# Machine Learning in Bioinformatics

## Alexandros Stamatakis[1,2,3] and Franziska Reden[1]

1. Institute of Computer Science, Foundation for Research and Technology - Hellas

2. Heidelberg Institute for Theoretical Studies

3. Institute of Theoretical Informatics, Karlsruhe Institute of Technology
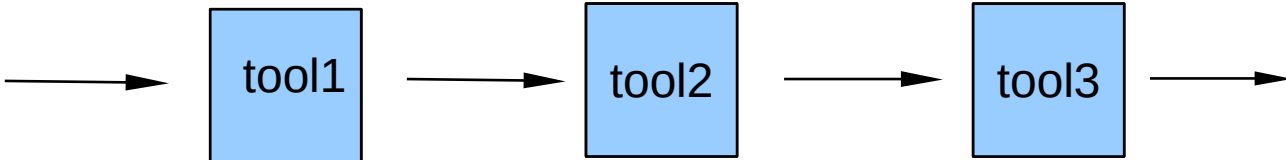
[www.biocomp.gr](www.biocomp.gr) (Crete lab)

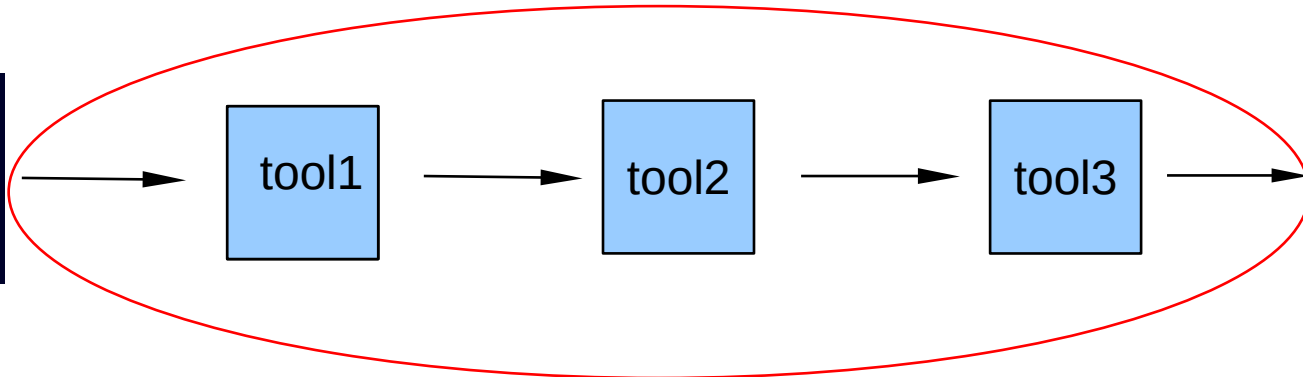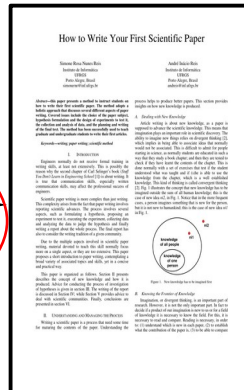[www.exelixis-lab.org](www.exelixis-lab.org) (Heidelberg lab)

# Disclaimer

- I never wanted to do machine learning

  → Somebody must keep working on algorithms, HPC, hardware architectures, `C++`

- Current generation of CS students

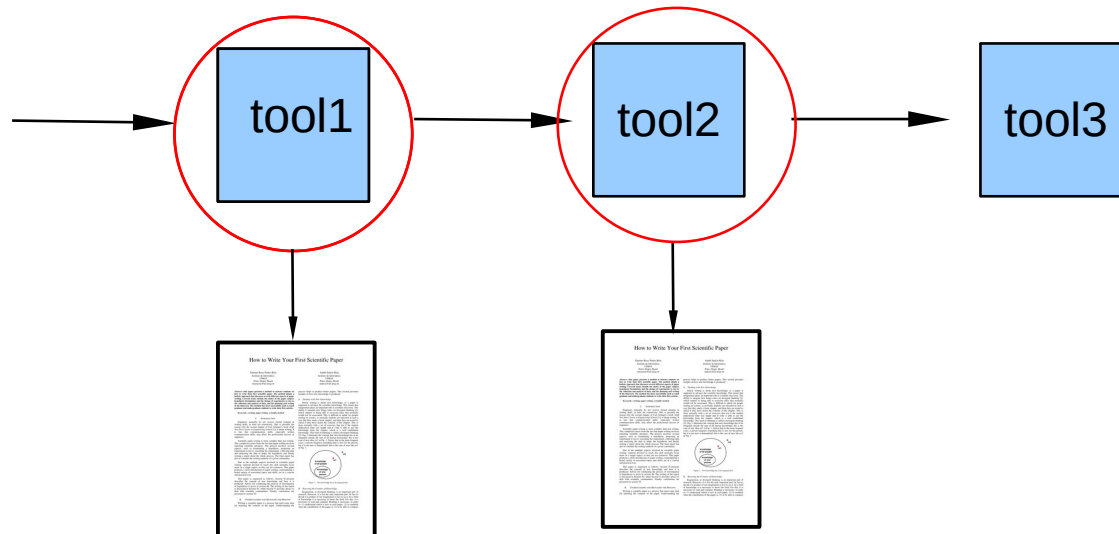  *"I want to do something with data science and/or machine learning"*

# Bioinformatics

# Bioinformatics



**Data-centric:** pipeline building



**Method-centric:** tool building

# Example: Tree Inference Pipeline

```
Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT
```

→

**MSA Program**

→

```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

→

**Tree inference program**

↓

Taxon 1    Taxon 3

Taxon 2    Taxon 4

# Tree Inference Pipeline

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

→

**MSA Program**

→

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

→

**Tree inference program**

**Multiple Sequence Alignment:**
Mostly *ad hoc* methods →
no widely used uncertainty
quantification approach, but
Muscle 5 tool → ensembles

Taxon 1

Taxon 3

Taxon 2

Taxon 4

# Tree Inference Pipeline

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

→

**MSA Program**

→

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

→

**Tree inference program**

**Phylogenetic Inference:**
A long history of explicit uncertainty models
Bootstrap Methods for Maximum Likelihood
Posterior Probabilities for Bayesian Inference using MCMC

Taxon 1    Taxon 3

Taxon 2    Taxon 4

# Naively Propagating Uncertainty

```
Unaligned FASTA
LOP001-11    AACTTTATATTTTATTTTTGGAATTTGAGCAGGAATAGTAGGAACCTCTT
LOP002-11    ATATTTTATTTTTGGAATTTGAGCTGGATTAATTGGAACTTCATT
LOP003-11    AACTCTATATTTTATTTTTGGAATTTGAGCAGGATTACTAGGAACT
LOP004-11    TCTATATTTTATTTTTGGAATTTGAGCAGGTTTAGTTGGAACTTCATT
LOP005-11_F  AATGATGTTCCTAACATACCTGCTCAAATACCAAAAATAAAATATAAAGT
```
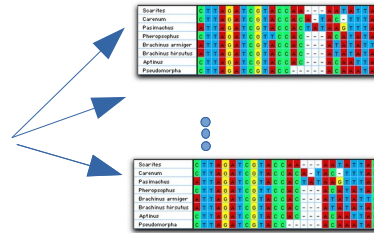
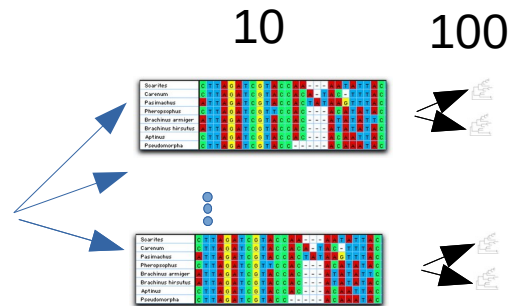# Naively Propagating Uncertainty



10

**Unaligned FASTA**

| | |
|---|---|
| LOP001-11 | AACTTTATATTTTATTTTTGGAATTTGAGCAGGAATAGTAGGAACCTCTT |
| LOP002-11 | ATATTTTATTTTTGGAATTTGAGCTGGATTAATTGGAACTTCATT |
| LOP003-11 | AACTCTATATTTTATTTTTGGAATTTGAGCAGGATTACTAGGAACT |
| LOP004-11 | TCTATATTTTATTTTTGGAATTTGAGCAGGTTTAGTTGGAACTTCATT |
| LOP005-11_F | AATGATGTTCCTAACATACCTGCTCAAATACCAAAAATAAAATATAAAGT |

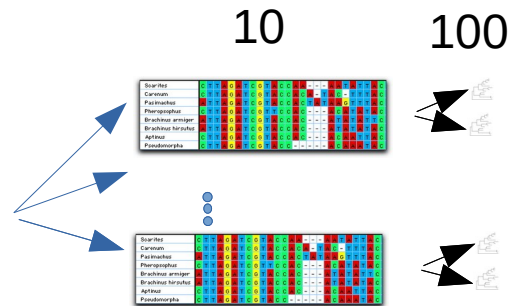# Naively Propagating Uncertainty



10    100

**Unaligned FASTA**
LOP001-11    AACTTTATATTTTATTTTTGGAATTTGAGCAGGAATAGTAGGAACCTCTT
LOP002-11    ATATTTTATTTTTGGAATTTGAGCTGGATTAATTGGAACTTCATT
LOP003-11    AACTCTATATTTTATTTTTGGAATTTGAGCAGGATTACTAGGAACT
LOP004-11    TCTATATTTTATTTTTGGAATTTGAGCAGGTTTAGTTGGAACTTCATT
LOP005-11_F    AATGATGTTCCTAACATACCTGCTCAAATACCAAAAATAAAATATAAAGT
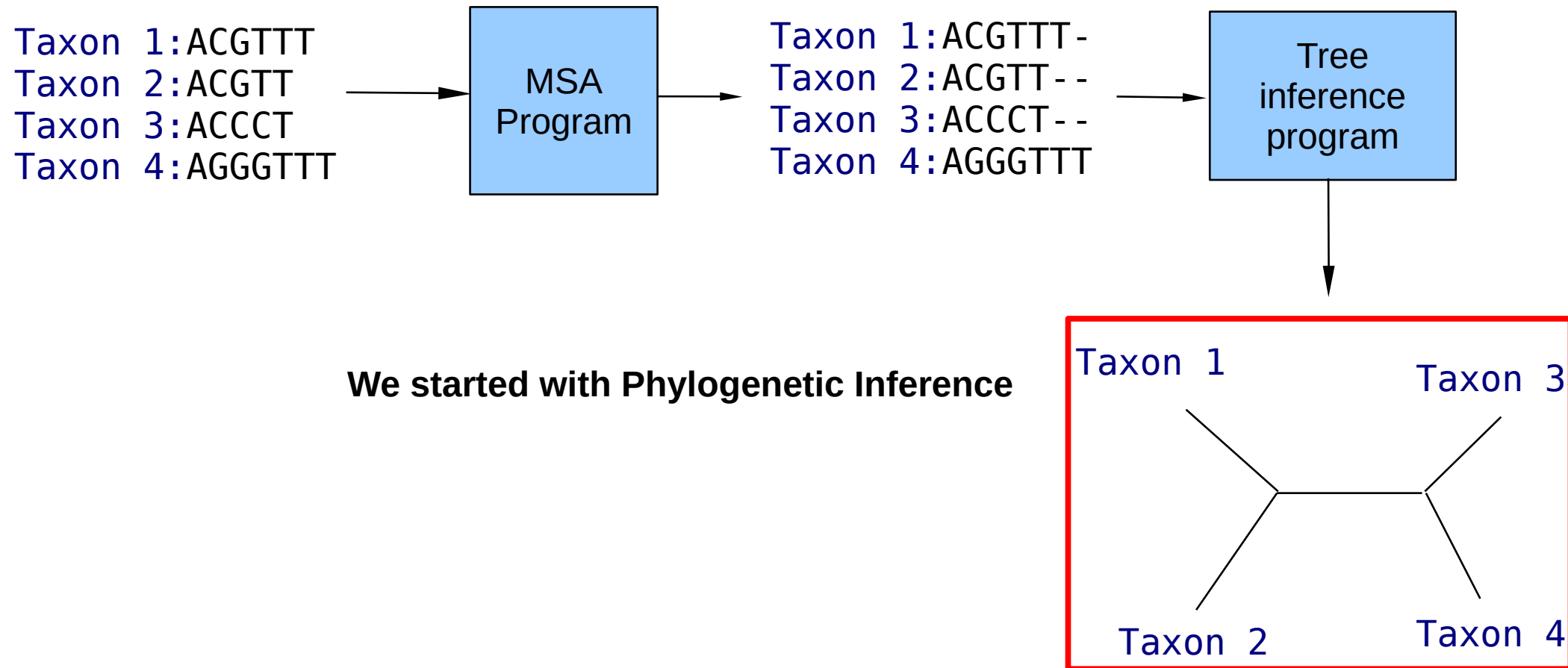
# Naively Propagating Uncertainty

# Key Idea

- Given the **input** data
  - e.g., unaligned or aligned sequences
- Predict the variance of the **output**
  - → decide if we need to propagate uncertainty
  - → decide how to run the analysis
    - generate a single result or ensemble of results?

# Tree Inference Pipeline

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

→

MSA Program

→

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

→

Tree inference program

↓

**We started with Phylogenetic Inference**

Taxon 1

Taxon 3

Taxon 2

Taxon 4

# From Easy to Hopeless—Predicting the Difficulty of Phylogenetic Analyses 🔓

Julia Haag ✉, Dimitri Höhler, Ben Bettisworth, Alexandros Stamatakis

# SARS-CoV-2 data

# SARS-CoV-2 data

```
The predicted difficulty for MSA examples/covid.fasta is: 0.84.

FEATURES:

num_taxa: 4869

num_sites: 28361

[ ... ]

num_sites/num_taxa: 5.82

[ ... ]

avg_rfdist_parsimony: 0.79

proportion_unique_topos_parsimony: 1.0

Feature computation runtime:     1830.182 seconds

[ ... ]
```

JOURNAL ARTICLE

**Phylogenetic Analysis of SARS-CoV-2 Data Is Difficult**

Benoit Morel, Pierre Barbera, Lucas Czech, Ben Bettisworth, Lukas Hübner, Sarah Lutteropp, Dora Serdari, Evangelia-Georgia Kostaki, Ioannis Mamais, Alexey M Kozlov ... Show more
   Author Notes

*Molecular Biology and Evolution*, Volume 38, Issue 5, May 2021, Pages 1777–1791, https://doi.org/10.1093/molbev/msaa314
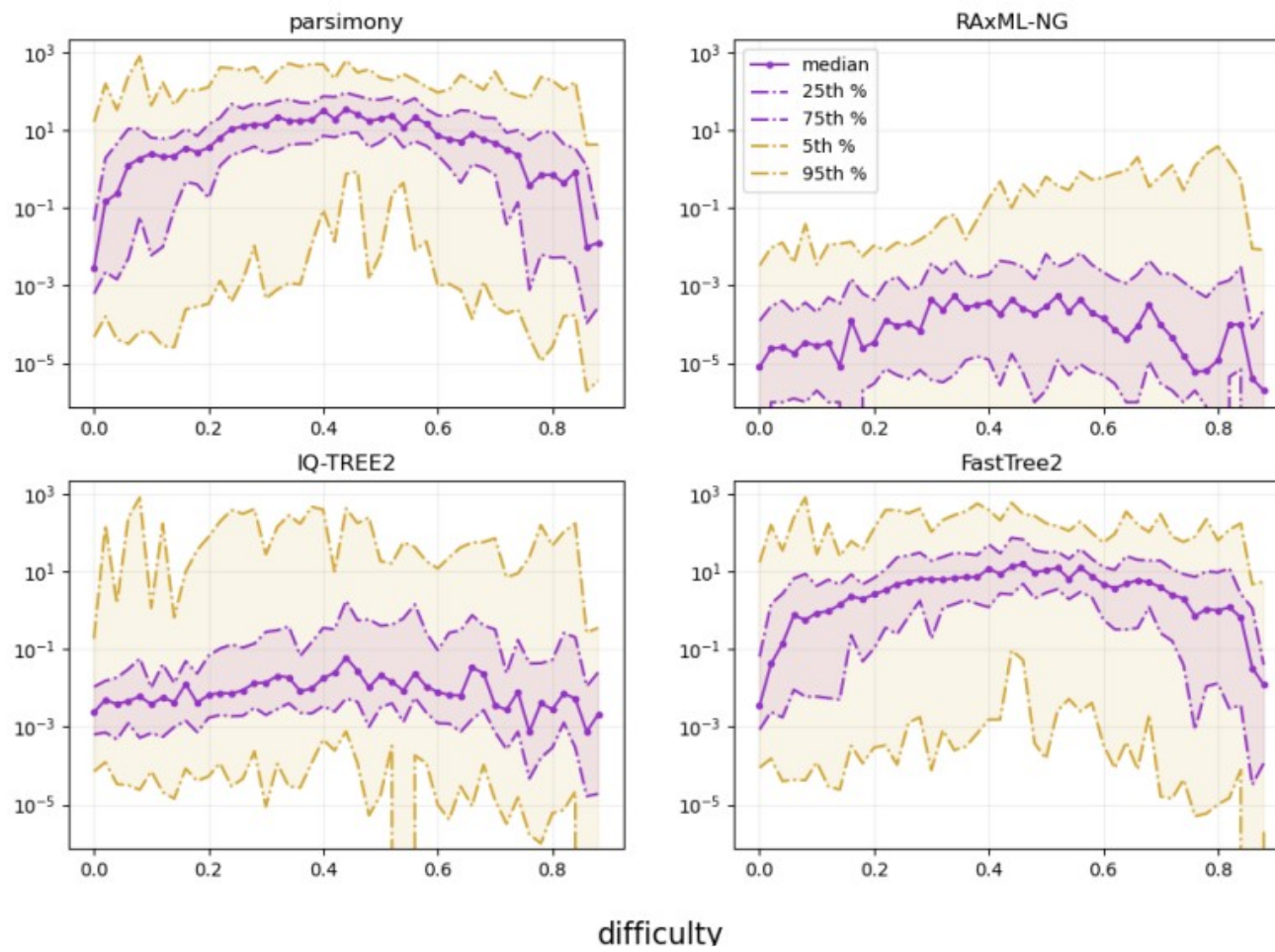**Published:** 15 December 2020

# Use case 1:
## Tool/Model performance as a Function of Difficulty

# Use Case 2:
## Adaptive Search Algorithms

Tune heuristic search parameters as a function of difficulty

→ equally accurate results

→ much faster than difficulty-agnostic algorithm

JOURNAL ARTICLE

**Adaptive RAxML-NG: Accelerating Phylogenetic Inference under Maximum Likelihood using Dataset Difficulty** 🔓

Anastasis Togkousidis ✉, Oleksiy M Kozlov, Julia Haag, Dimitri Höhler, Alexandros Stamatakis    Author Notes

*Molecular Biology and Evolution*, Volume 40, Issue 10, October 2023, msad227, https://doi.org/10.1093/molbev/msad227

**Published:** 06 October 2023    **Article history** ▾

# Use Case 3:
## Biased Experimental Setup



Accuracy with data from **our** paper

# But ...



Accuracy from paper by others – using different data

# Skewed Difficulty Distribution



21

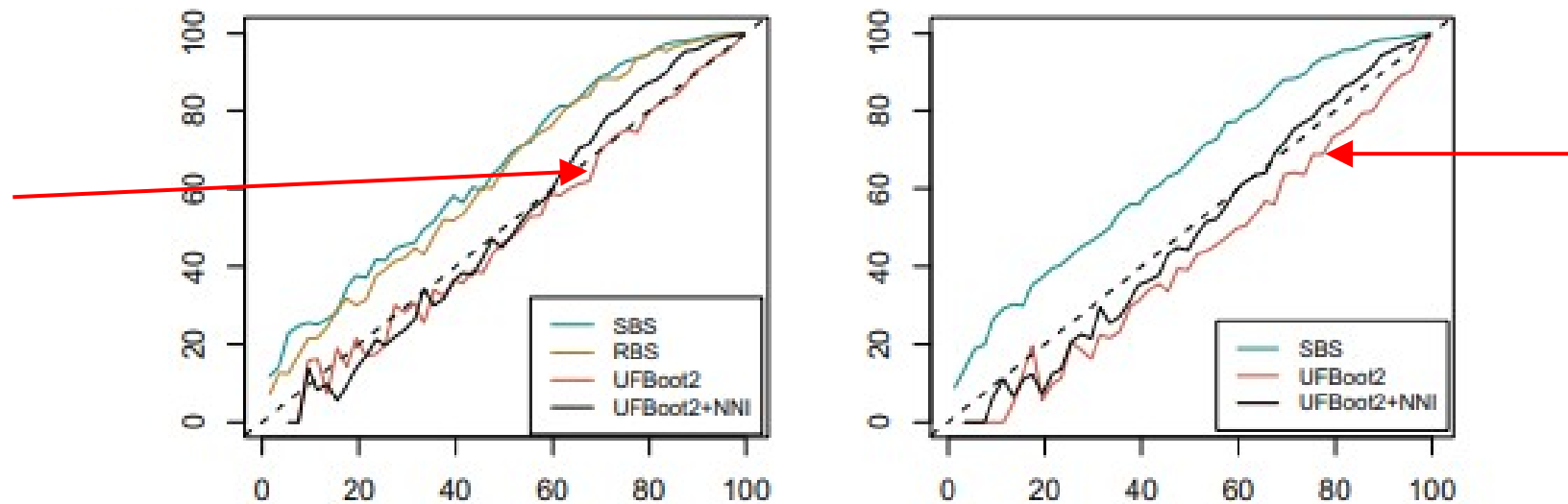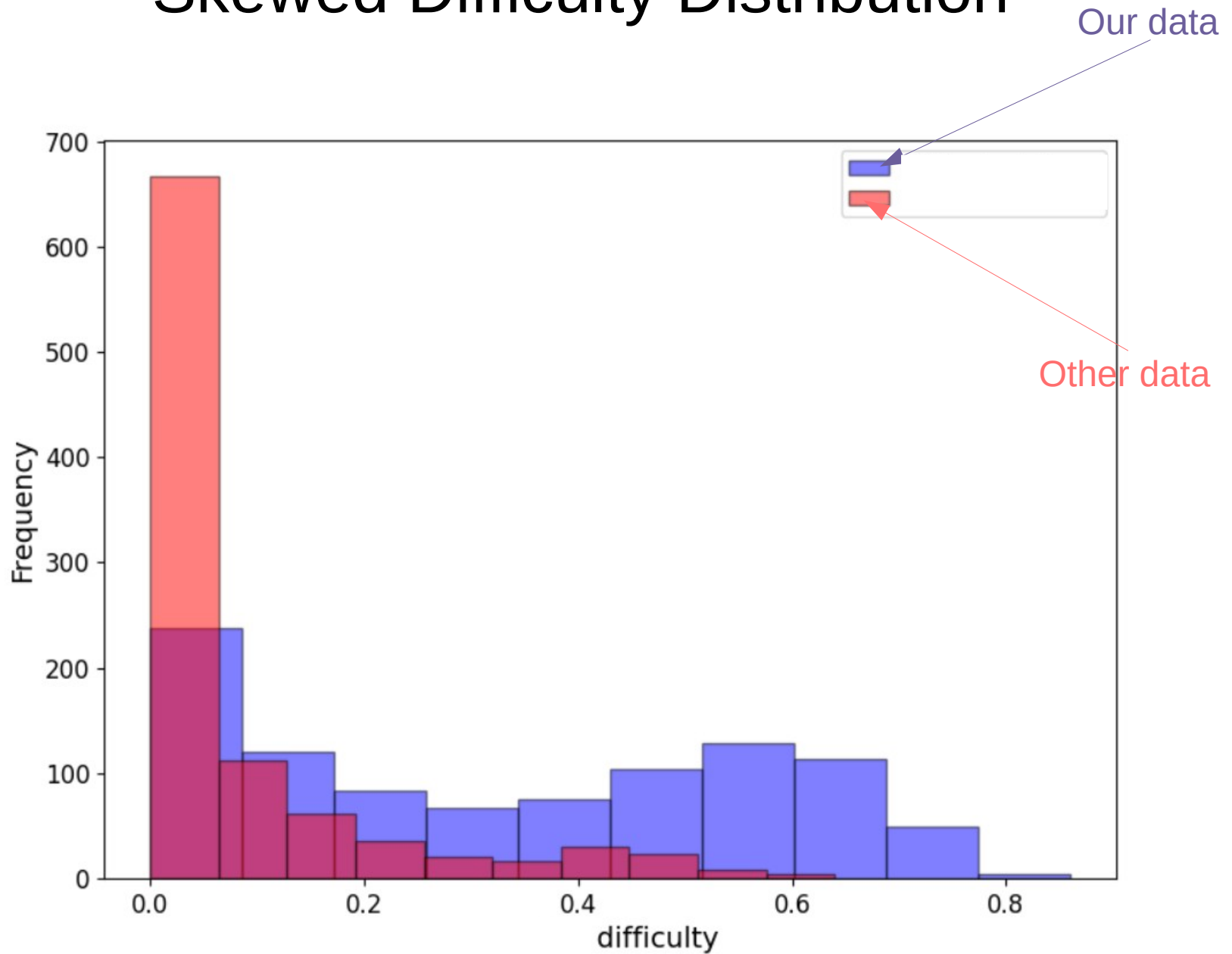# Related & Ongoing Work

- Rapidly predict phylogenetic support values

JOURNAL ARTICLE

**Predicting Phylogenetic Bootstrap Values via Machine Learning**

Julius Wiegert ✉, Dimitri Höhler, Julia Haag, Alexandros Stamatakis    Author Notes

*Molecular Biology and Evolution*, Volume 41, Issue 10, October 2024, msae215,
https://doi.org/10.1093/molbev/msae215
**Published:** 17 October 2024    **Article history** ▾

- Simulated DNA data sucks!

JOURNAL ARTICLE

**Simulations of Sequence Evolution: How (Un)realistic They Are and Why**

Johanna Trost, Julia Haag ✉, Dimitri Höhler, Laurent Jacob, Alexandros Stamatakis,
Bastien Boussau    Author Notes

*Molecular Biology and Evolution*, Volume 41, Issue 1, January 2024, msad277,
https://doi.org/10.1093/molbev/msad277
**Published:** 20 December 2023    **Article history** ▾

- Predict difficulty of Multiple Sequence Alignment
  - The step before phylogenetic inference
    - → almost done
- Franziska's part of the presentation

# Thank you for your attention



Listaros village, Crete

# Context-Aware Modeling of Phylogenetic Inference Difficulty

Franziska Reden

Biodiversity Computing Group
Mini Symposium, 27.02.2026

# Phylogenetic Pipeline using Difficulty Prediction

```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

Phylogenetic inference

Pythia

0.50
difficulty score

T1    T2
      T3
T4

# Phylogenetic Pipeline using Difficulty Prediction



```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

Pythia

0.50
difficulty score

Phylogenetic
inference

T1    T2

T3

T4

# Phylogenetic Pipeline using Difficulty Prediction

```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

Phylogenetic inference →

Pythia ↓

0.50

difficulty score

T1 T2

T3

T4

0     1

Difficulty

# Phylogenetic Pipeline using Difficulty Prediction

```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

Phylogenetic inference

T1    T2

T4    T3

Pythia

0.50
difficulty score

???

0                                    1

Difficulty

# Phylogenetic Pipeline using Difficulty Prediction



Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

Phylogenetic inference

T1    T2
T4    T3

Pythia

0.50
difficulty score

#columns

#rows

Taxon 1:ACG-TT-
Taxon 2:ACGTTT-
Taxon 3:ACCCT-G
Taxon 4:A-GGTTT

0.51
0.47
0.60

???

0        Difficulty        1

# Phylogenetic Pipeline using Difficulty Prediction
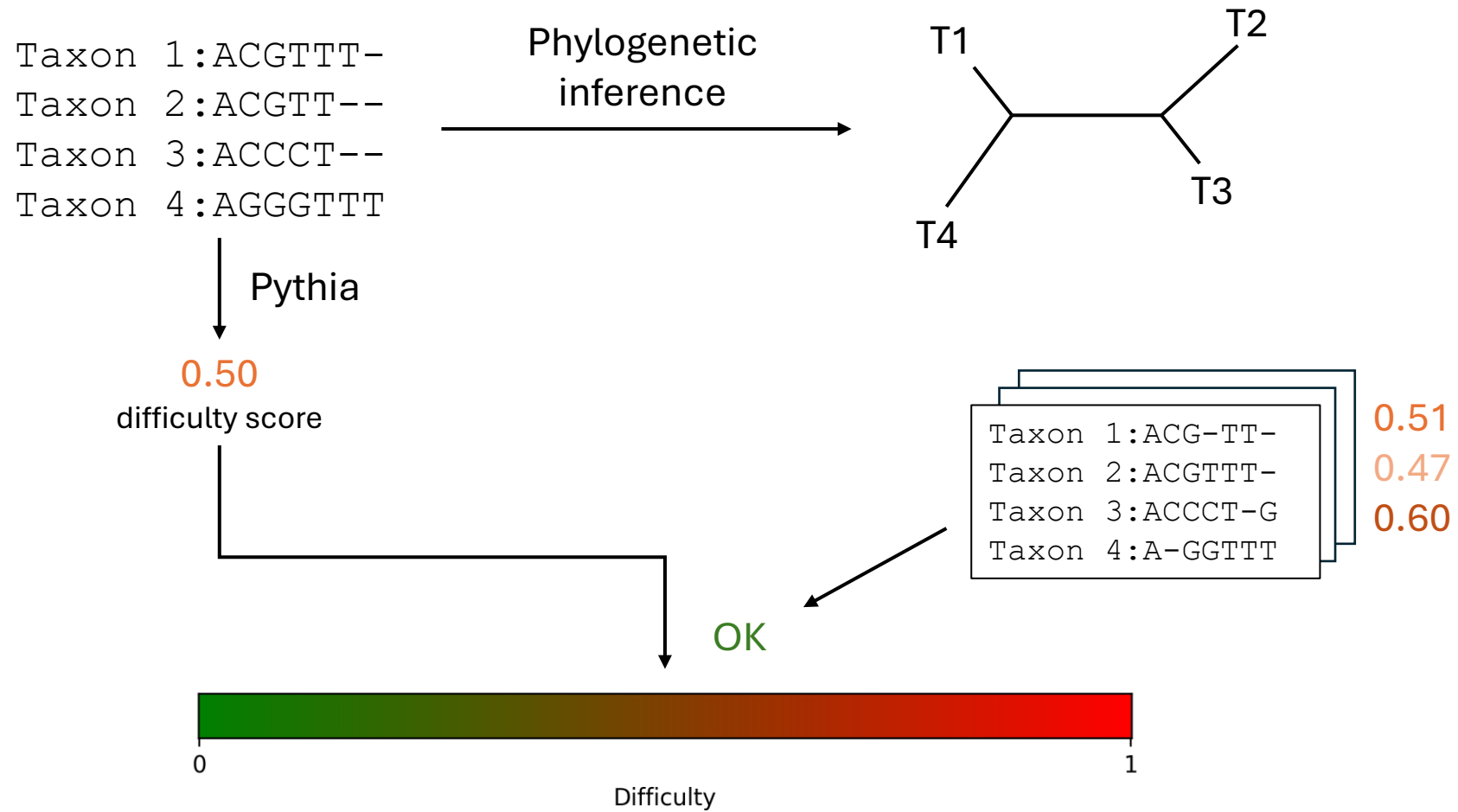
# Phylogenetic Pipeline using Difficulty Prediction



Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

Phylogenetic inference

T1  T2
T3
T4

Pythia

0.50
difficulty score

Taxon 1:ACG-TT-
Taxon 2:ACGTTT-
Taxon 3:ACCCT-G
Taxon 4:A-GGTTT

0.11
0.09
0.14

???

0                                    1
Difficulty

# Phylogenetic Pipeline using Difficulty Prediction

```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

Phylogenetic inference

T1      T2

T4      T3

Pythia

0.50
difficulty score

???

```
Taxon 1:ACG-TT-
Taxon 2:ACGTTT-
Taxon 3:ACCCT-G
Taxon 4:A-GGTTT
```
0.11
0.09
0.14

???

0                                    1
Difficulty

# Phylogenetic Pipeline using Difficulty Prediction

# Phylogenetic Pipeline using Difficulty Prediction
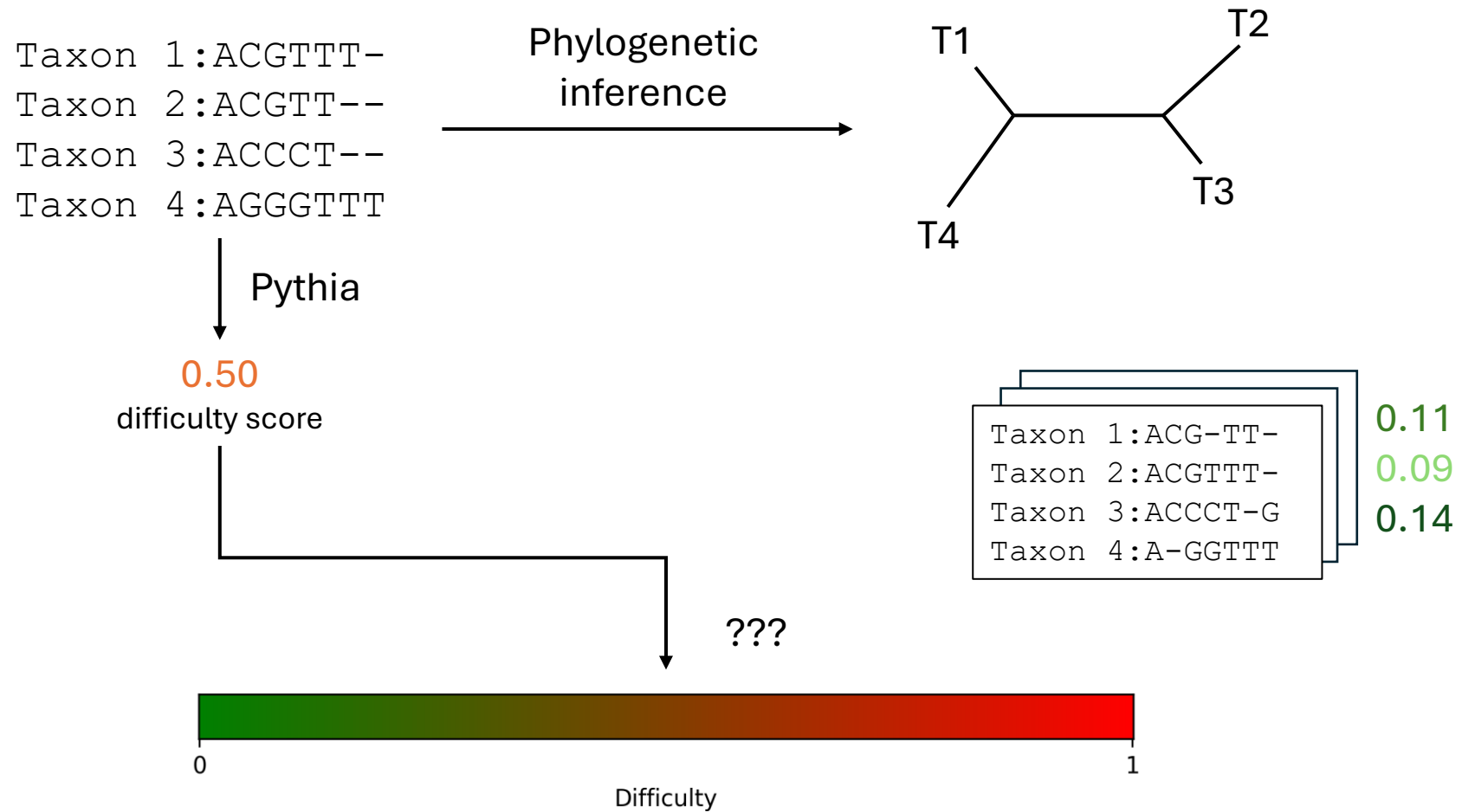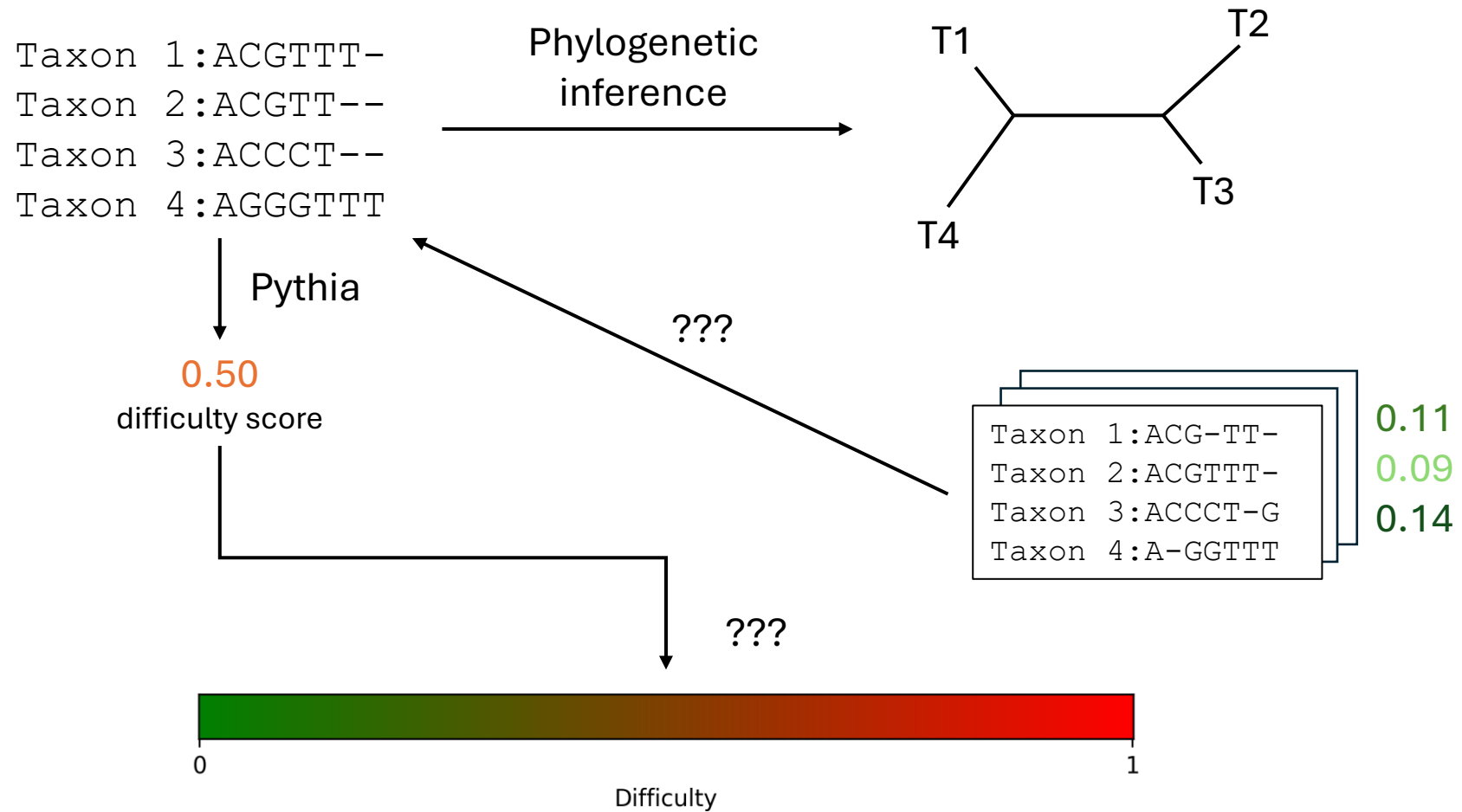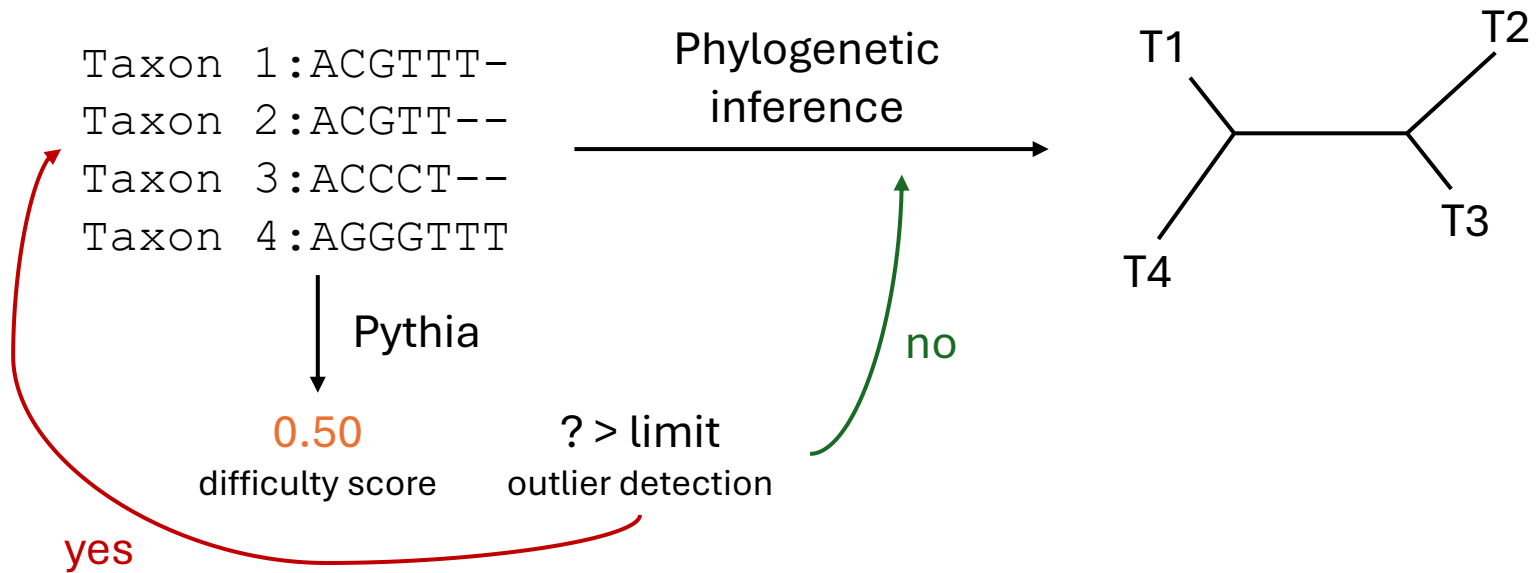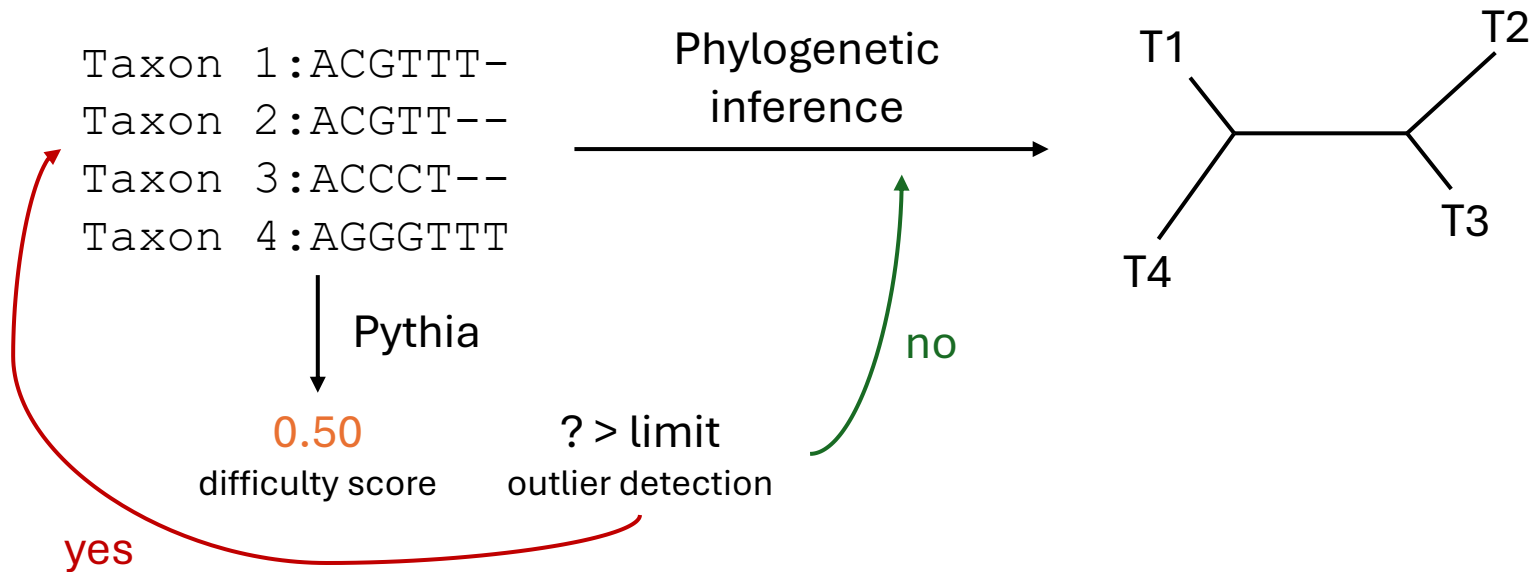
```
Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT
```

Pythia

0.50

difficulty score

? > limit

outlier detection

yes

no

Phylogenetic inference

T1
T2
T3
T4

The aim is to create a **context-aware (in terms of alignment features)** expected difficulty baseline

# Outlier detection using Quantile Regression Models

- **Aim:** Estimate conditional difficulty thresholds

- **Approach**: Model phylogenetic inference difficulty as a function of alignment-level features using quantile regression

$$Q_\tau(difficulty \mid alignment\ features)$$

- $\tau$ ...... quantile level
- $difficulty$ ...... difficulty as predicted with Pythia

# Outlier detection using Quantile Regression Models

- **Aim:** Estimate conditional difficulty thresholds

- **Approach**: Model phylogenetic inference difficulty as a function of alignment-level features using quantile regression

$$Q_\tau(difficulty \mid alignment\ features)$$

- $\tau$ ...... quantile level
- $difficulty$ ...... difficulty as predicted with Pythia

Example with $\tau = 0.95$:
"Given the feature set, 95% of comparable alignments have difficulty below X"

# Workflow

- Gather representative data sets

- Compute difficulty (target variable) using Pythia

- Choose and extract alignment features

- Choose and train regression models

- Evaluate models

# Workflow

- **Gather representative data sets**

- Compute difficulty (target variable) using Pythia

- Choose and extract alignment features

- Choose and train regression models

- Evaluate models

# Workflow: Gather representative data sets

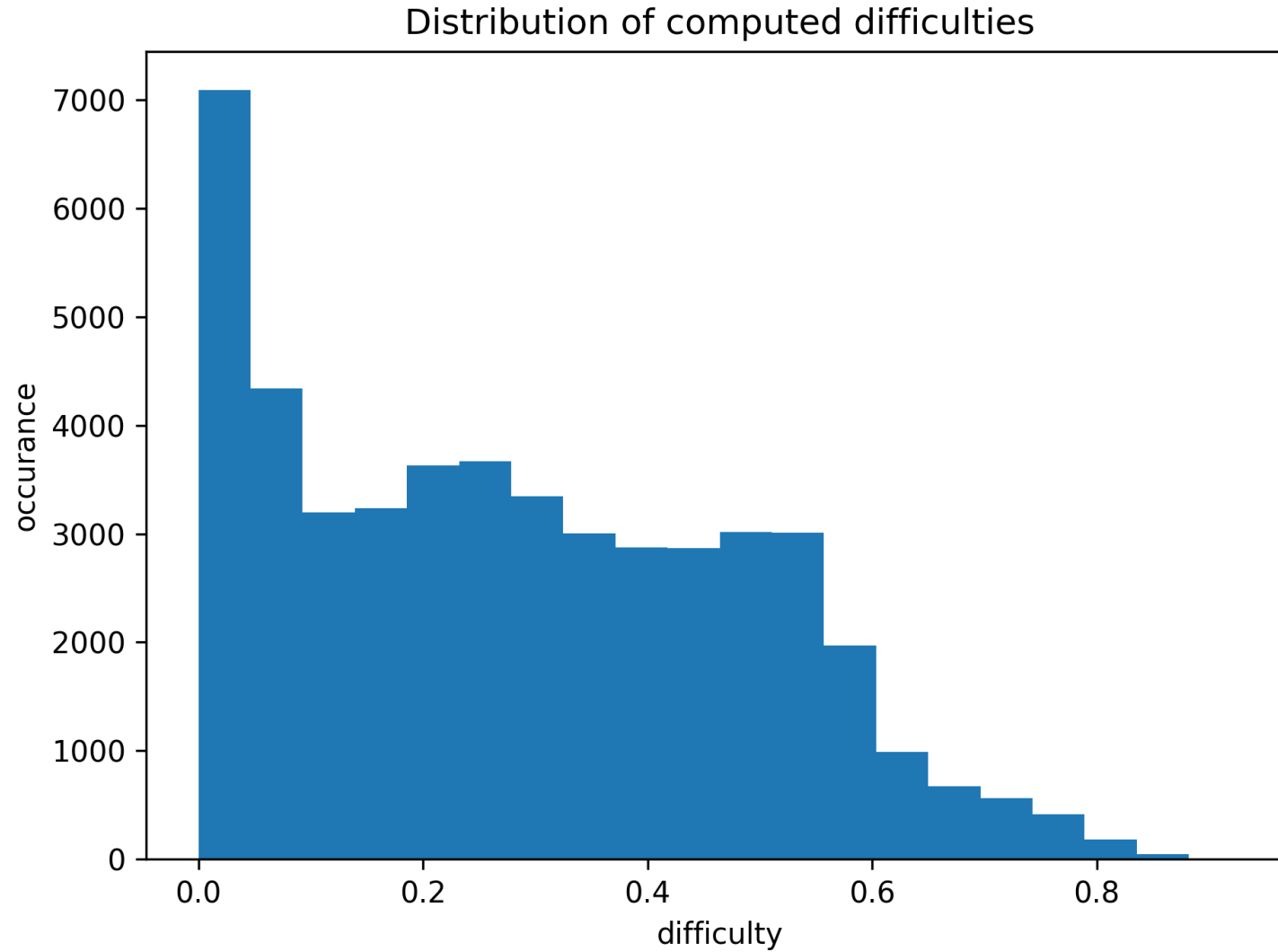Gathered **~50,000** data sets from published sources:

- TreeBASE database (Piel et al., 2002)

- PANDIT database (Whelan, 2006)

- Benchmark alignments of Rob Lanfear

- OrthoMaM database (Scornavacca et al., 2019; Allio et al., 2023)

# Workflow

- Gather representative data sets

- **Compute difficulty (target variable) using Pythia**

- Choose and extract alignment features

- Choose and train regression models

- Evaluate models

# Workflow: Compute difficulty



Distribution of computed difficulties

# Workflow

- Gather representative data sets

- Compute difficulty (target variable) using Pythia

- **Choose and extract alignment features**

- Choose and train regression models

- Evaluate models

# Workflow: Choose and extract alignment features

- Chosen features:
  - num_patterns/num_taxa
  - proportion_gaps
  - num_taxa
  - proportion_invariant
  - num_patterns
  - num_sites


- With an increased number of features:
  - Problem of sparsity in feature space
  - Tail estimation becomes unstable

# Workflow

- Gather representative data sets

- Compute difficulty (target variable) using Pythia

- Choose and extract alignment features

- **Choose and train regression models**

- Evaluate models

# Workflow: Choose and train regression models

Estimate conditional upper percentiles:

$$Q_\tau(difficulty \mid alignment\ features)$$

- Focus on:
  - Upper tail ($\tau=\{0.950, 0.975, 0.990\}$)
  - No distributional assumptions
  - Preserve empirical extremes
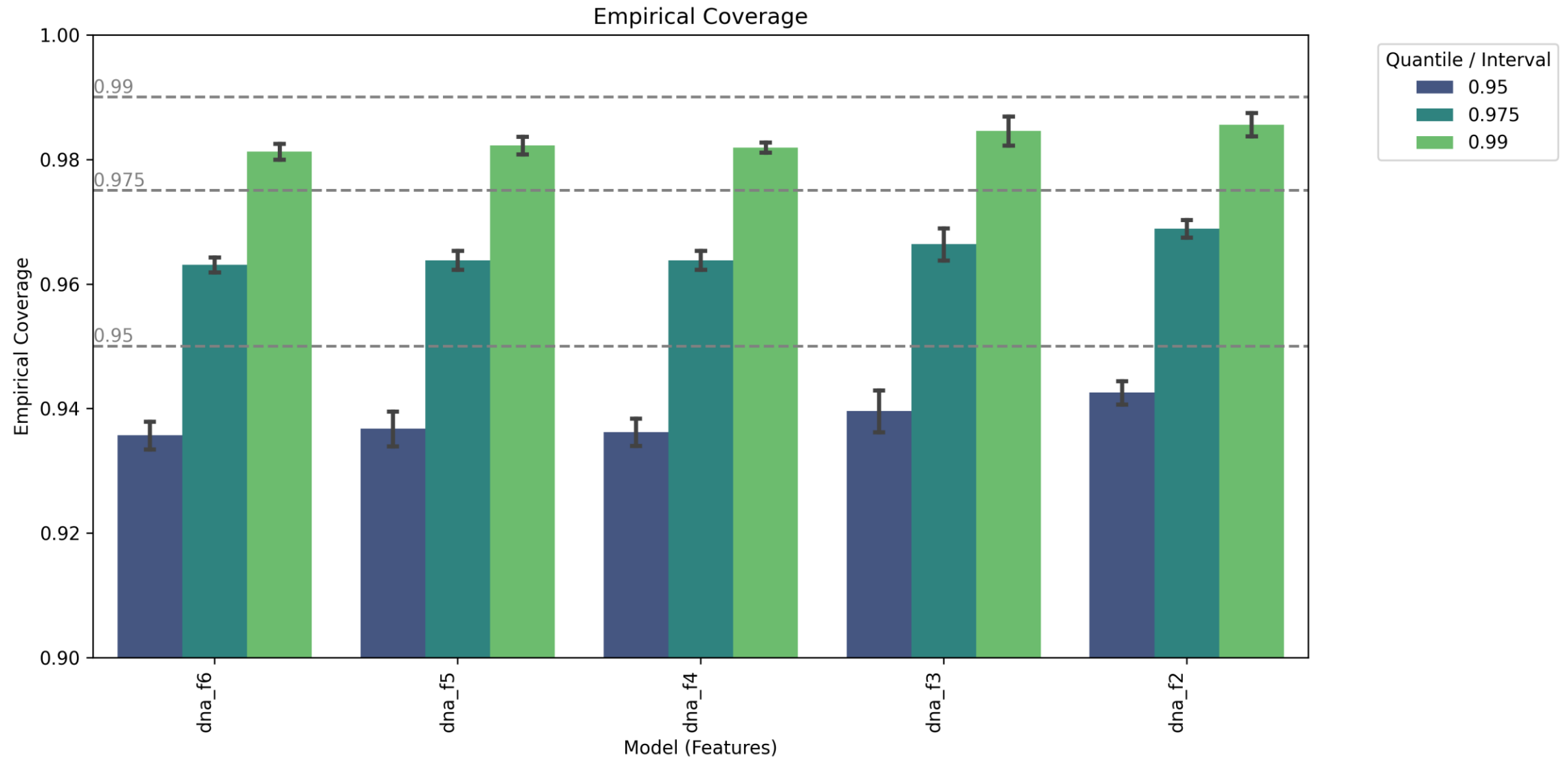
# Workflow: Choose and train regression models

Estimate conditional upper percentiles:

$$Q_\tau(difficulty \mid alignment\ features)$$

- Focus on:
  - Upper tail ($\tau=\{0.950, 0.975, 0.990\}$)
  - No distributional assumptions
  - Preserve empirical extremes

- **Gradient Boosted Trees**:
  - flexible, non-parametric, robust to correlated predictors

# $Q_\tau(\textit{difficulty} \mid \textit{alignment features}) \approx$ observed $\tau$−th percentile in reference data
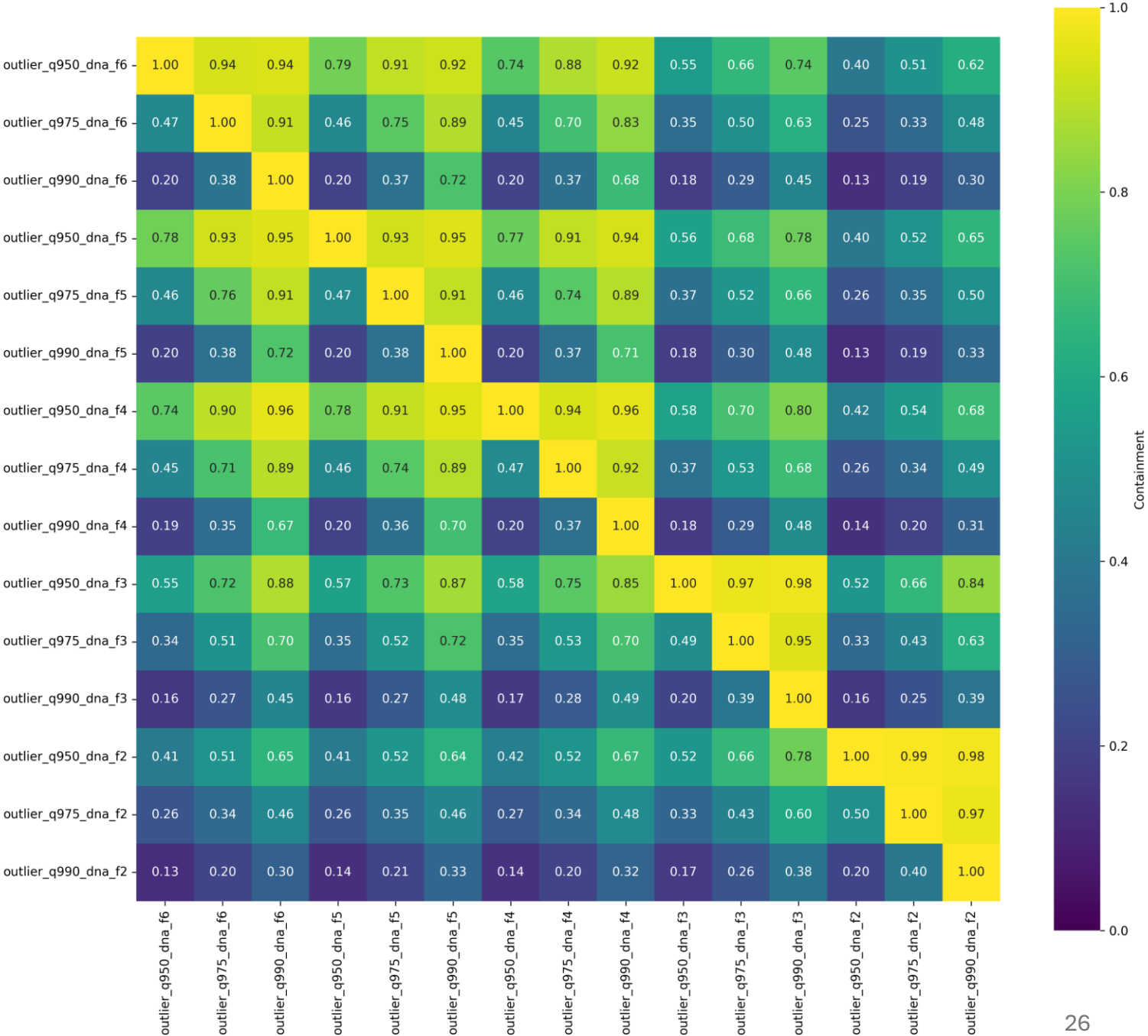
**Outlier containment:**

$$C(\text{X} \subseteq \text{Y}) = \frac{|X \cap Y|}{|X|}$$

**Outlier containment:**

Comparison between models

# Workflow

- Gather representative data sets

- Compute difficulty (target variable) using Pythia

- Choose and extract alignment features

- Choose and train regression models

- **Evaluate models**

# Workflow

- Gather representative data sets

- Compute difficulty (target variable) using Pythia

- Choose and extract alignment features

- Choose and train regression models

- Evaluate models: Work in progress
  - First results show high potential of models trained on **3 features**

# Overview

- **Objective**: Estimate conditional difficulty thresholds based on empirical data using quantile regression models

# Overview

- **Objective**: Estimate conditional difficulty thresholds based on empirical data using quantile regression models

- **Concerns**:

  - Extreme quantile instability (very few effective observations)
  - Feature-space sparsity (choice of features)
  - Quantile crossing risk (violation of monotonicity)
  - Empirical bound limitation

# Overview

- **Objective**: Estimate conditional difficulty thresholds based on empirical data using quantile regression models

- **Concerns**:

  - Extreme quantile instability (very few effective observations)
  - Feature-space sparsity (choice of features)
  - Quantile crossing risk (violation of monotonicity)
  - Empirical bound limitation

- **Outlook**:
  - Complete evaluation of models
  - Increase training data size
  - Additional measure of data representation