

Microsoft Research

# Faculty Summit



FUTURE WORLD

2011 ← → 2031

# Information Technology Applied to Bioenergy Genomics: Probabilistic Annotation using Artificial Intelligence

Advisor: Prof. Dr Ricardo Z. N. Vêncio

Danillo C. Almeida-e-Silva, grad student  
Department of Computing and Mathematics - FFCLRP  
University of Sao Paulo – Brazil

adapted from:

*Introduction: Next generation biofuels*

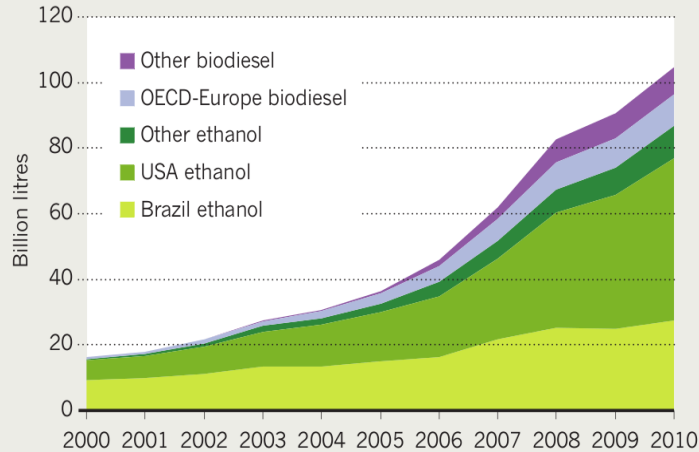
Fairley P., Nature (2011)



## THE RISE OF BIOFUELS

Biofuel production now tops 100 billion litres per year. Different fuel types vary in their costs, carbon emissions and impact on land use.

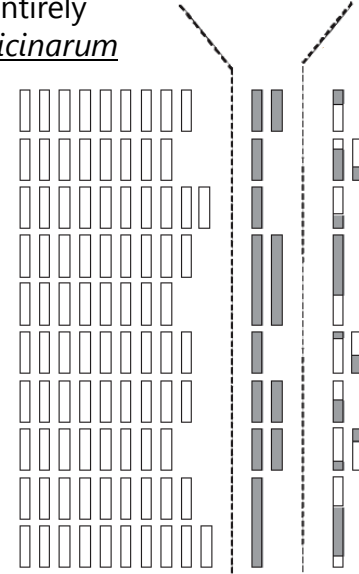
### GLOBAL BIOFUEL PRODUCTION



why?

10% to 20% of their chromosomes are inherited in their entirety from *S. spontaneum*

70% to 80% are inherited entirely from *S. officinarum*



Around 10% are the result of recombination between chromosomes from the two ancestral species.

adapted from:

*Genomics of tropical crop plants*

Moore & Ming (2008)

Function: **unknown**  
limited technological use



Function: **known**  
possible technological use

Function: **unknown**  
limited technological use



Function: **known**  
possible technological use

**Information Technology Applied to Bioenergy**  
**Genomics: Probabilistic Annotation using**  
**Artificial Intelligence**

*Goal*

Function: **unknown**  
limited technological use



Function: **known**  
possible technological use

The aim is to develop user-friendly software to **rationally guess** the biological functions of genes from sugarcane

Information Technology Applied to Bioenergy  
Genomics: Probabilistic **Annotation** using  
Artificial Intelligence

*Goal*

Function: **unknown**  
limited technological use



Function: **known**  
possible technological use

We propose to estimate the probability:  
 $p = \mathbf{P}(\text{gene X has function A} \mid \text{evidence E})$   
instead of just saying gene X is A.

Information Technology Applied to Bioenergy  
Genomics: **Probabilistic** Annotation using  
Artificial Intelligence

*Goal*

Function: **unknown**  
limited technological use



Function: **known**  
possible technological use

We propose to estimate the probability:

$$p = \mathbf{P}(\text{gene X has function A} \mid \text{evidence E})$$

instead of just saying gene X is A.

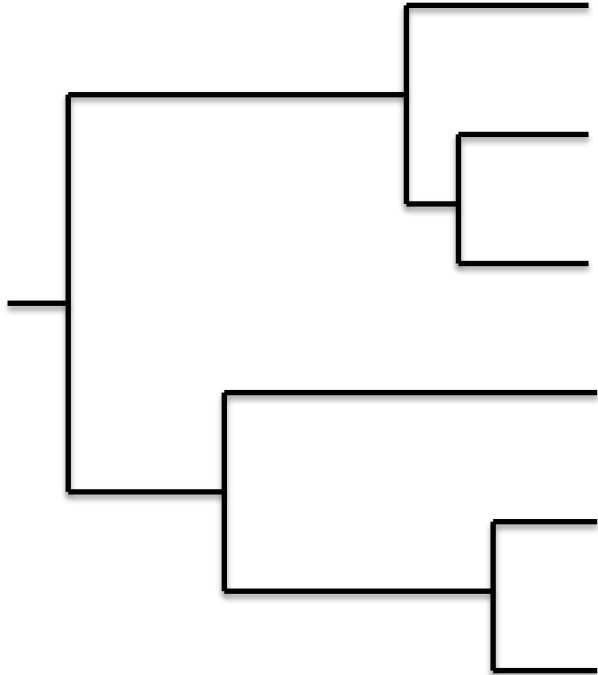
modelling quantitative and  
qualitative uncertainty

Information Technology Applied to Bioenergy  
Genomics: **Probabilistic** Annotation using  
Artificial Intelligence

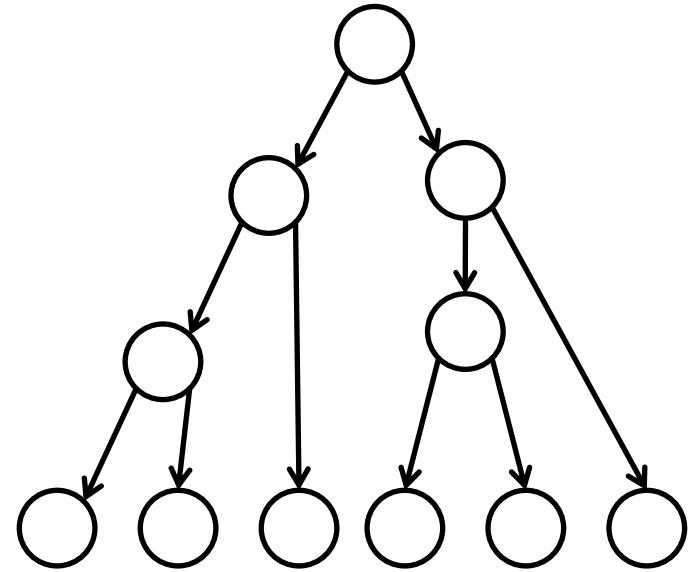
*Goal*



## Phylogenomic



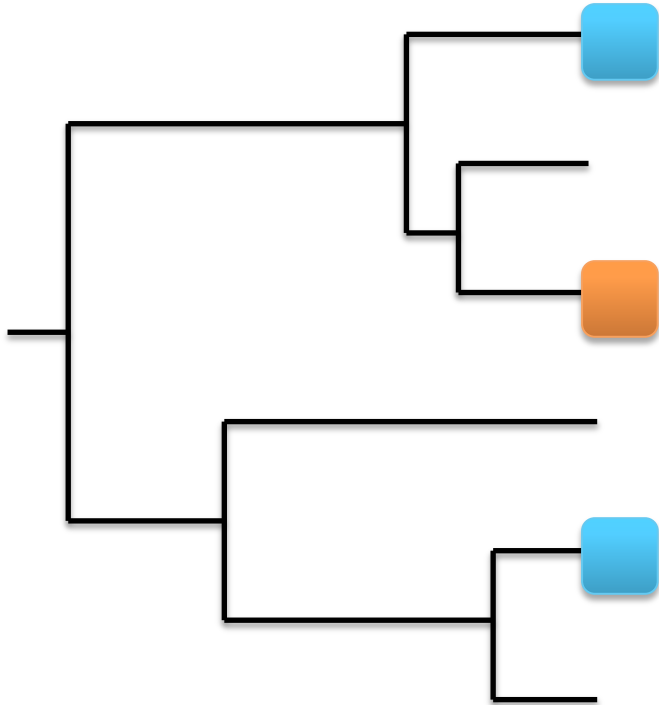
## Probabilistic



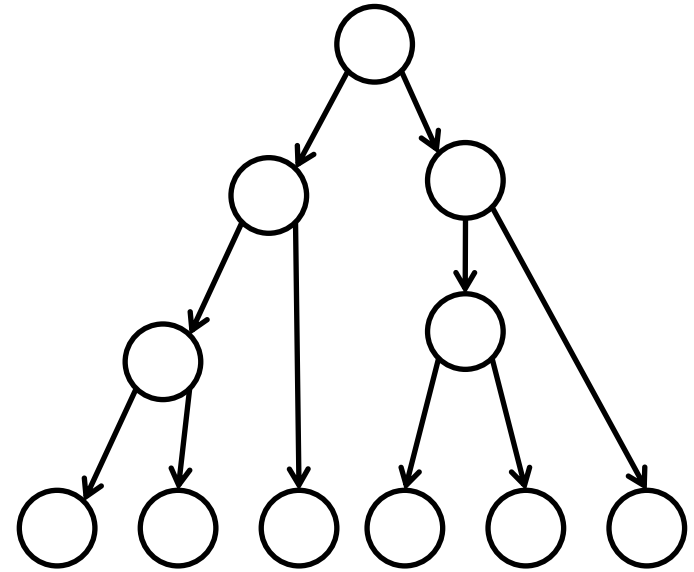
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

## Phylogenomic



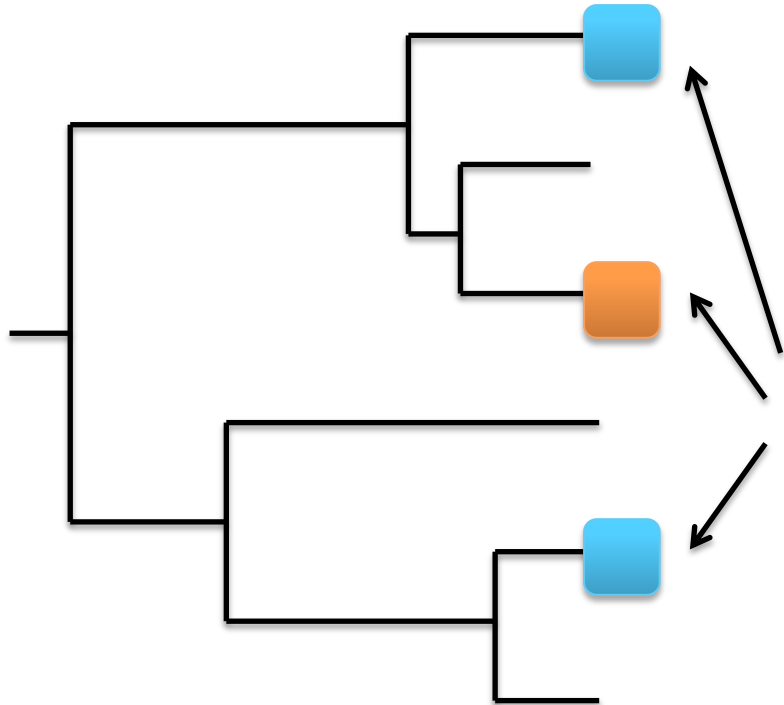
## Probabilistic



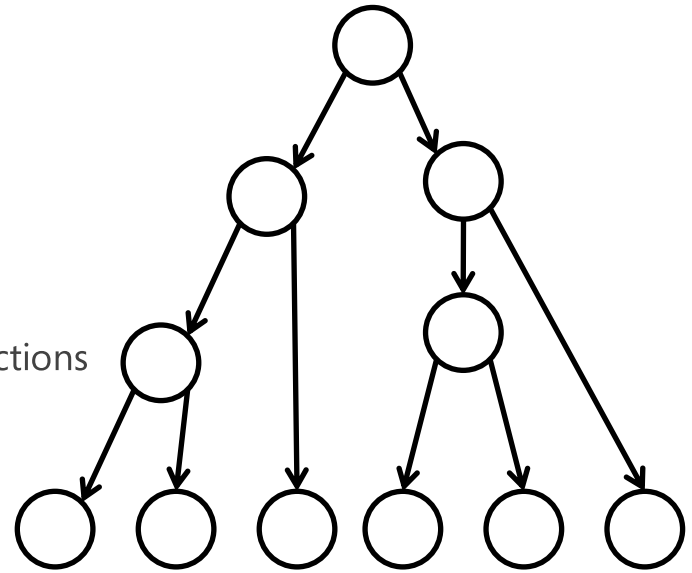
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

# Phylogenomic



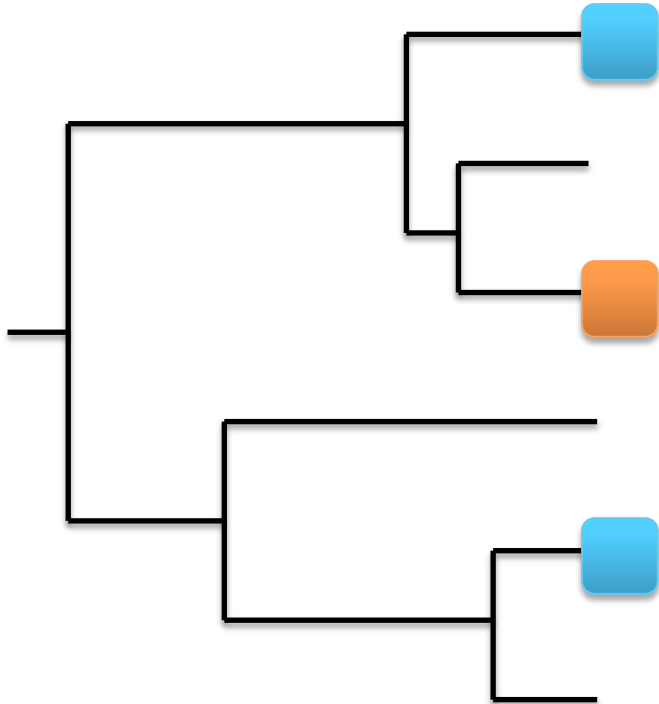
# Probabilistic



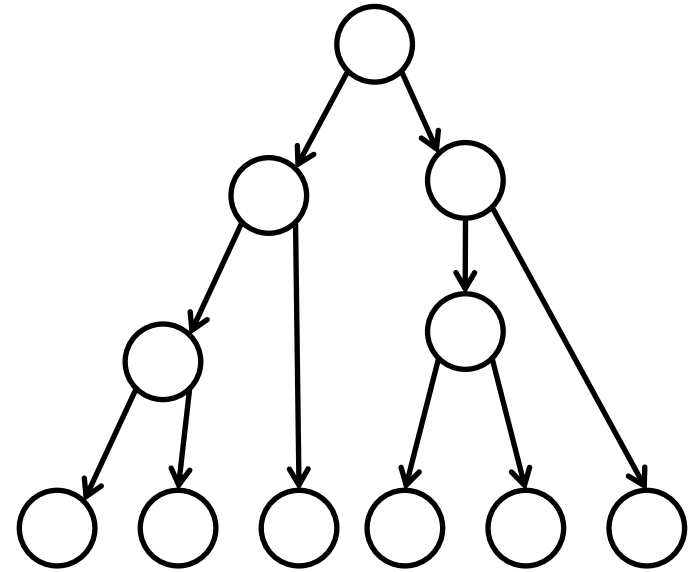
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

## Phylogenomic



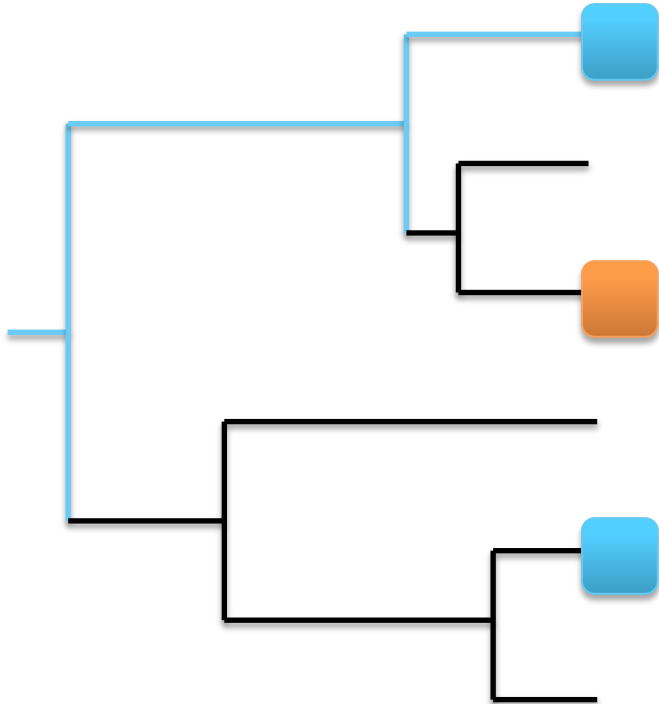
## Probabilistic



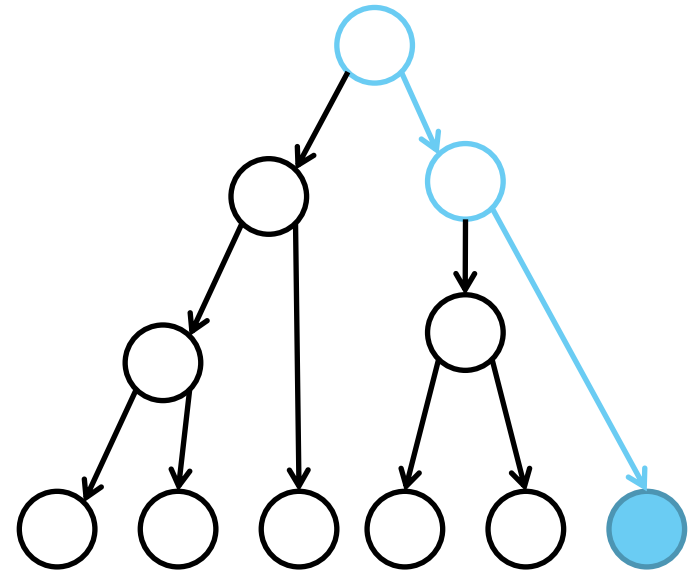
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

## Phylogenomic



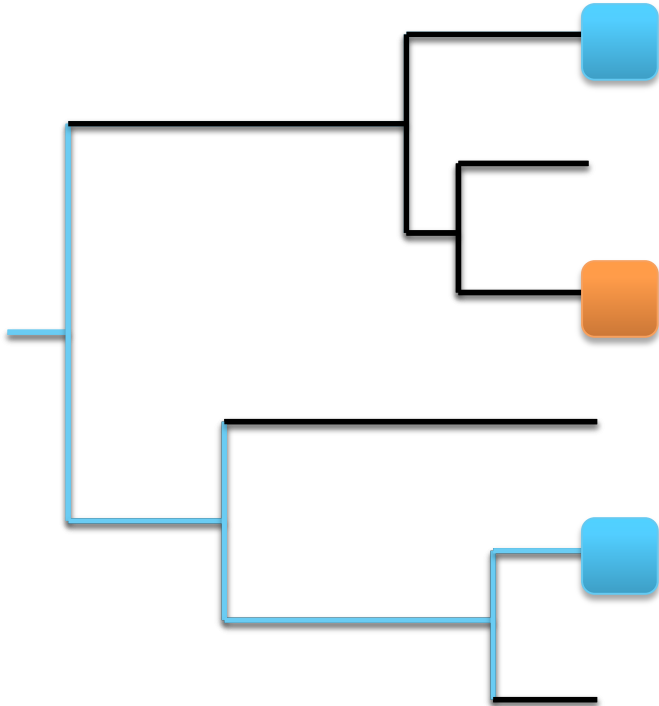
## Probabilistic



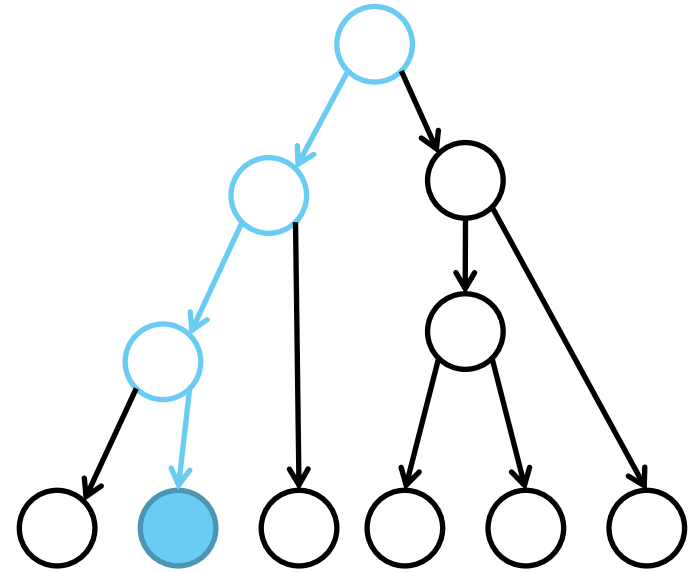
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

## Phylogenomic



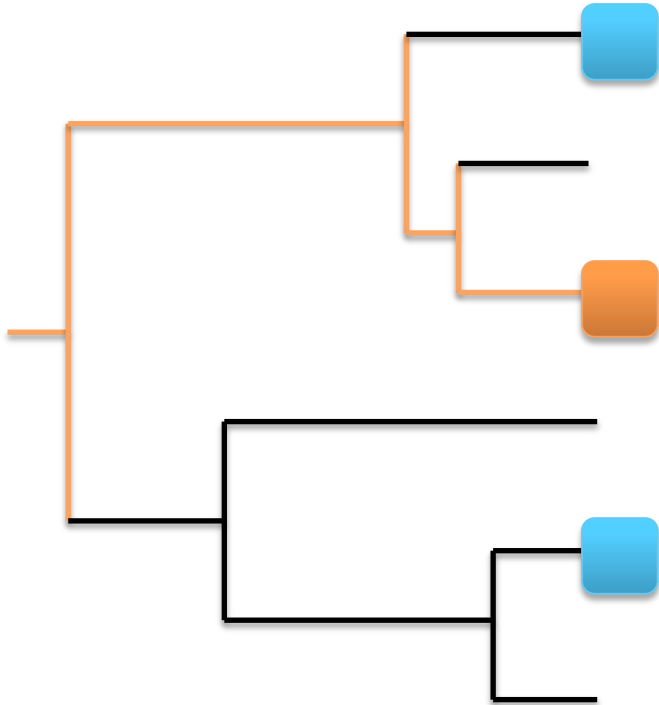
## Probabilistic



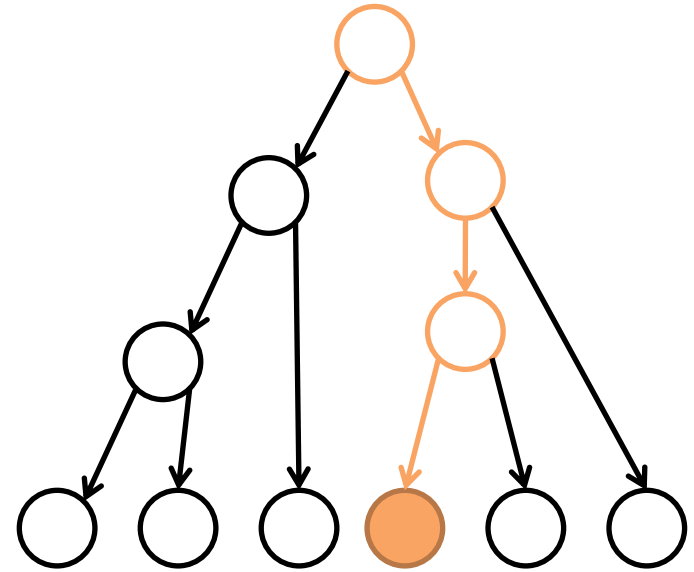
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

## Phylogenomic



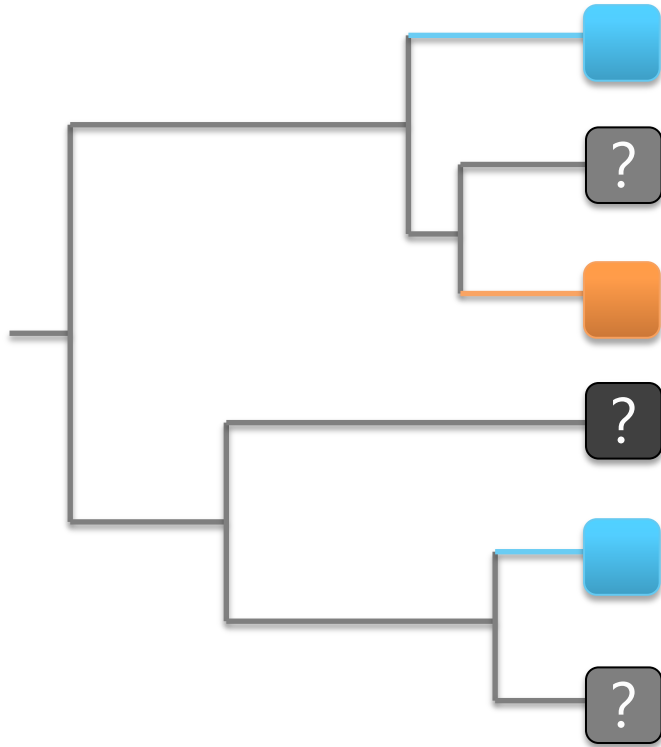
## Probabilistic



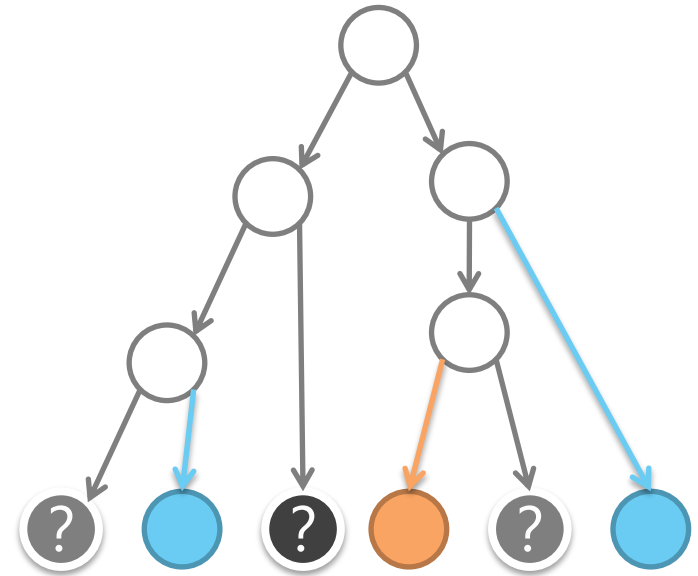
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005

## Phylogenomic



## Probabilistic

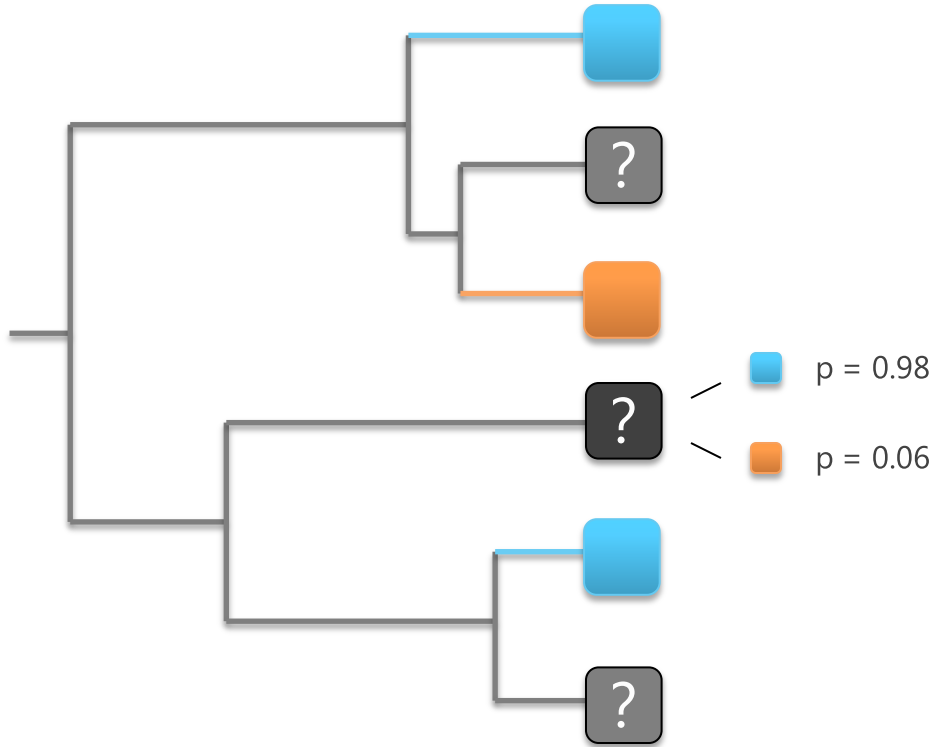


Bayesian Networks

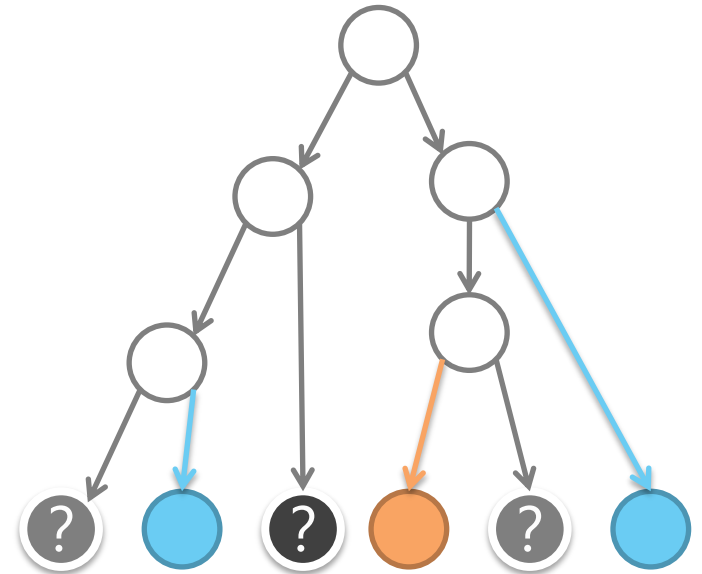
Adapted from: Engelhardt *et al.*, 2005



# Phylogenomic

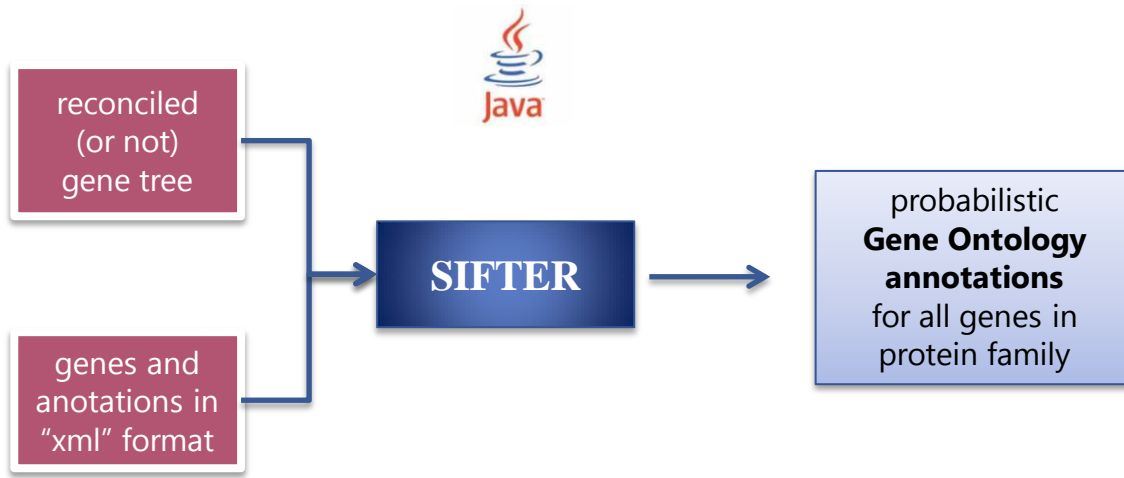


# Probabilistic



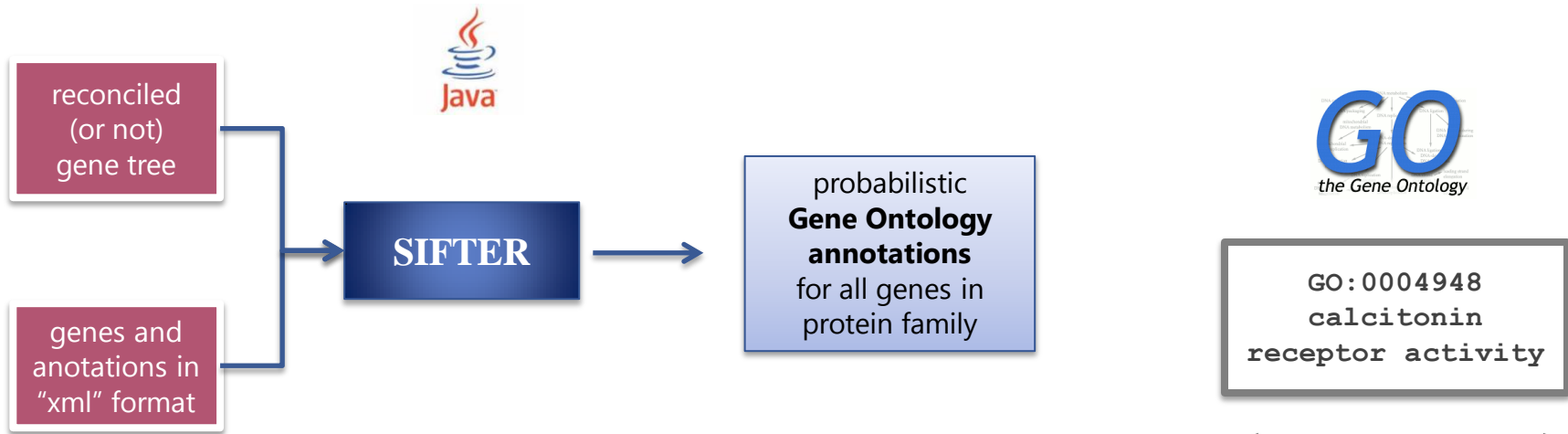
Bayesian Networks

Adapted from: Engelhardt *et al.*, 2005



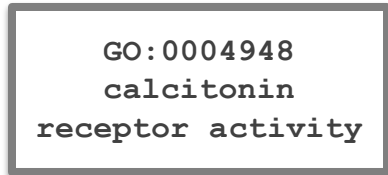
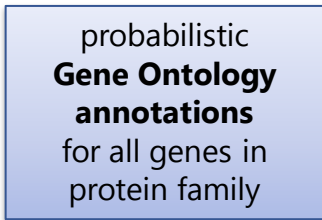
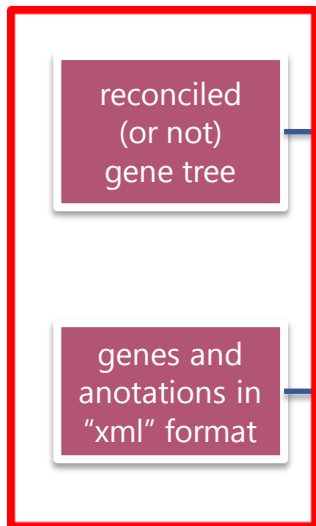
Example of annotation of genes from a given protein family:

a7ryr2_nemve	0.9960239633105679	0.07628977950793243	0.07950043523590085	<b>4948</b>
b3rjg4_triad	0.9819255688728563	0.15671589849437728	0.15943461283254143	<b>4948</b>



Example of annotation of genes from a given protein family:

a7ryr2_nemve	0.9960239633105679	0.07628977950793243	0.07950043523590085	<b>4948</b>
b3rjg4_triad	0.9819255688728563	0.15671589849437728	0.15943461283254143	<b>4948</b>



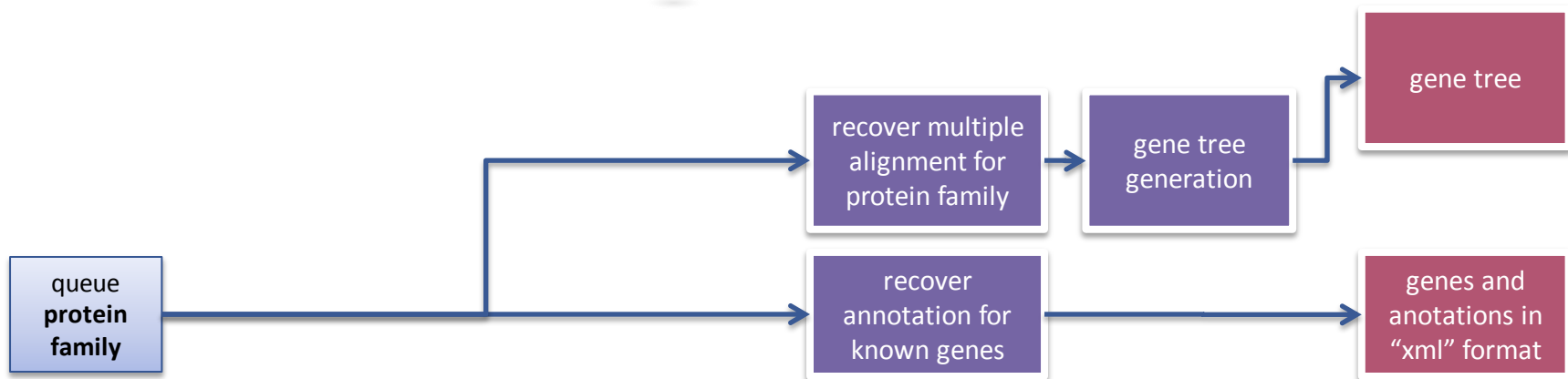
## Proof of Principle

Not suitable for actual use in real genome project

Example of annotation of genes from a given protein family:

a7ryr2_nemve	0.9960239633105679	0.07628977950793243	0.07950043523590085	<b>4948</b>
b3rjg4_triad	0.9819255688728563	0.15671589849437728	0.15943461283254143	<b>4948</b>

# State of the Art



Pfam

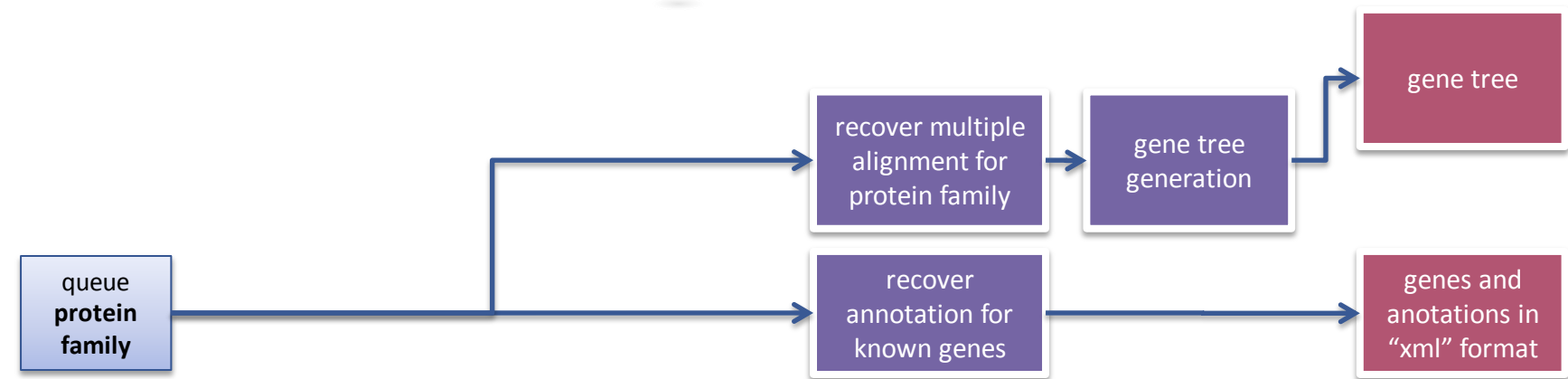
GO  
the Gene Ontology

UniProt

Fast Tree

*Starting Point for Improvement*

# State of the Art

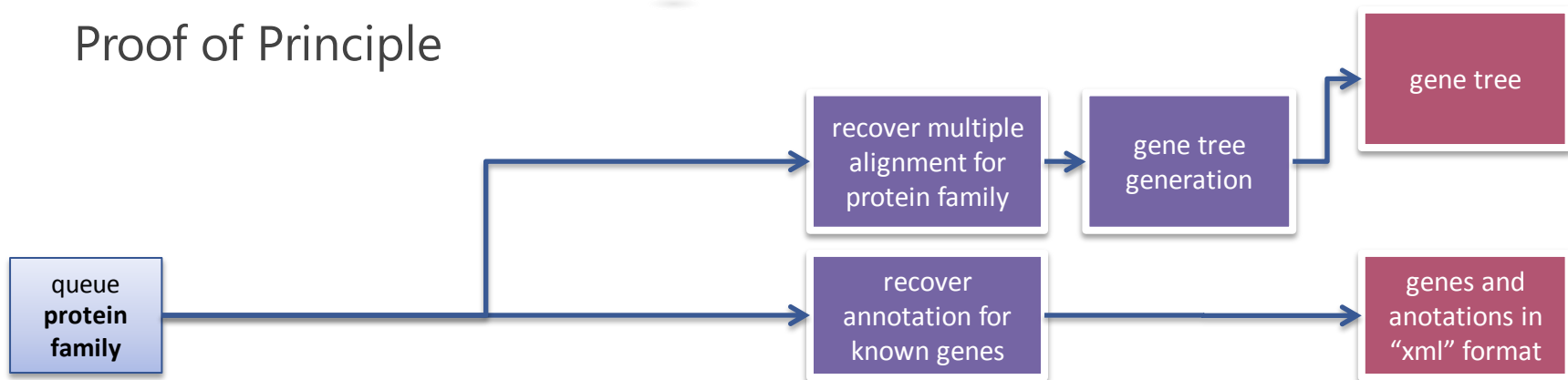


PF00258  
PF05681  
(...)



*Starting Point for Improvement*

# State of the Art Proof of Principle



Pfam

GO  
the Gene Ontology

UniProt

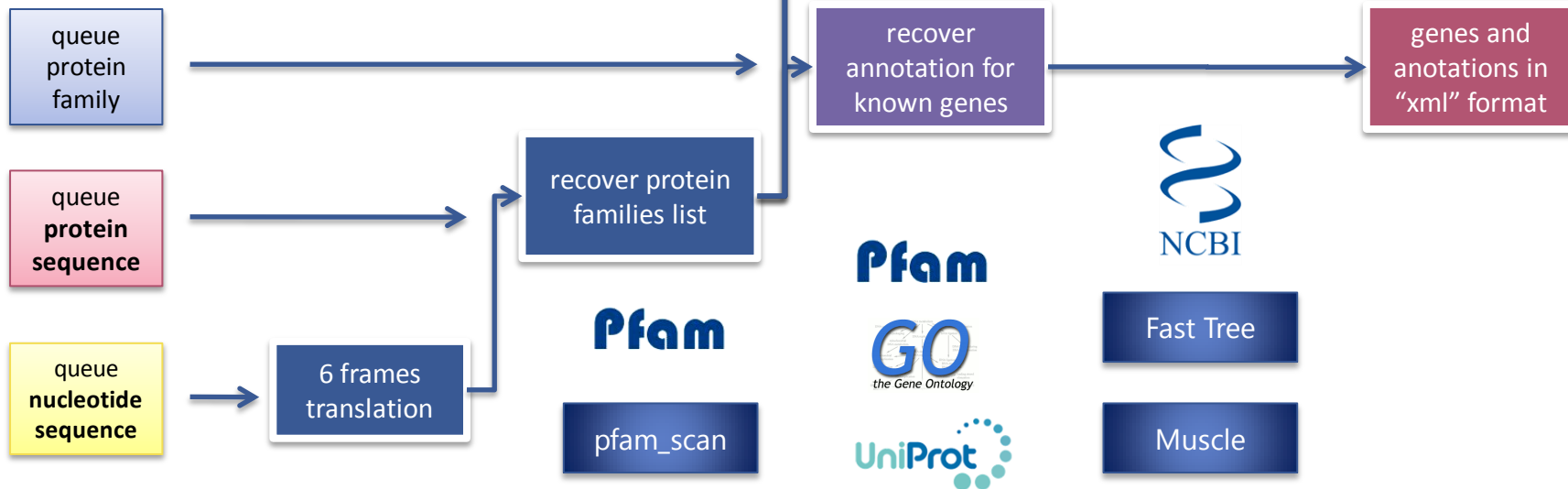
Fast Tree

*Starting Point for Improvement*



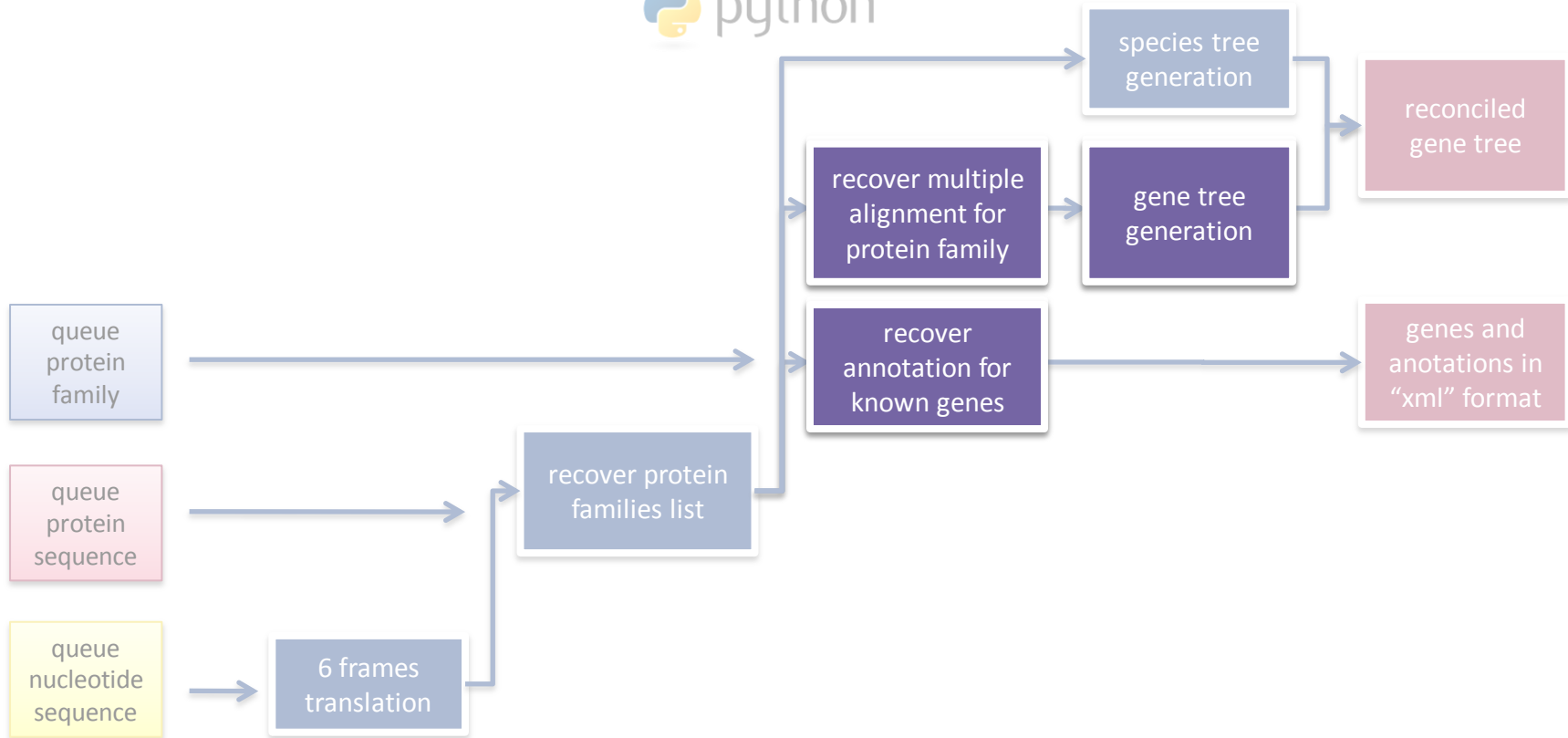
# Proof of Principle

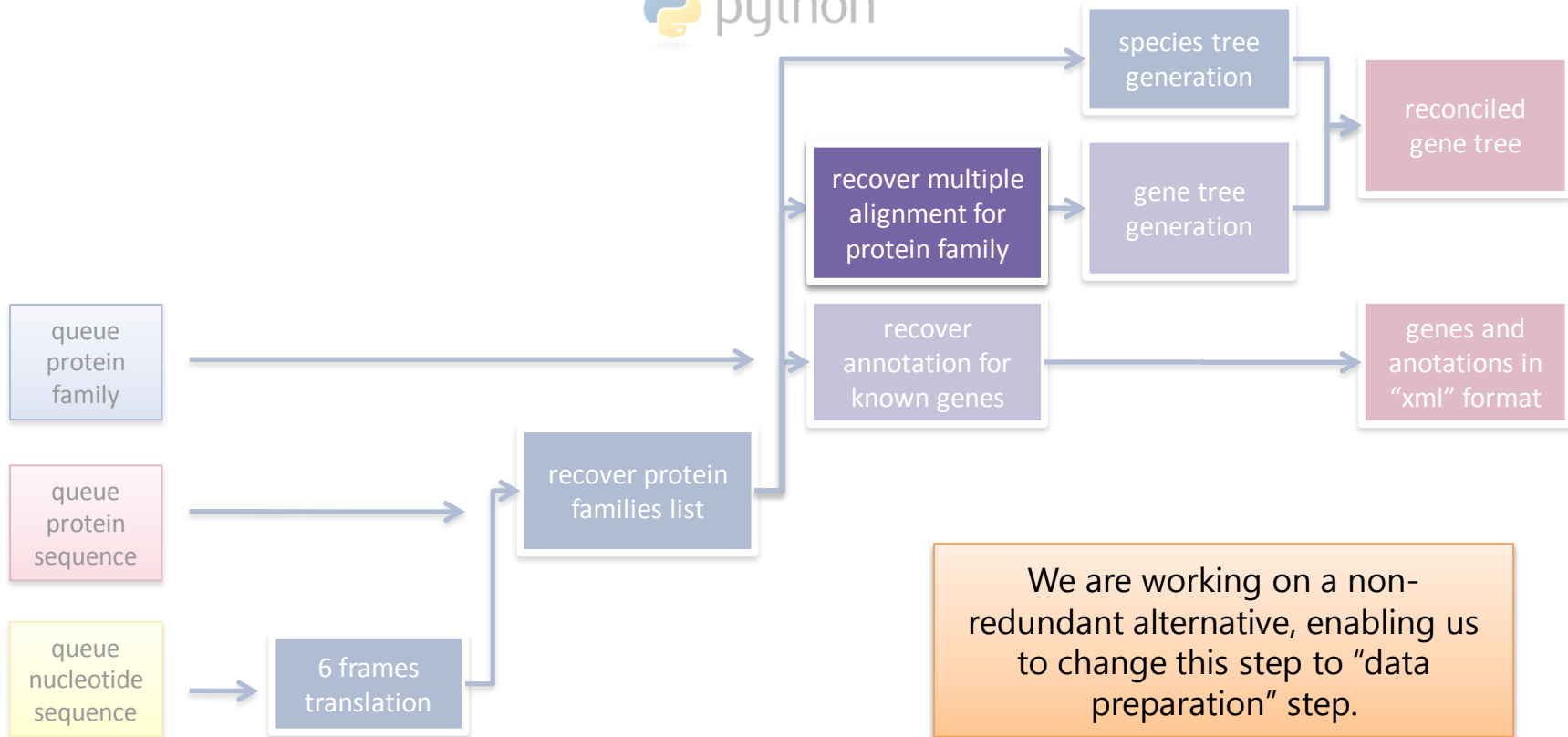
Applicable to whole genomes



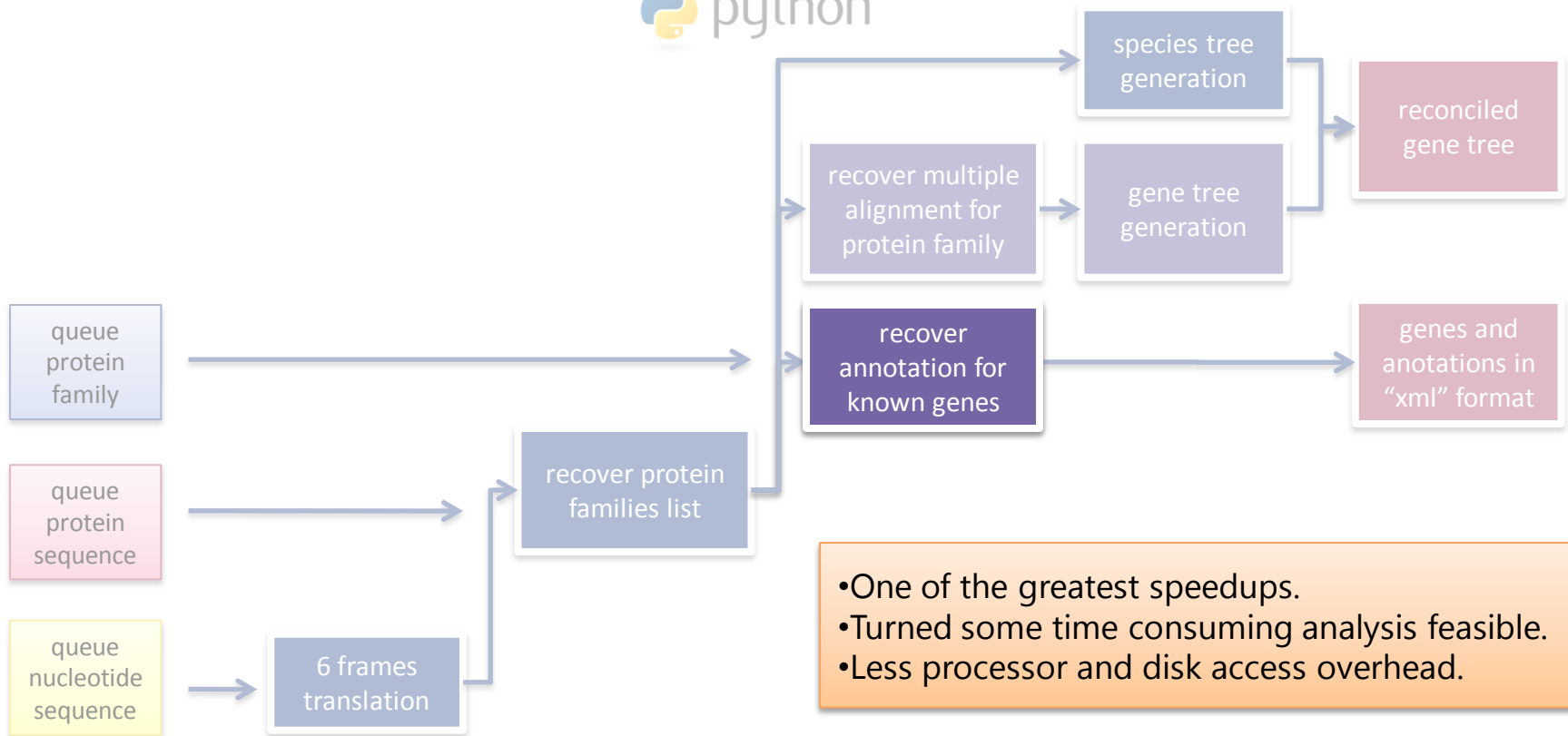
## Current Results



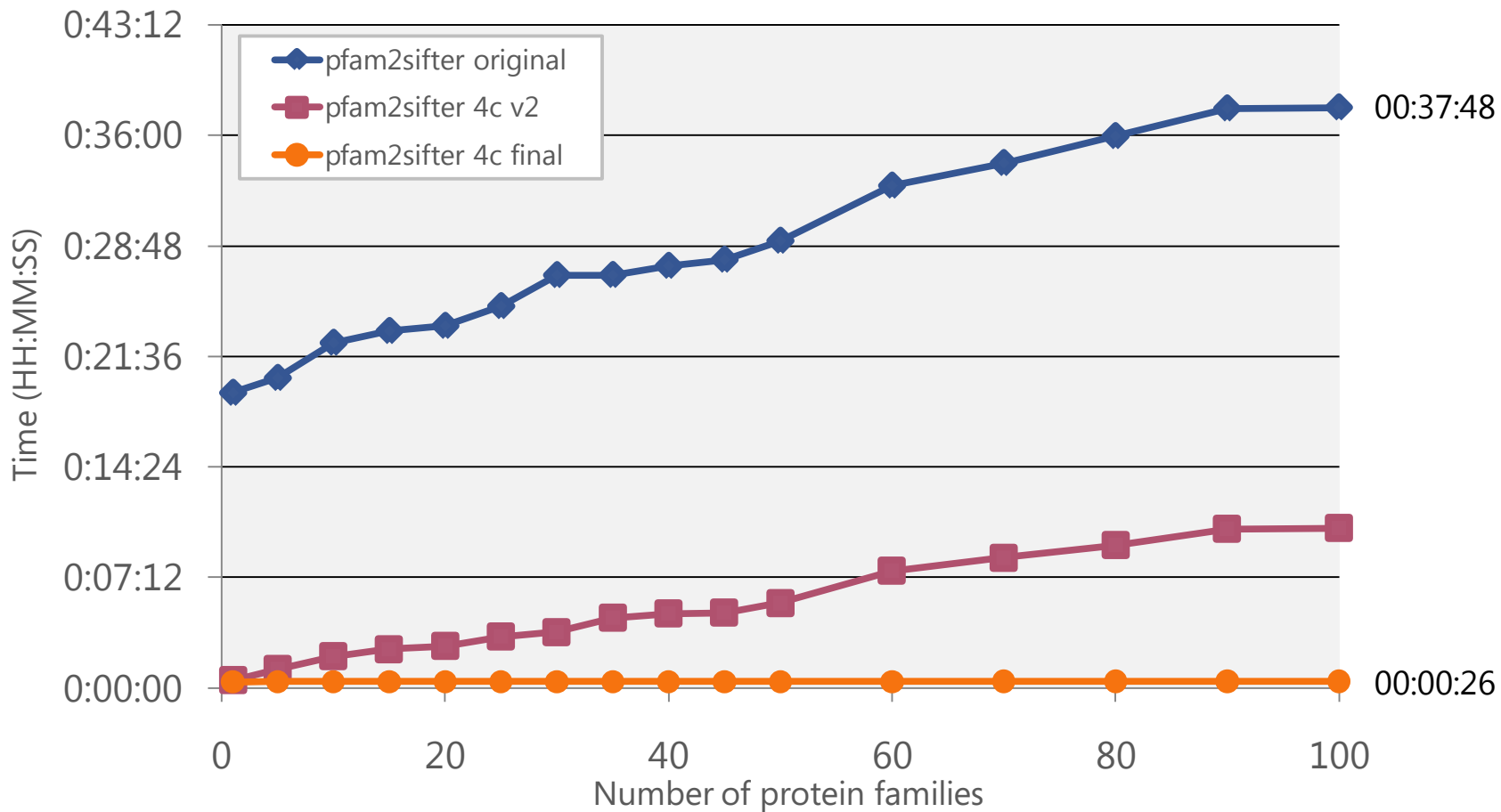


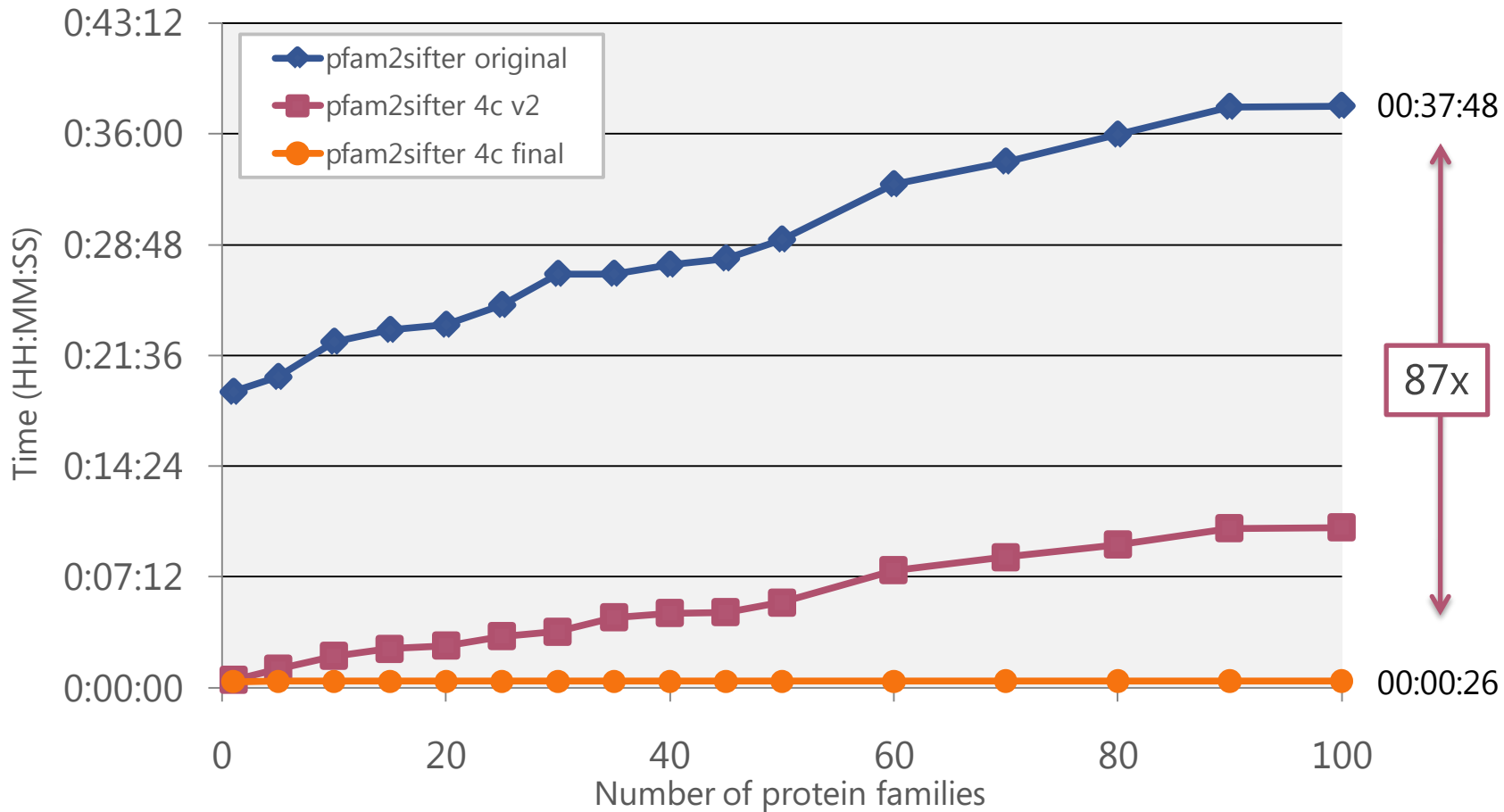


We are working on a non-redundant alternative, enabling us to change this step to "data preparation" step.

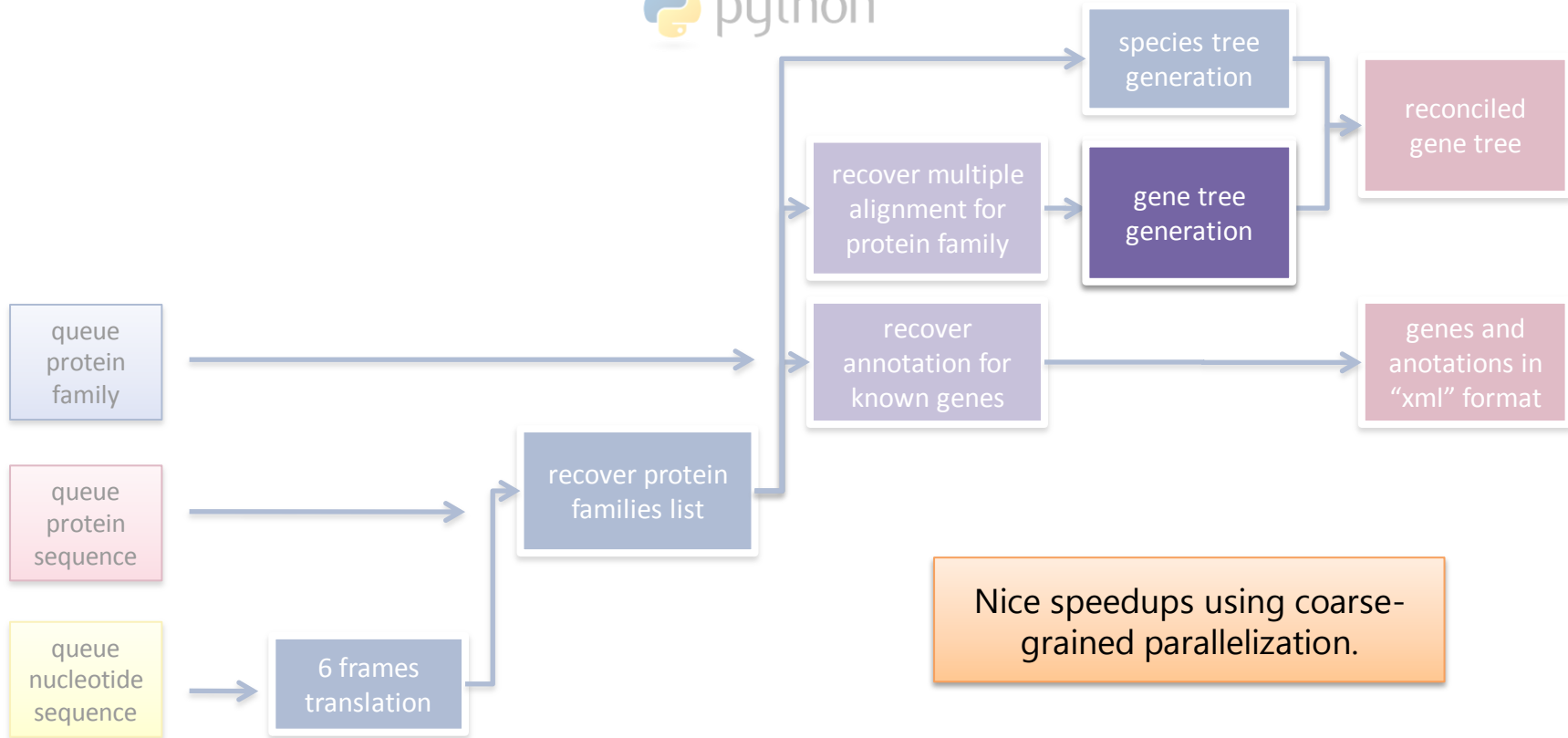


- One of the greatest speedups.
- Turned some time consuming analysis feasible.
- Less processor and disk access overhead.

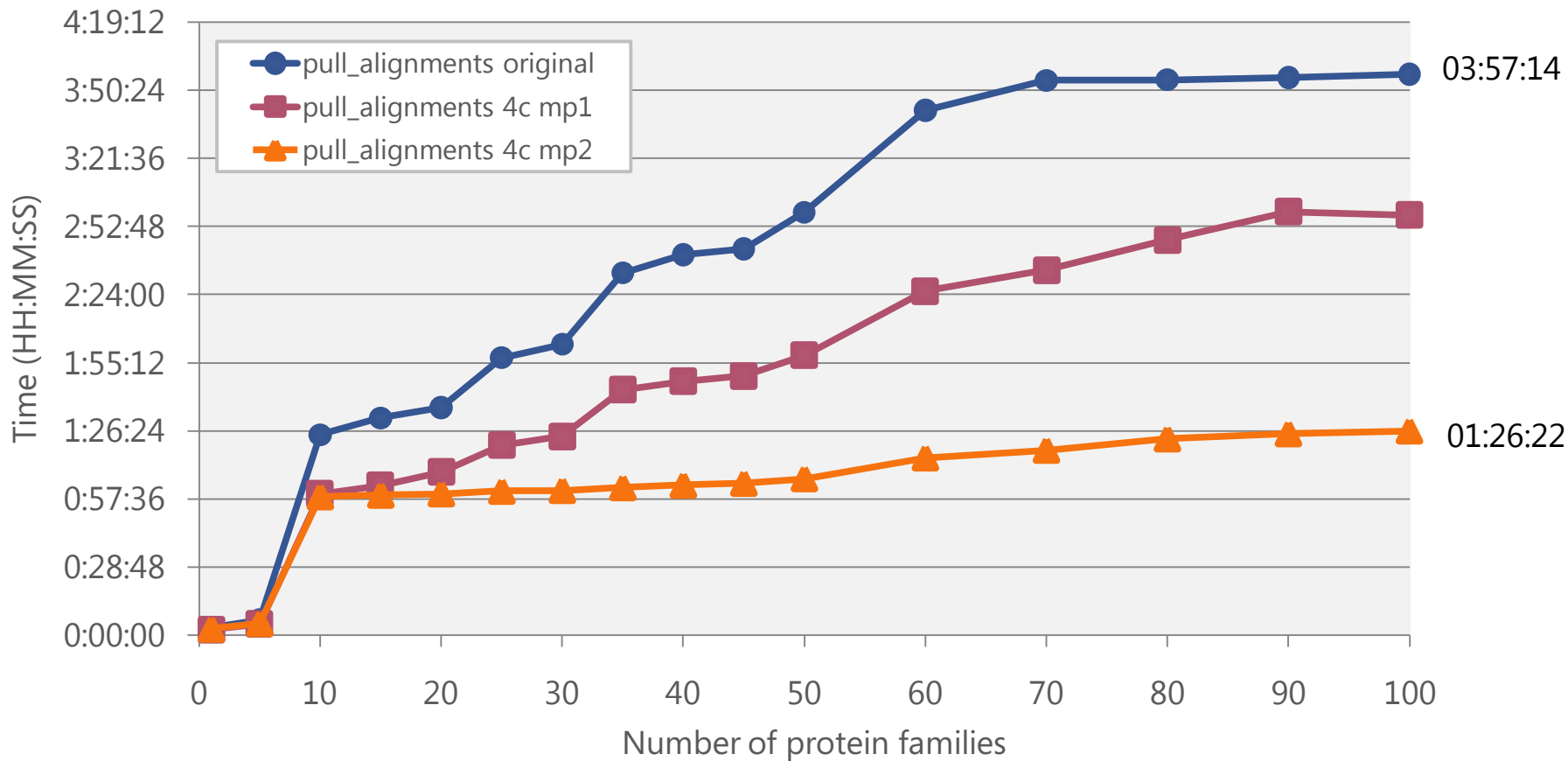


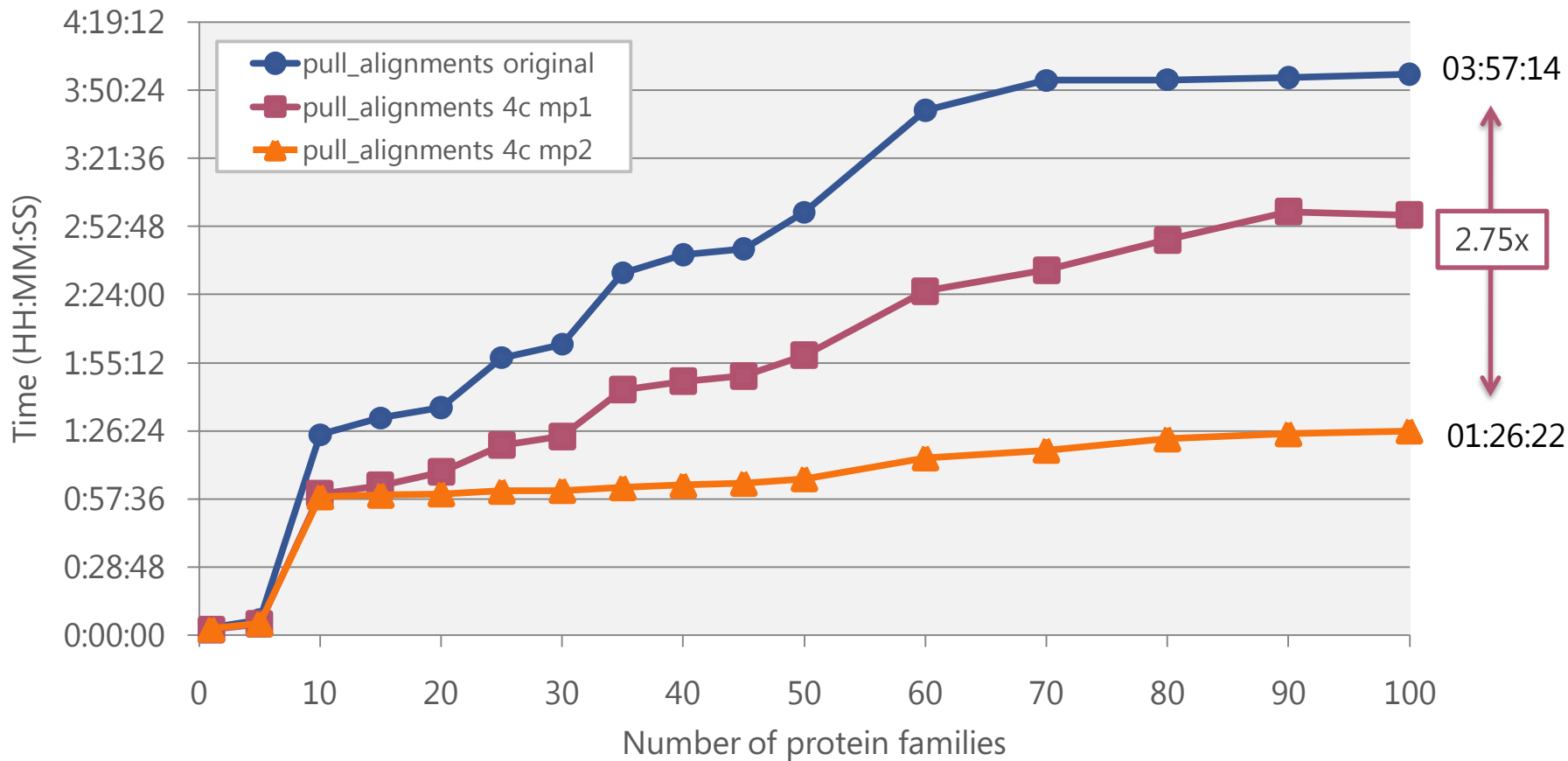


# Current Results



Nice speedups using coarse-grained parallelization.









## Microsoft Biology Foundation

[Microsoft Biology Initiative](#)

[Microsoft Biology Foundation](#)

[Microsoft Biology Tools](#)

[MBF Training](#)

[MBF Sample Applications](#)



The Microsoft Biology Foundation (MBF) is a language-neutral bioinformatics toolkit built as an extension to the Microsoft .NET Framework to help researchers and scientists work together and explore new discoveries. This open-source platform serves as a library of commonly-used bioinformatics functions. MBF facilitates collaboration and accelerates scientific research by enabling different data sets to communicate. Several universities and corporations use MBF tools, which help reduce processing time and enable scientists to focus on research.

### Microsoft Biology Foundation 2.0 Beta 1 Video

Simon Mercer, director in the Microsoft Research Connections group, provides an overview of the improvements and features of the beta release of MBF version 2.0.

Have you used this  
Research Accelerator?

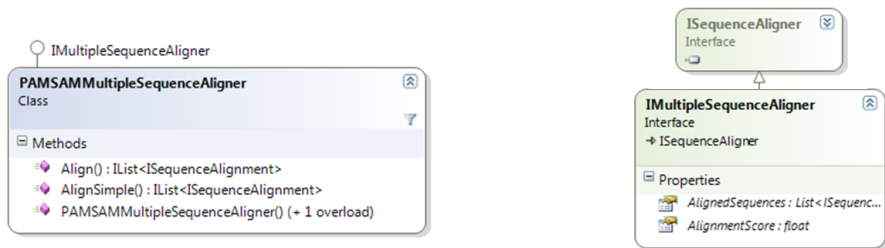
- Take a quick survey

### Downloads

- [MBF 2.0 Beta 1](#)
- [Sequence Assembler 2.0 Beta 1](#)
- [Biology Extension for Excel 2.0 Beta 1](#)
- [MBF 2.0 Beta 1 Source Code](#)
- [MBF 1.0](#)

# Multiple alignment sequences

- **PAMSAMMultipleSequenceAligner implements optimized multiple sequence alignment algorithm**
  - produces a sequence alignment using three or more sequences
  - based on MUSCLE algorithm
  - uses .NET 4.0 Parallel Extensions to take advantage of multicores
  - included in `Add-ins\Bio.PamSam.dll`
- **Algorithm defined by `IMultipleSequenceAligner` interface**
  - defines score as floating point value



Python scripting environment

adapted from:

*Microsoft Biology Foundation v2.0  
Training Material (Module 06)*

# Bioinformatics

## Sequence (re)annotation:

- pathogenic bacteria genome (*Leifsonia xyli xyli*)
- cellulosic fungi genome (*Neurospora crassa*)  
    useful “sandboxes”
- sugarcane expressed genome (SUCEST 2001)
- sugarcane genome (what is available 2012)

## Information Technology:

- automatic probabilistic functional annotation tool
- “technology transfer” from bioenergy to global health

# BioInformatics

## Sequence (re)annotation:

- pathogenic bacteria genome (*Leifsonia xyli xyli*)
- cellulosic fungi genome (*Neurospora crassa*)  
    useful “sandboxes”
- sugarcane expressed genome (SUCEST 2001)
- sugarcane genome (what is available 2012)

**SUCEST**  
The Sugar Cane EST Project

Vettore *et al.*, 2001

## Information Technology:

- automatic probabilistic functional annotation tool
- “technology transfer” from bioenergy to global health

# Bioinformatics

## Sequence (re)annotation:

- pathogenic bacteria genome (*Leifsonia xyli xyli*)
- cellulosic fungi genome (*Neurospora crassa*)  
useful “sandboxes”
- sugarcane expressed genome (SUCEST 2001)
- sugarcane genome (what is available 2012)

**SUCEST**  
The Sugar Cane EST Project

Vettore *et al.*, 2001

## Information Technology:

- automatic probabilistic functional annotation tool
- “technology transfer” from bioenergy to global health

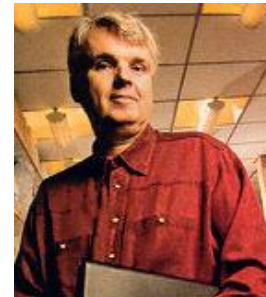
Easy to use  
Lightweight  
Broader public

*Expected Results*

# Partnership

- Microsoft Research–FAPESP Institute for IT Research

- Dr. David Heckerman
- <http://research.microsoft.com/>



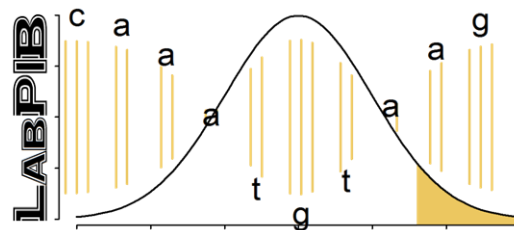
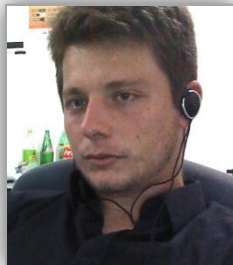
- FAPESP Bioenergy Program

- Prof. Dr. Gláucia Mendes Souza
  - Program Coordinator
  - Project: *Sugarcane Signaling And Regulatory Networks*
  - <http://bioenfapesp.org/>



# Special Thanks

- Laboratory for Biological Information Processing - **LABPIB**
  - Prof. Dr. Ricardo Z.N. Vêncio
  - Undergrad Ricardo Cacheta Waldemarin



<http://labpib.fmrp.usp.br>

# Thank you for your attention!

Danillo C. Almeida-e-Silva, grad student  
Department of Computing and Mathematics - FFCLRP  
University of Sao Paulo – Brazil  
Contact: [danillosilva@usp.br](mailto:danillosilva@usp.br)





Microsoft Research  
**FacultySummit**

© 2011 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries.

The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft cannot guarantee the accuracy of any information provided after the date of this presentation.

MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.