Current Size: 100%

**MySASA**

> MySPUDS

**NAVIGATION**

- Site A-Z
- Site map
- How to use our site
- Scientific Publications

**CONTACT US**

SASA Rockinglaw Road Edinburgh EH12 9JU

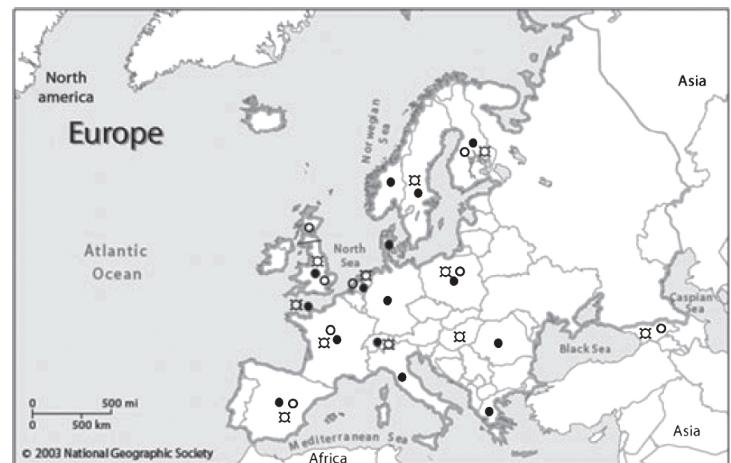
T: +44(0)131 244 8890  
F: +44(0)131 244 8940  
E: [info@sasa.gsi.gov.uk](mailto:info@sasa.gsi.gov.uk)  
W: [www.sasa.gov.uk](http://www.sasa.gov.uk)

- Freedom of Information requests  
- Map of our location  
- Customer feedback  
- Website feedback



# *Dickeya* spp. are a threat in Europe

- *D. dianthicola* an established pathogen across Europe
- ***D. solani* an emerging threat**
  - Increased virulence
  - Different temperature preference
- Must be able to distinguish and identify, to legislate and/or quarantine
  - Zero-tolerance in Scotland
- No qPCR primers existed to distinguish among *Dickeya* spp.
- Was asked to find a solution, starting from our draft quality genome sequences





# qPCR primer design

- Typical approaches use *rational design*
  - Select region that is (thought to be):
    - Similar in target organisms
    - Divergent or absent in off-target organisms
    - Conserved enough to be amplified by single primer set
  - Design primers to amplify that region
  - Frequently-used regions: intergenic transcribed spacers (ITS), ribosomal DNA, “housekeeping” genes, “virulence” genes.
- Alignment-Based Process:
  1. Identify common region(s) (“alignment”)
  2. Design primers

**Availability of whole pathogen genomes gives new options...**

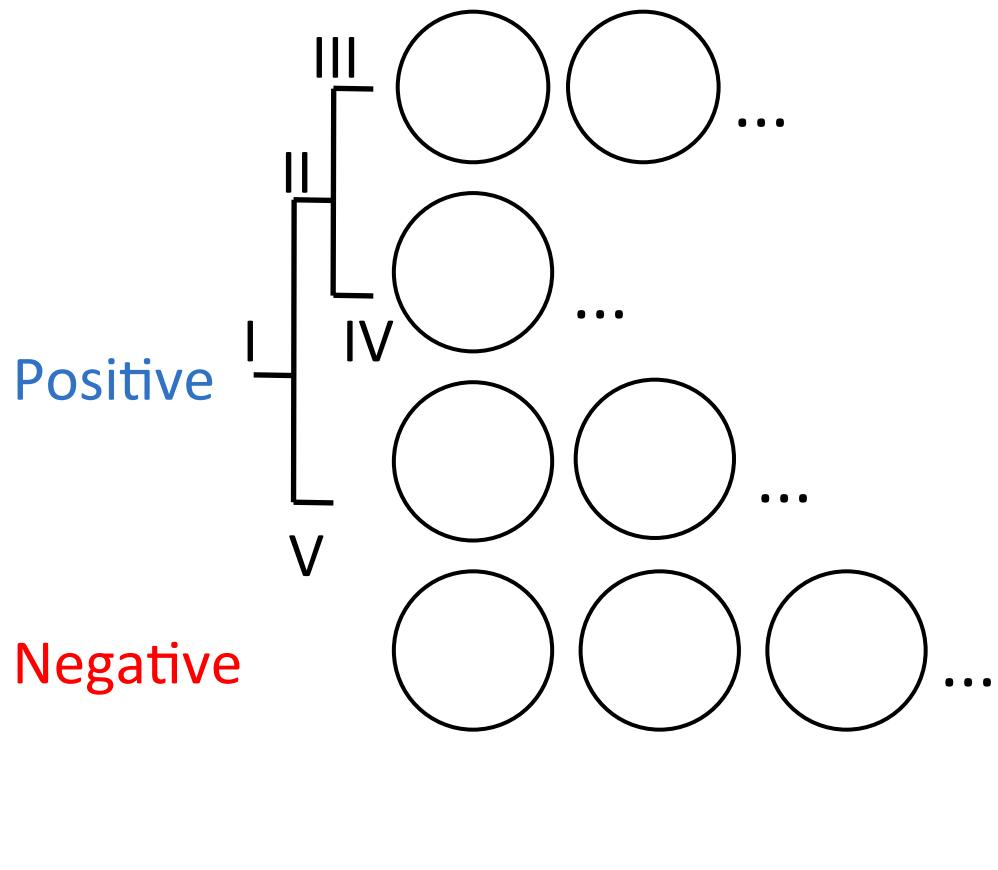


# “Alignment-free” Genomic PCR Design

1. Define **positive**, **negative** sets, and **classes**

2. Convert to single (pseudo)chromosomes:

3. Genome feature locations from GBK file or predicted



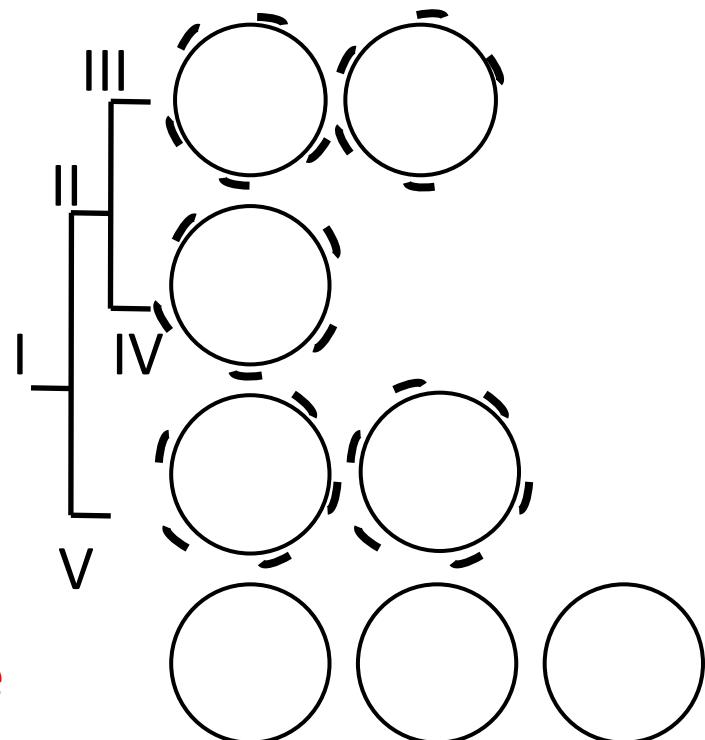


# Primer prediction (on positive set)

4. Predict **> 1000**  
thermodynamically plausible  
primer sets on each  
(pseudo)chromosome

Positive

Negative



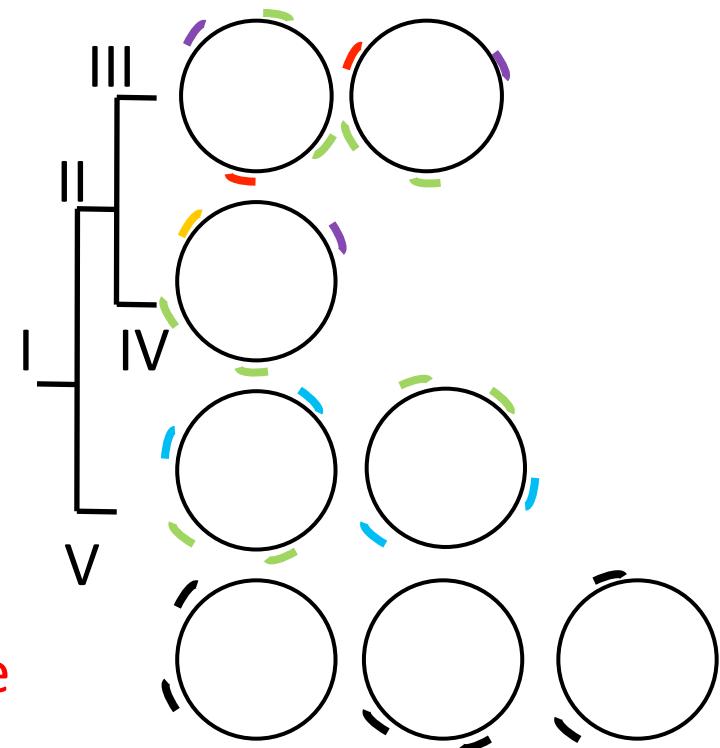
# Test cross-amplification *in silico*

5. Check cross-amplification:  
All primer sets tested against  
other organisms, colour by  
amplified class(es).

6. BLAST screen:  
All primers screened for off-  
target sequences with BLAST

Positive

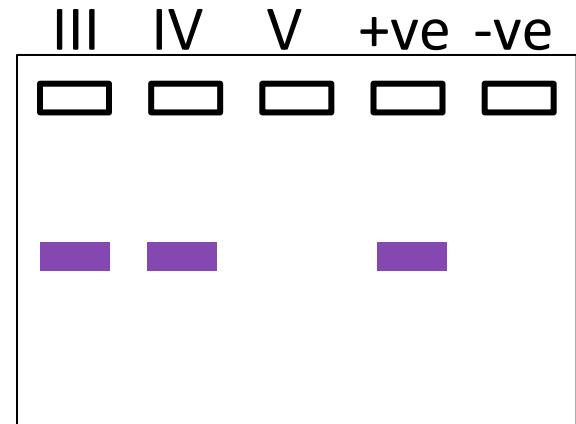
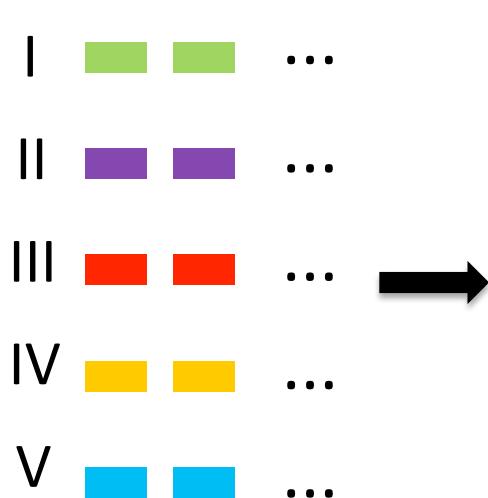
Negative



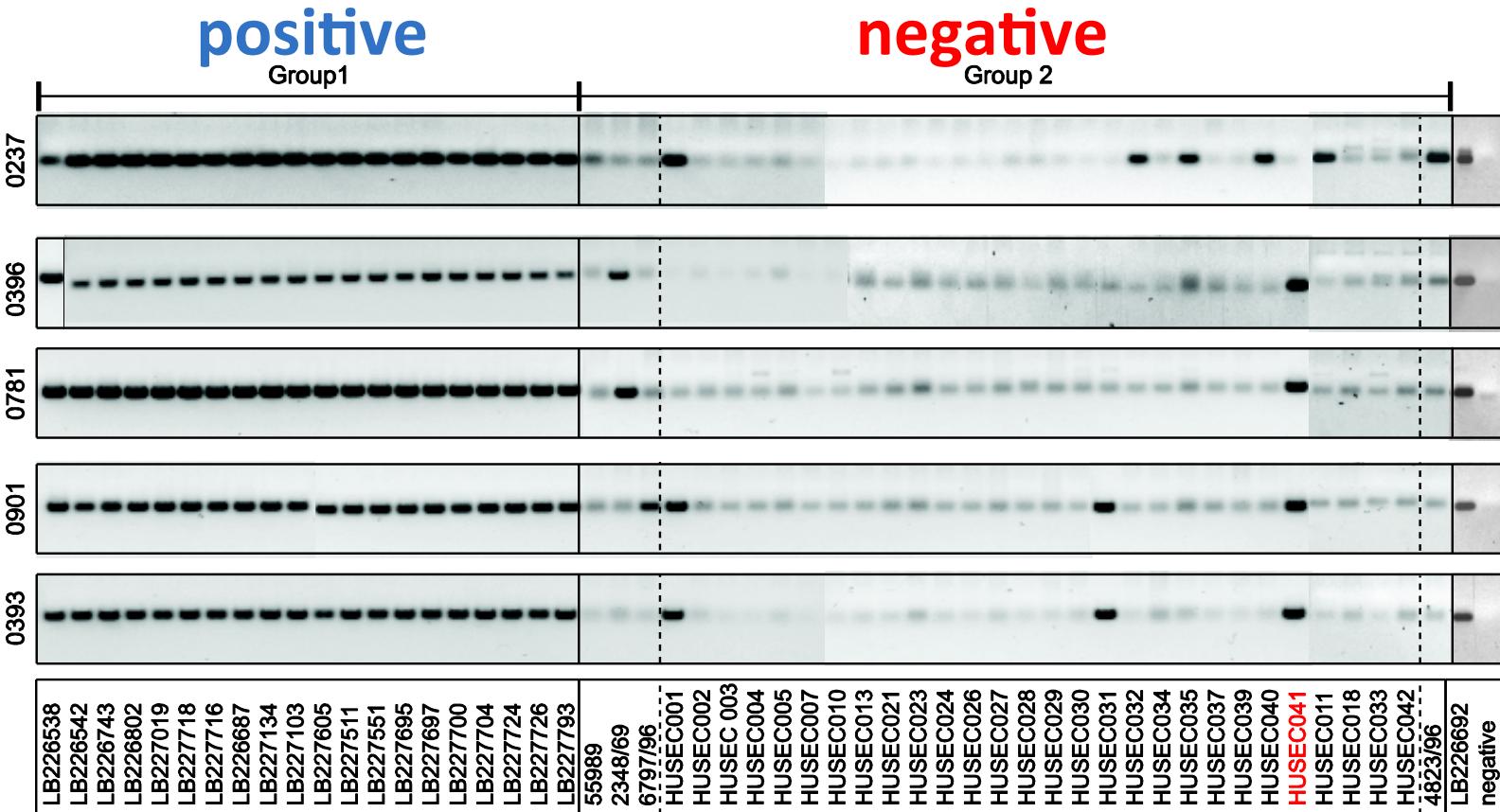
# Classify primers and validation

7. Classify primers:  
according to the ability  
to amplify specific  
classes of input  
sequence.

8. Validate primers:  
Primer set validated on  
positive and negative  
targets *in vitro*.



# Primer design: validated *in vitro* (*E.coli*)



**Sub-serotype specificity:**  
*impB* ( $\beta$ -lactamase plasmid),  
prophage gp20 transfer protein

OPEN  ACCESS Freely available online

**Alignment-Free Design of Highly Discriminatory Diagnostic Primer Sets for *Escherichia coli* O104:H4 Outbreak Strains**

Leighton Pritchard<sup>1\*</sup>, Nicola J. Holden<sup>2\*</sup>, Martina Bielaszewska<sup>3</sup>, Helge Karch<sup>3</sup>, Ian K. Toth<sup>2</sup>



# Primer design: validated *in vitro* (*Dickeya*)

Table 5 Specificity of real-time qPCR assays predicted by the pipeline compared with existing qPCR and conventional PCR assays

Test species	Isolates tested	Number of isolates detected, by assay								
		DIC	DIA-A	DIA-C	SOL-C	SOL-D	PEC	ECH	ECA	ADE
<i>D. dianthicola</i>	7	7	7	7	0	0	7	7	0	7
' <i>D. solani</i> ' (DUC-1)	16	16	0	0	16	16	16	16	0	16
DUC-2	5	5	0	0	0	0	5	5	0	5
DUC-3	1	1	0	0	0	0	1	1	0	1
<i>D. dadantii</i>	11	11	0	0	1	1	11	11	0	11
<i>D. dieffenbachiae</i>	6	6	0	0	0	0	6	6	0	6
<i>D. chrysanthemi</i> bv. <i>chrysanthemi</i>	7	7	0	0	0	0	7	7	0	7
<i>D. chrysanthemi</i> bv. <i>parthenii</i>	3	3	0	0	0	0	3	3	0	3
<i>D. paradisiaca</i>	1	1	0	0	0	0	1	1	0	1
<i>D. zeae</i>	11	11	0	0	0	0	11	11	0	11
New <i>Dickeya</i> species level clade (I)	1	1	0	0	0	0	1	1	0	1
New <i>Dickeya</i> species level clade (II)	1	1	0	0	0	0	1	1	0	1
<i>Pectobacterium atrosepticum</i>	1	1	0	0	0	0	1	0	1	0
<i>P. carotovorum</i> subsp. <i>carotovorum</i>	1	1	0	0	0	0	1	0	0	0
<i>P. betavasculorum</i>	1	1	0	0	0	0	1	0	0	0
<i>P. carotovorum</i> subsp. <i>odoriferum</i>	1	1	0	0	0	0	1	0	0	0
<i>P. wasabiei</i>	1	1	0	0	0	0	1	0	0	0
<i>Pantoea agglomerans</i>	1	1	0	0	0	0	0	0	0	0
<i>Brenneria quercina</i>	1	1	0	0	0	0	0	0	0	0
<i>Erwinia amylovora</i>	1	1	0	0	0	0	0	0	0	0

Primers tested against 70 isolates



Plant Pathology (2012)

Doi: 10.1111/j.1365-3059.2012.02678.x

## Detection of phytopathogens of the genus *Dickeya* using a PCR primer prediction pipeline for draft bacterial genome sequences

L. Pritchard<sup>a</sup>, S. Humphris<sup>a</sup>, G. S. Saddler<sup>b</sup>, N. M. Parkinson<sup>c</sup>, V. Bertrand<sup>c</sup>, J. G. Elphinstone<sup>c</sup> and I. K. Toth<sup>a\*</sup>

# Dickeya genecalling

- **No free lunch for gene-calling**
- Tested several automated methods against “gold standard” annotations
- **Systematic errors:** *rpmDHJ*, *prfB*, *outE*, *secD*, *hrpD*, *virB7*, etc.

Manual annotation, *Pectobacterium*

Genecaller:	Glimmer3	GeneMark	Prodigal (SD)	Prodigal (motif)	RAST
Prediction count	4679	4606	4467	4468	4648
Absolutely missed CDS	112	133	156	153	202

**strict assessment (start and end base match)**

Sensitivity	0.73	0.75	0.86	0.85	0.79
False discovery rate (FDR)	0.31	0.27	0.14	0.14	0.24
Positive predictive value (PPV)	0.69	0.73	0.86	0.86	0.76

**lenient assessment (end base matches)**

Sensitivity	0.97	0.97	0.97	0.97	0.95
False discovery rate (FDR)	0.07	0.06	0.03	0.03	0.08
Positive predictive value (PPV)	0.93	0.94	0.97	0.97	0.92

Community annotation, *Dickeya*

Genecaller:	Glimmer3	GeneMark	Prodigal (SD)	Prodigal (motif)	RAST
Prediction count	4752	4440	4287	4291	4604
Absolutely missed CDS	284	393	407	408	461

**strict assessment (start and end base match)**

Sensitivity	0.62	0.59	0.71	0.71	0.64
False discovery rate (FDR)	0.41	0.39	0.25	0.25	0.37
Positive predictive value (PPV)	0.59	0.61	0.75	0.75	0.63

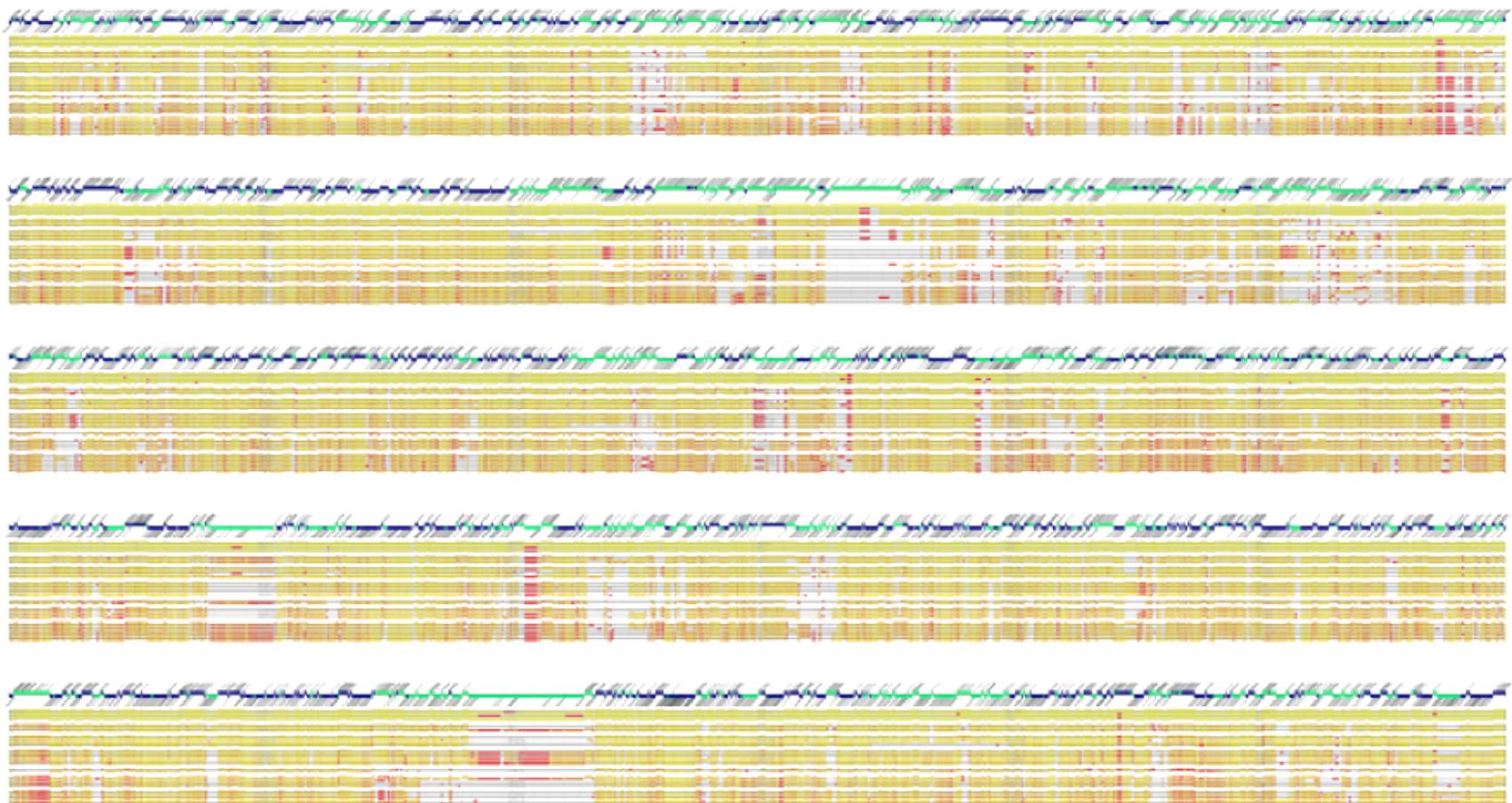
**lenient assessment (end base matches)**

Sensitivity	0.94	0.91	0.91	0.91	0.90
False discovery rate (FDR)	0.10	0.06	0.03	0.03	0.09
Positive predictive value (PPV)	0.90	0.94	0.97	0.97	0.91



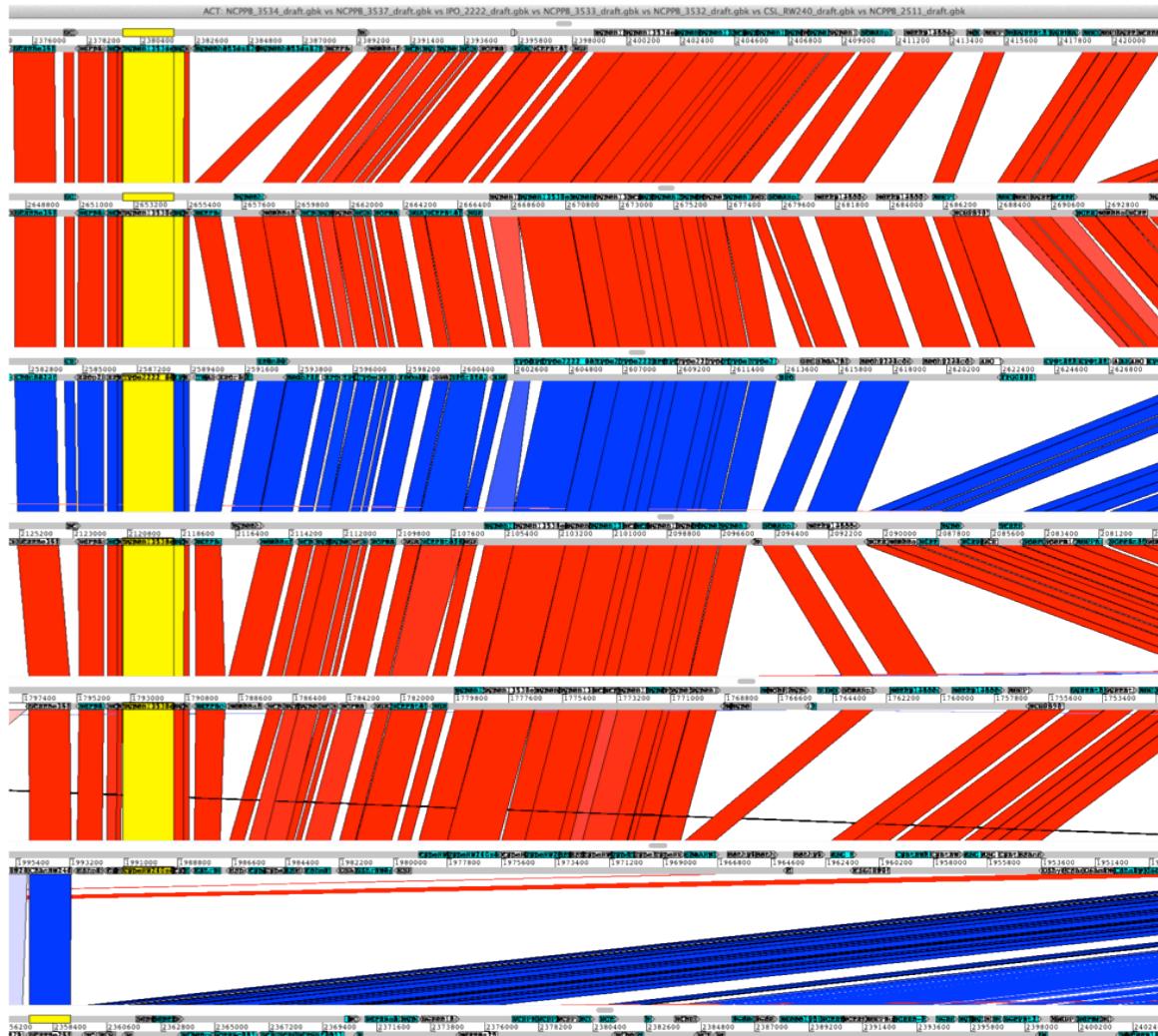
# Dickeya RBH

- Comparison of predicted gene complements
- Identifies sequence-similar ‘equivalence’, *not orthology*

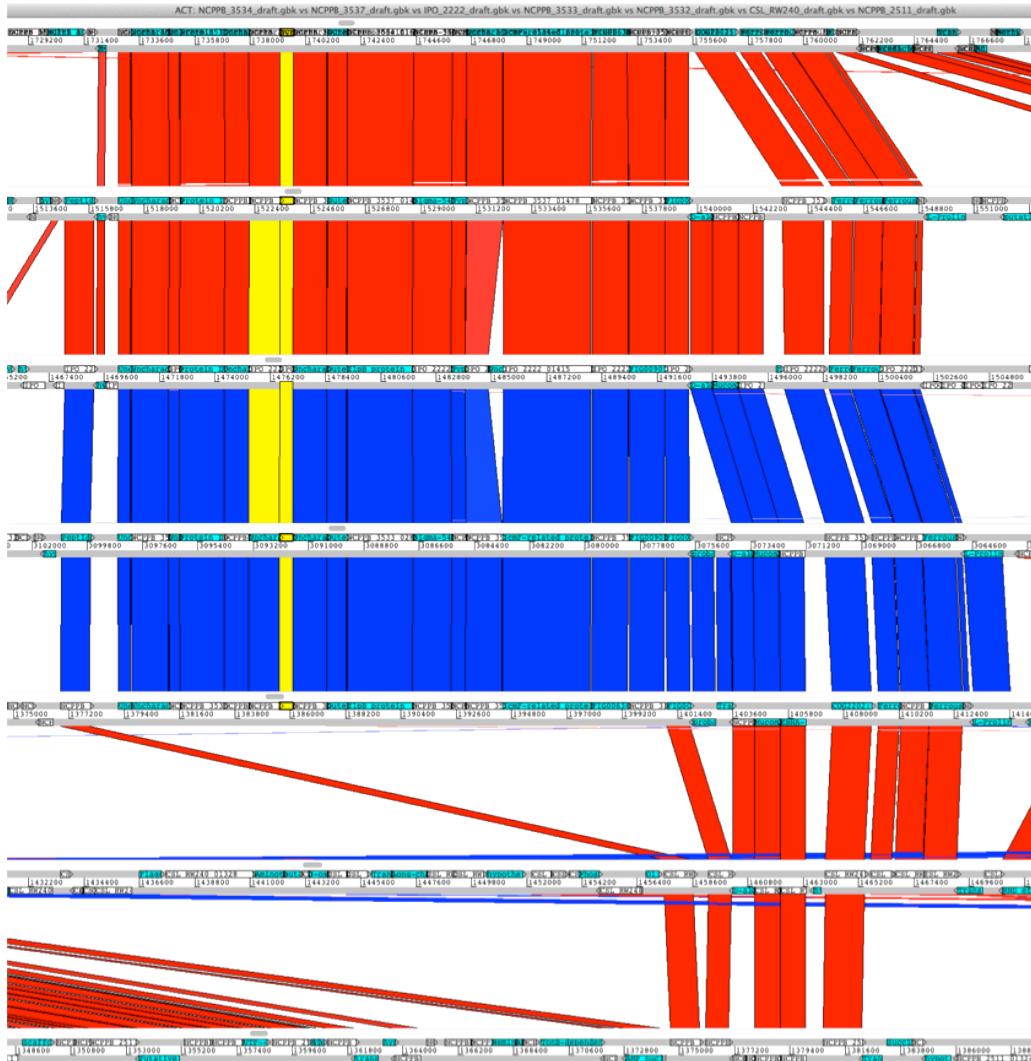




# Type III Secretion System

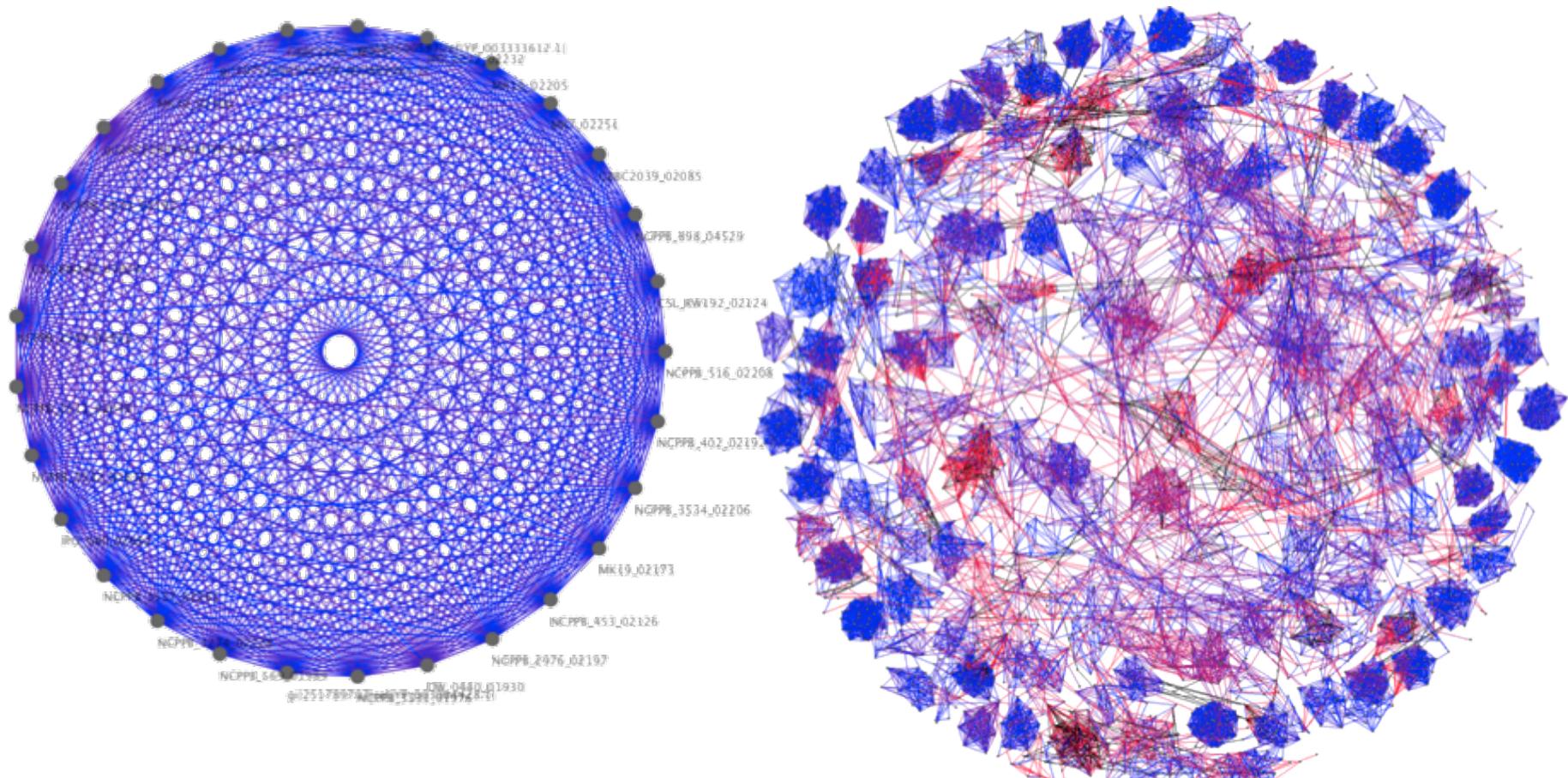


# Type VI Secretion System



# Pangenome identification

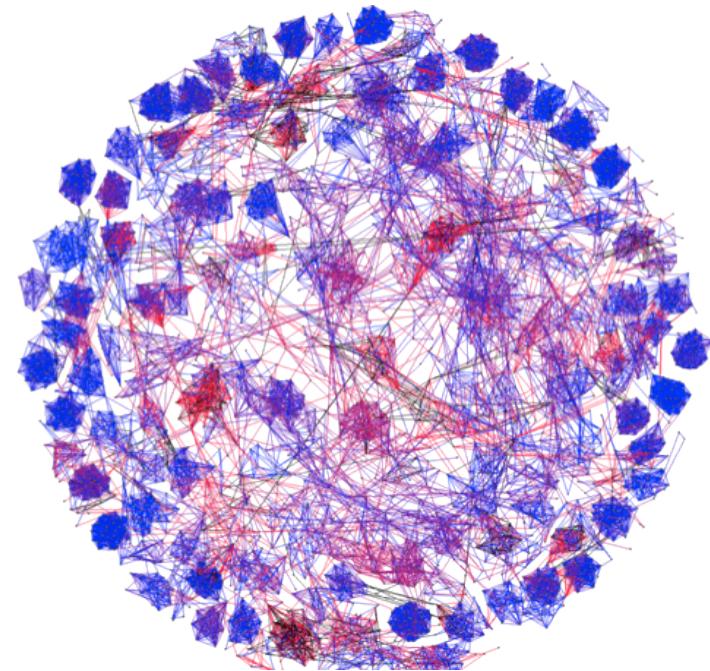
- **Clique:** each CDS makes RBH with every other member of the same *clique*



# Pangenome identification

- Prune non-clique graphs to identify accessory genomes for each species

Species	Weak Pruning	Full Pruning
Core Genome	2201	2201
<i>D. chrysanthemi</i>	32	36
<i>D. dadantii</i>	11	14
<i>D. dianthicola</i>	102	127
<i>D. paradisiaca</i>	404	441
<i>D. solani</i>	120	157
<i>D. zeae</i>	33	40



- Accessory: RBBH with all other members of same species, but no other *Dickeya*
- Weak pruning: remove all RBBH <80% identity, <40% coverage
- Full pruning: trim graph (by Mahalanobis distance) until minimal cliques found



# Collinearity

- Using conserved order of RBH pairs to identify collinear regions
  - i-ADHoRe (<http://bioinformatics.psb.ugent.be/software/details/i--ADHoRe>)

29 *Dickeya* genomes: core genome



*Dickeya zae*: core and accessory genome



Simillion *et al.* (2008) *Bioinformatics* [doi:10.1093/bioinformatics/btm449](https://doi.org/10.1093/bioinformatics/btm449)  
<http://armchairbiology.blogspot.co.uk/2012/09/the-colours-man-colours.html>  
<https://github.com/widdowquinn/pyADHoRe>

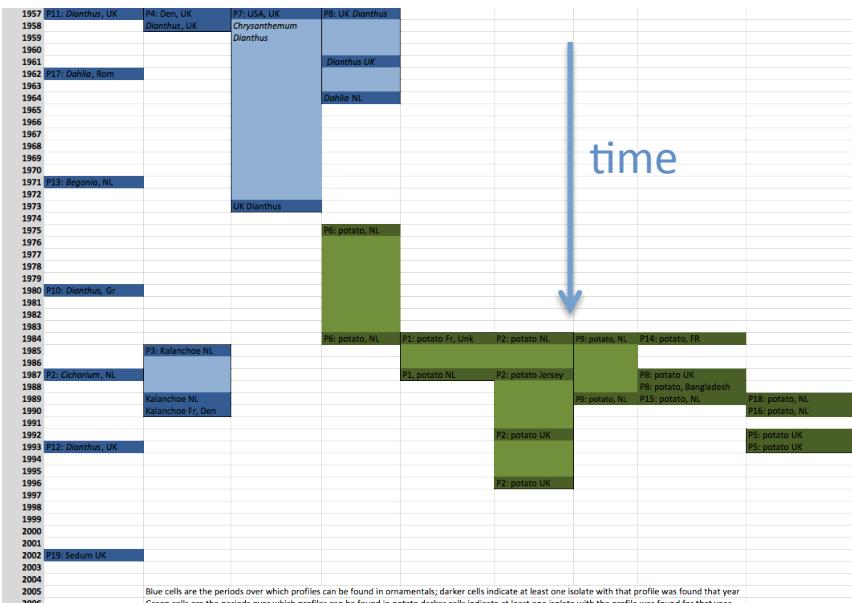


# Epidemiology

- *Dickeya infections are well-controlled on potato imports. But Dickeya has a broad host range. How did/does it get into Europe?*

# Variable nucleotide tandem repeats (VNTR)

- Identified in sequenced *D. dianthicola* genomes
- variable repeat numbers in historical isolates

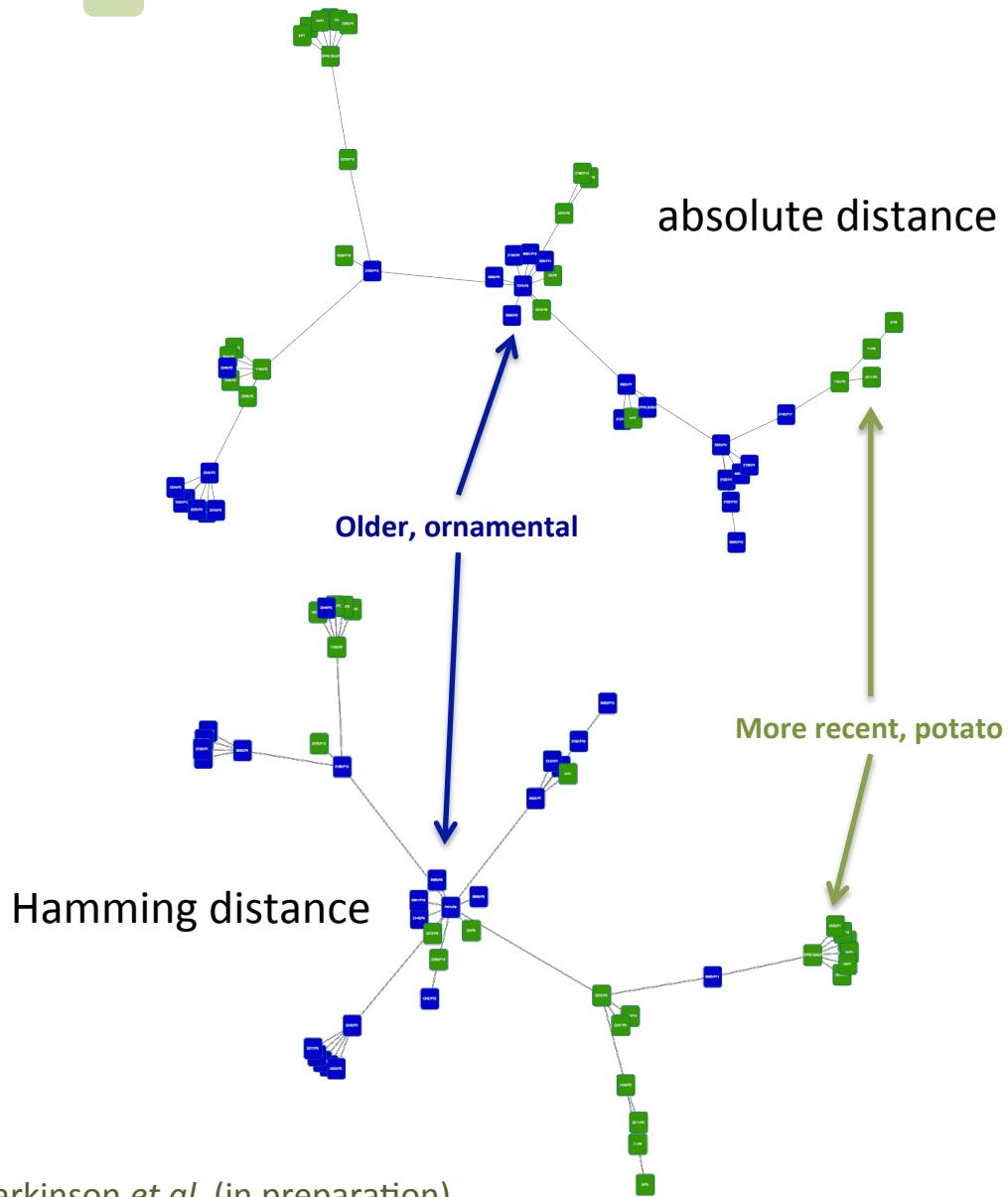


Blue: ornamental  
Green: potato

Culture collection ref	Fera ref	Host	Country	Accession Date	19	19b	29	10	11b	28	VNTR Profile
PD 502	6	Potato	Unknown	Pre 1987	4	5	6	5	5	6	P1
PD 769	10	Potato	Unknown	Pre 1987	4	5	6	5	5	6	
PD 767	15	Potato	Unknown	Pre 1987	4	5	6	5	5	6	
NCPPB 3534	2205	Potato	Netherlands	1987	4	5	6	5	5	6	
PD 1022	2248	Potato	Netherlands	1987	4	5	6	5	5	6	
NCPPB 3345		Potato	France	1984	4	5	6	5	5	6	
NCPPB 3530	2209	Potato	Jersey	1987	3	4	7	5	3	5	P2
	1644	Potato	UK	1996	3	4	7	5	3	5	
	1649	Potato	UK	1996	3	4	7	5	3	5	
PD 484	2243	Potato	Netherlands	1984	3	4	7	5	3	5	
PD 788	2246	<i>Cichorium intybus</i>	Netherlands	1987	3	4	7	5	3	5	
	1106	Potato	UK	1992	3	4	7	5	3	5	
PD 554	1234	<i>Kalanchoe blossfeldiana</i>	Netherlands	1985	3	4	9	9	2	3	P3
PD 593	2245	<i>Kalanchoe blossfeldiana</i>	Netherlands	1985	3	4	9	9	2	3	
PD 1325	2250	<i>Kalanchoe</i> sp.	Netherlands	1989	3	4	9	9	2	3	
PD 1343	2251	<i>Kalanchoe</i> sp.	Netherlands	1989	3	4	9	9	2	3	
NCPPB 3730	2216	<i>Kalanchoe blossfeldiana</i>	Denmark	1990	3	4	9	9	2	3	
NCPPB 3729	2217	<i>Kalanchoe blossfeldiana</i>	France	1990	3	4	9	9	2	3	
NCPPB 518	2132	<i>Dianthus caryophyllus</i>	Denmark	1957	3	4	3	2	3	5	P4
NCPPB 393	2138	<i>Dianthus caryophyllus</i>	UK	1957	3	4	3	2	3	5	
NCPPB 429	6953	<i>Dianthus caryophyllus</i>	UK	1957	3	4	3	2	3	5	
NCPPB 430	6954	<i>Dianthus caryophyllus</i>	UK	1958	3	4	3	2	3	5	
	1104	Potato	UK	1992	3	2	3	2	3	3	
NCPPB 3881	2211	Potato	UK	1993	3	2	3	2	3	3	
PD 771	5	Potato	Netherlands	1984	1	2	3	2	3	3	P5
PD 256	11	Potato	Netherlands	1975	1	2	3	2	3	3	
CPPB 394	2137	<i>Chrysanthemum morifolium</i>	USA	1957	3	4	3	3	3	6	P7
NCPPB 426	6952	<i>Dianthus caryophyllus</i>	UK	1957	3	4	3	3	3	6	
NCPPB 2536		<i>Dianthus caryophyllus</i>	UK	1973	3	4	3	3	3	6	
PD 765	4	Potato	Unknown	Pre 1987	3	4	3	3	3	6	
NCPPB 1111	2146	<i>Dianthus caryophyllus</i>	UK	1961	3	4	3	4	3	6	P8
NCPPB 453	7074	<i>Dianthus caryophyllus</i>	UK	1957	3	4	3	4	3	6	
NCPPB 1609	6968	<i>Dahlia</i> sp.	Netherlands	1964	3	4	3	4	3	6	
NCPPB 1956	6969	<i>Dahlia</i> sp.	Netherlands	1964	3	4	3	4	3	6	
	2272	Potato	Bangladesh	1988	3	4	3	4	3	6	
NCPPB 3529	23	Potato	UK	1987	3	4	3	4	3	6	
PD 1406	2247	Potato	Netherlands	1989	3	2	3	5	3	6	P9
PD 482	2275	Potato	Netherlands	1984	3	2	3	5	3	6	
NCPPB 3138	2182	<i>Dianthus caryophyllus</i>	Greece	1980	4	4	3	3	4	4	P10
NCPPB 3139	6983	<i>Dianthus caryophyllus</i>	Greece	1980	4	4	3	3	4	4	
NCPPB 452	6955	<i>Dianthus caryophyllus</i>	UK	1957	3	5	3	5	3	6	P11
	1242	<i>Dianthus</i> sp.	UK	1993	3	3	9	8	3	2	
NCPPB 2421	2166	<i>Begonia bertainii</i>	Netherlands	1971	3	4	3	6	3	5	P13
NCPPB 3344	2198	Potato	France	1984	3	2	2	6	3	6	P14
	2279	Potato	Netherlands	1989	4	4	3	5	5	5	P15
NCPPB 3710	2220	Potato	Netherlands	1990	4	5	4	5	6	6	P16
NCPPB 1385	2145	<i>Dahlia</i> sp.	Romania	1962	3	3	2	3	5	6	P17
PD 1405	2269	Potato	Netherlands	1989	3	4	3	8	3	6	P18
NCPPB 4305	6991	Sedum	UK	2002	3	4	3	4	6	6	P19



# VNTR Minimum Spanning Trees

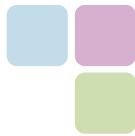


- Central nodes are associated with:
  - older accessions
  - ornamental accessions
- Recent accessions are associated with distal nodes
- Potato accessions are associated with distal nodes



# VNTR suggests ornamental introduction

- Earliest strains identified (1957) show VNTR diversity
- No *Dickeya* infection of potato observed until 1975
- Ornamental isolate diversity is greater than potato isolate diversity
- Cross-infection of ornamentals and potato supported by three VNTR profiles, found in both host types
- There is strict prohibition and regulation of potato import into the EU; there is not equivalent regulation of ornamental imports
- Ornamental plants may have provided a route for introduction of the broad host range pathogen *Dickeya* to EU potato crops.



# Systems Biology

- *Dickeya* spp. have a wide host range, and a broad set of phenotypes. Can we infer organism-level properties from sequenced genomes, and gain insight into pathogenicity and host range?



# Mapping Genome Features to Metabolism

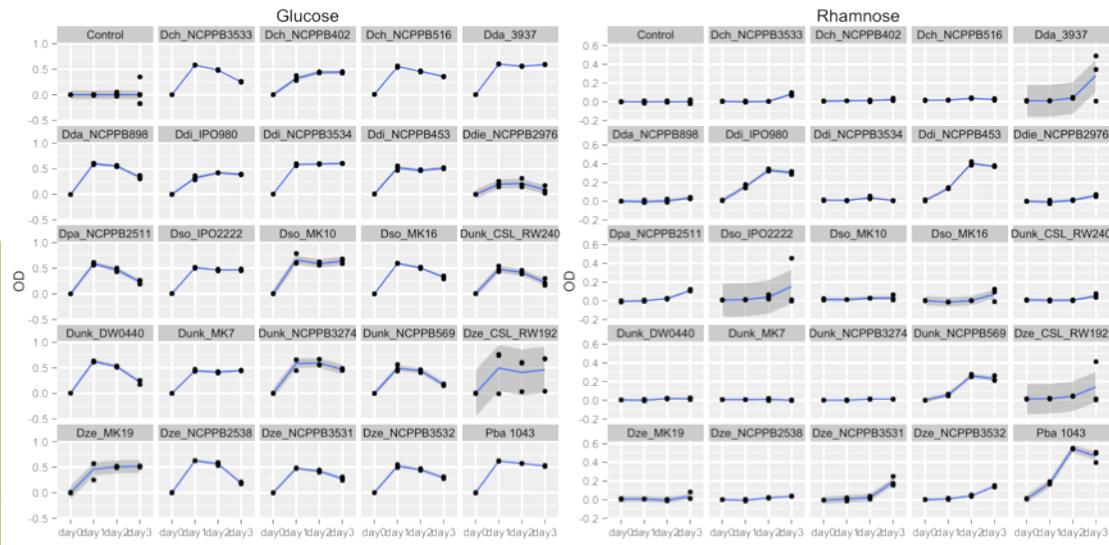
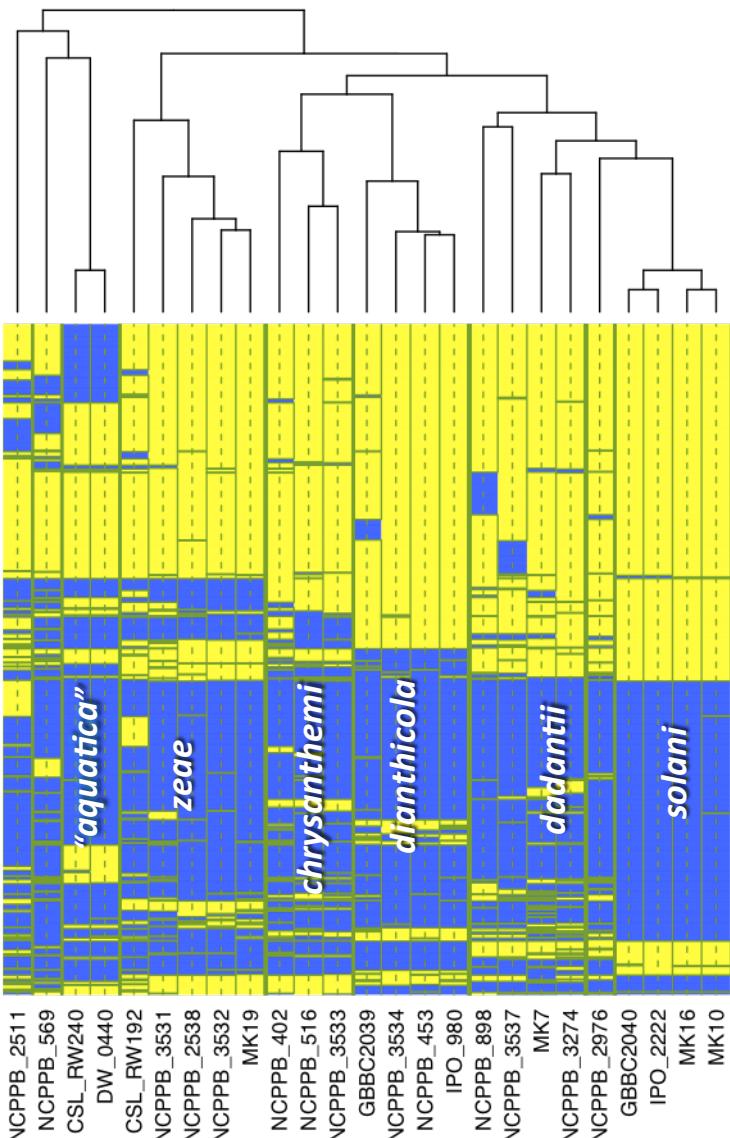
- Many annotation tools associate genome features with metabolic roles (but not all metabolic roles have known genes...):

- BioCyc/EcoCyc/PathwayTools  
(<http://biocyc.org/>)
- KEGG/KAAS  
(<http://www.genome.jp/kegg/kaas/>)
- CanOE  
(<http://www.genoscope.cns.fr/agc/microscope/metabolism/canoe.php>)
- UniPathway  
(<http://www.unipathway.org/>)

(this, again, is sequence classification)

- Having genome sequence also makes finding regulatory regions easier
  - potential for reconstructing regulatory networks

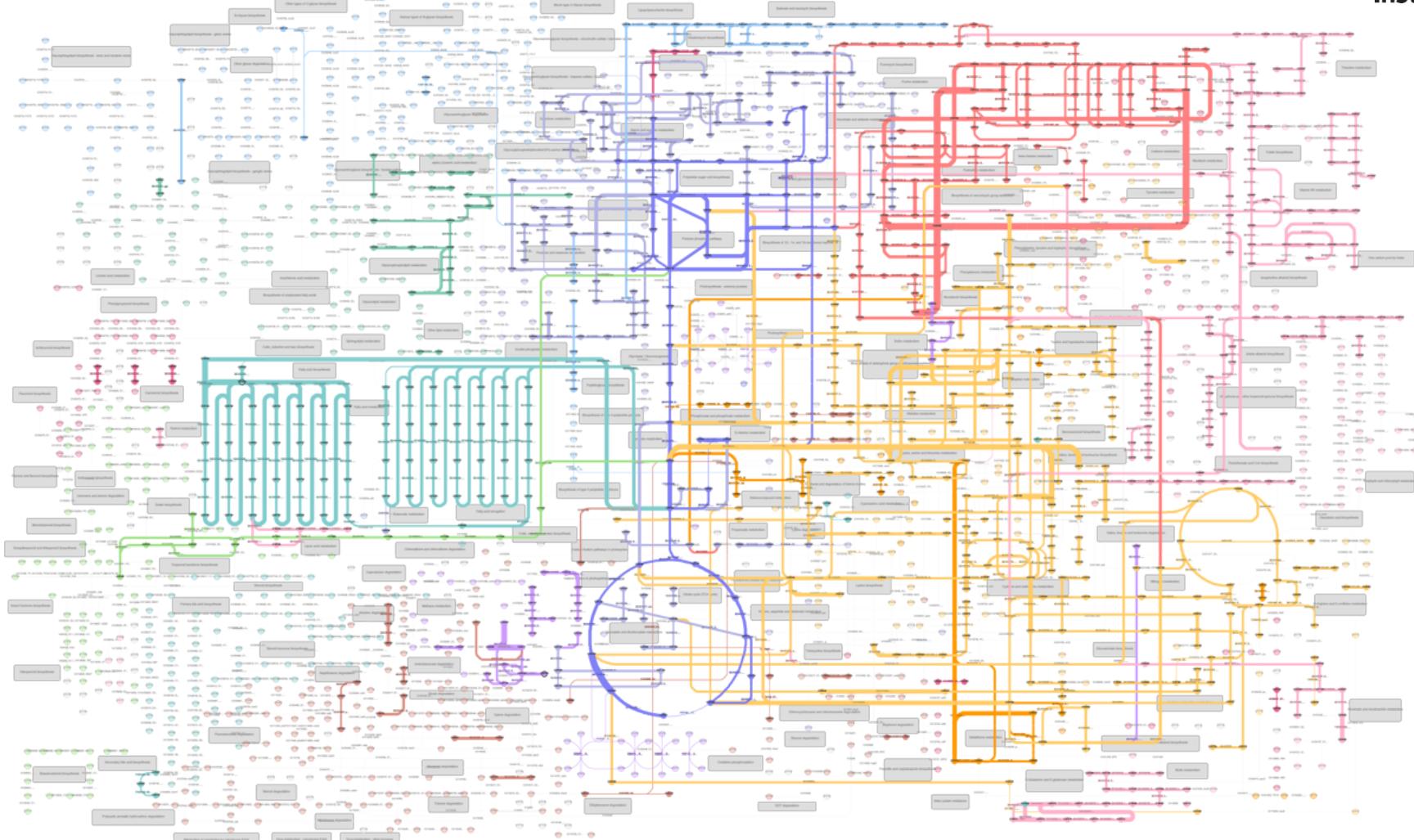
# Reactions differ by species



- Presence/absence of reactions associates with species
- Growth curves/substrate preference associates with species



# Presence/absence of KEGG reactions



<https://github.com/widdowquinn/KGML>

<http://dx.doi.org/10.6084/m9.figshare.767275>

# Flux Balance Analysis models

- FBA models represent whole-organism metabolic flux

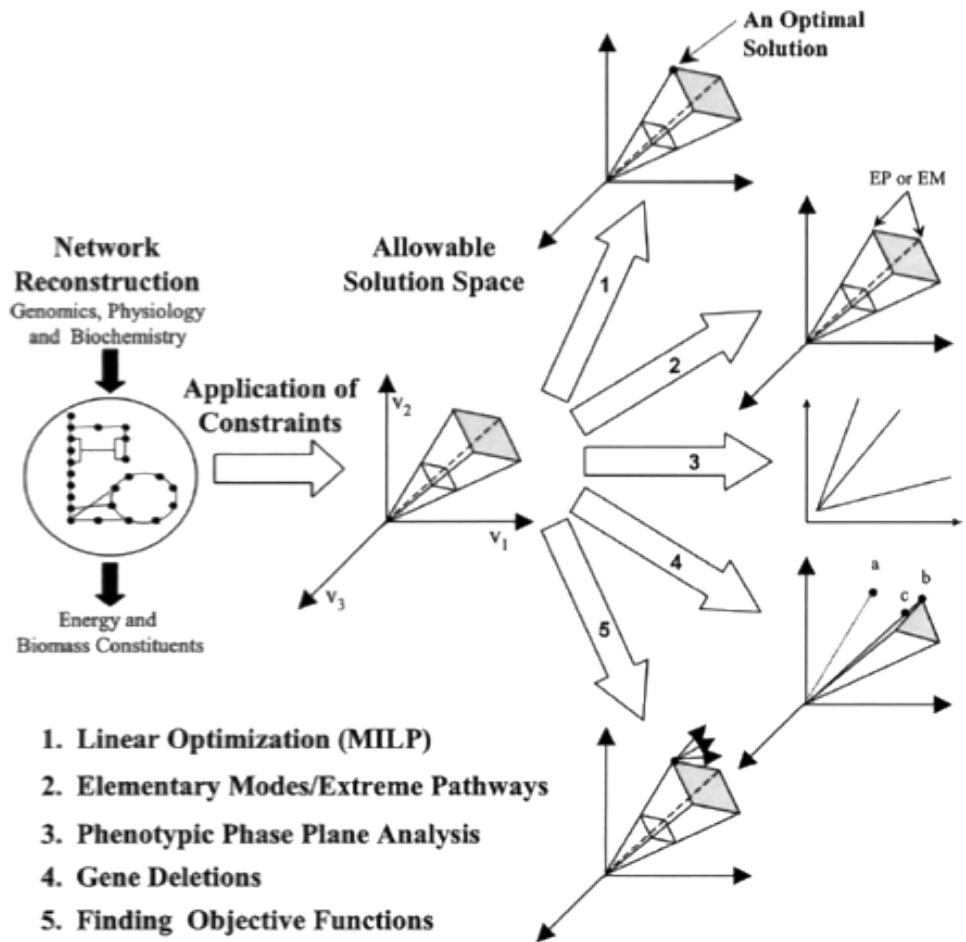
Elementary modes

Substrate usage

Predictive for knockouts

Basis for SynthBio &  
engineering

- FBA models constructed for each sequenced *Dickeya* isolate
- Transposon mutant libraries being constructed for representative isolates
- Regulatory network reconstruction in progress





# Conclusions III

- Microbial pathogen genomes are much easier to obtain, but (at the moment) typically less complete than they used to be
- **Automated annotations (and databases...) have errors**
  - be careful and critical
- **Imperfect sequences and annotations are still amazingly useful**
- Whole genome information allows for more precise taxonomic classification and diagnostics
- Genome analyses provide a springboard for understanding (and engineering) the biology of pathogens

# Acknowledgements

- **JHI *Dickeya***
  - Ian Toth
  - Sonia Humphris
  - Emma Campbell
  - Ashley Boath
  - Anne-Laure Lucquet
- **Edinburgh *Dickeya***
  - Ian Simpson
- **Oxford *Dickeya***
  - Gail Preston
  - Souvik Kusari
- **Fera *Dickeya***
  - Neil Parkinson
  - John Elphinstone
- **SASA *Dickeya***
  - Gerry Saddler
  - Vince Mulholland
- **ILVO *Dickeya***
  - Steve Baeyen
  - Johan Van Vaerenbergh
  - Martine Maes
- **Wageningen *Dickeya***
  - Jan van der Wolf
- **JHI Bioinformatics**
  - Peter Cock
  - Sue Jones
  - Linda Milne
- **JHI *E. coli***
  - Iain Milne
- **Muenster *E. coli***
  - Nicola Holden
  - Ian Toth
  - Ashleigh Holmes
- **TSL Potato**
  - Dan MacLean
  - Graham Etherington
- **Sci/Art**
  - Ian Toth
  - Elaine Shemilt
  - Genevieve Murphy
  - Danny Hill
  - Michel Perombelon
- **Aberdeen *Campylobacter***
  - Ken Forbes
  - Norval Strachan
- **JHI *Phytophthora***
  - Paul Birch
  - Steve Whisson
  - David Cooke
  - Miles Armstrong
  - Eleanor Gilroy
  - Hazel Mclellan
- **Warwick *Phytophthora***
  - Julie Squires
  - Susan Breen
- **JHI Potato**
  - Ingo Hein
  - Florian Jupe
- **TSL Potato**
  - Dan MacLean
  - Graham Etherington
- **JHI *Pectobacterium***
  - Ian Toth
  - Paul Birch
  - Kenny Bell
  - Beth Hyman
  - Maria Holeva
- **Sanger *Pectobacterium***
  - Julian Parkhill
  - Mo Sebaihia
  - Matt Holden
  - Steve Bentley
  - Nick Thomson

**AND MANY MORE...**



# How to get presentation materials

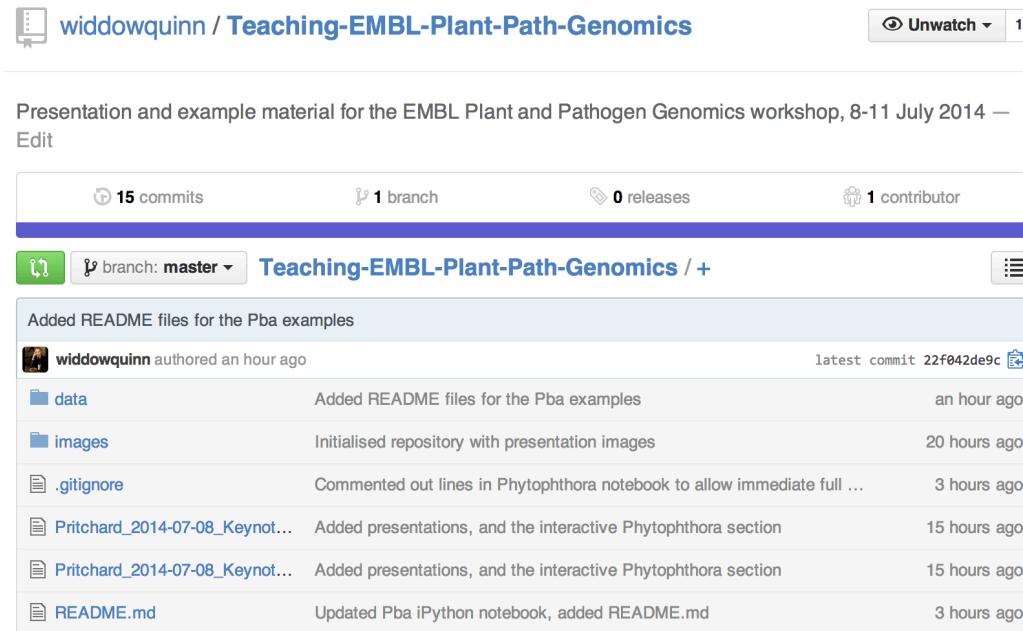
- All slides and example code/data/exercises at GitHub:

<https://github.com/widdowquinn/Teaching-EMBL-Plant-Path-Genomics>

- Use Git application or command line:

`git clone https://github.com/widdowquinn/Teaching-EMBL-Plant-Path-Genomics.git`

- Also in the course materials directory



The screenshot shows a GitHub repository page. At the top, it displays the repository name "widdowquinn / Teaching-EMBL-Plant-Path-Genomics" and a "Unwatch" button with a count of 1. Below this, a description states "Presentation and example material for the EMBL Plant and Pathogen Genomics workshop, 8-11 July 2014 — Edit". A summary bar shows "15 commits", "1 branch", "0 releases", and "1 contributor". The "branch: master" dropdown is selected. The main area shows a commit history with the following details:

File	Commit Message	Author	Date
data	Added README files for the Pba examples	widdowquinn	an hour ago
images	Initialised repository with presentation images	widdowquinn	20 hours ago
.gitignore	Commented out lines in Phytophthora notebook to allow immediate full ...	widdowquinn	3 hours ago
Pritchard_2014-07-08_Keynot...	Added presentations, and the interactive Phytophthora section	widdowquinn	15 hours ago
Pritchard_2014-07-08_Keynot...	Added presentations, and the interactive Phytophthora section	widdowquinn	15 hours ago
README.md	Updated Pba iPython notebook, added README.md	widdowquinn	3 hours ago