# Predicting Multiple Sequence Alignment Uncertainty via Machine Learning
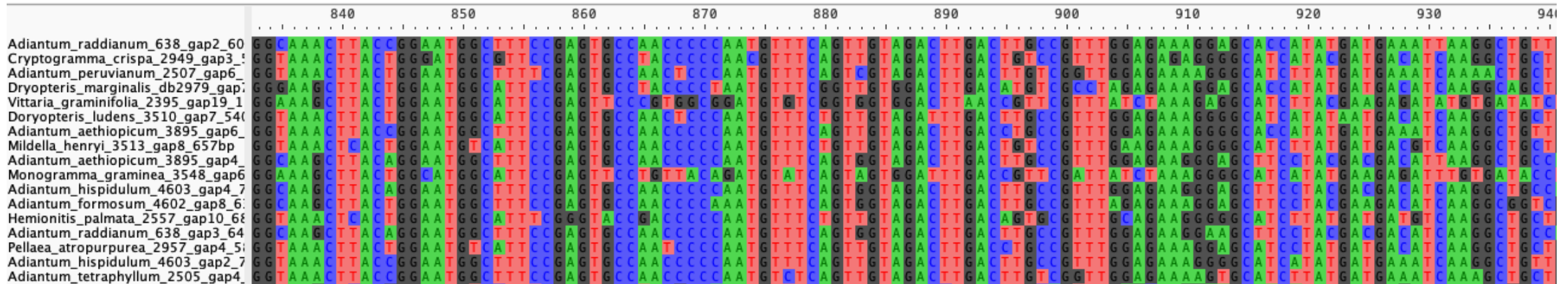
Mattis Bodynek, **Lucía Martín-Fernández**, Julia Haag,
Ben Bettisworth, Alexandros Stamatakis

LEGEND 2025. Machine Learning for Evolutionary Genomics Data

8th-12th December

# Multiple Sequence Alignment

**Definition**   Process to identify regions of similarity across three or more biological sequences



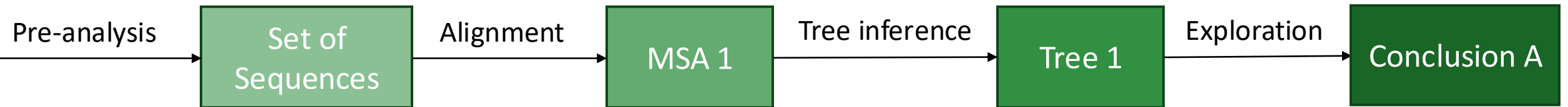**How?**   Heuristic and probabilistic algorithms.

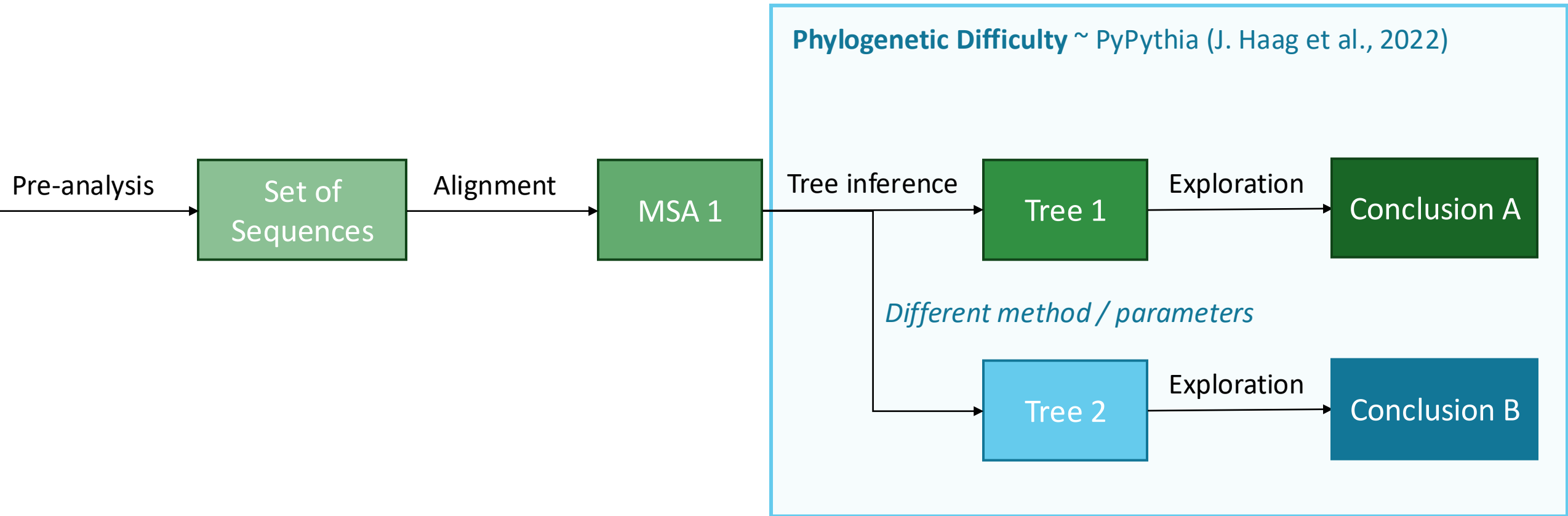Each tool relies on different model assumptions.

**Problem**   NP-Hard.

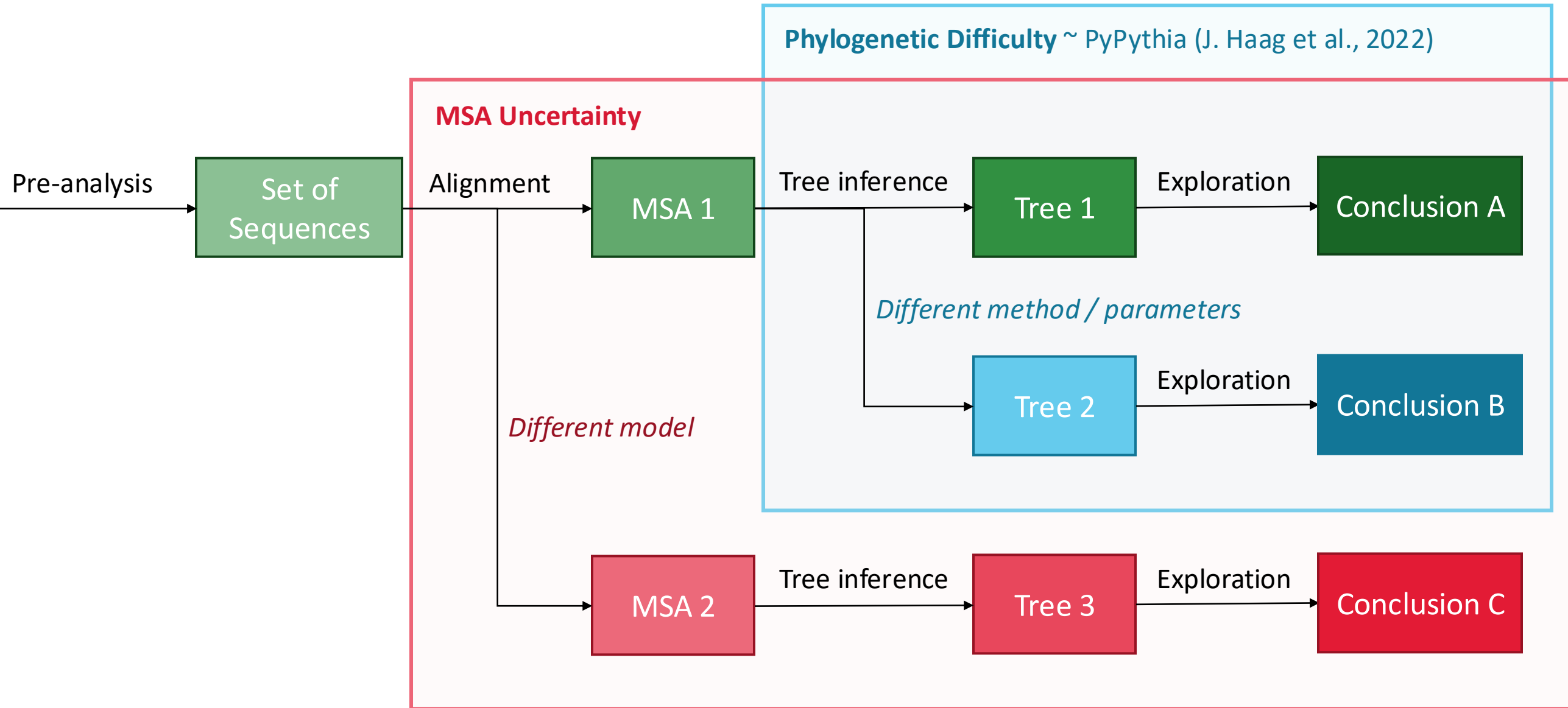We can only approximate the optimal solution.

# MSA Uncertainty
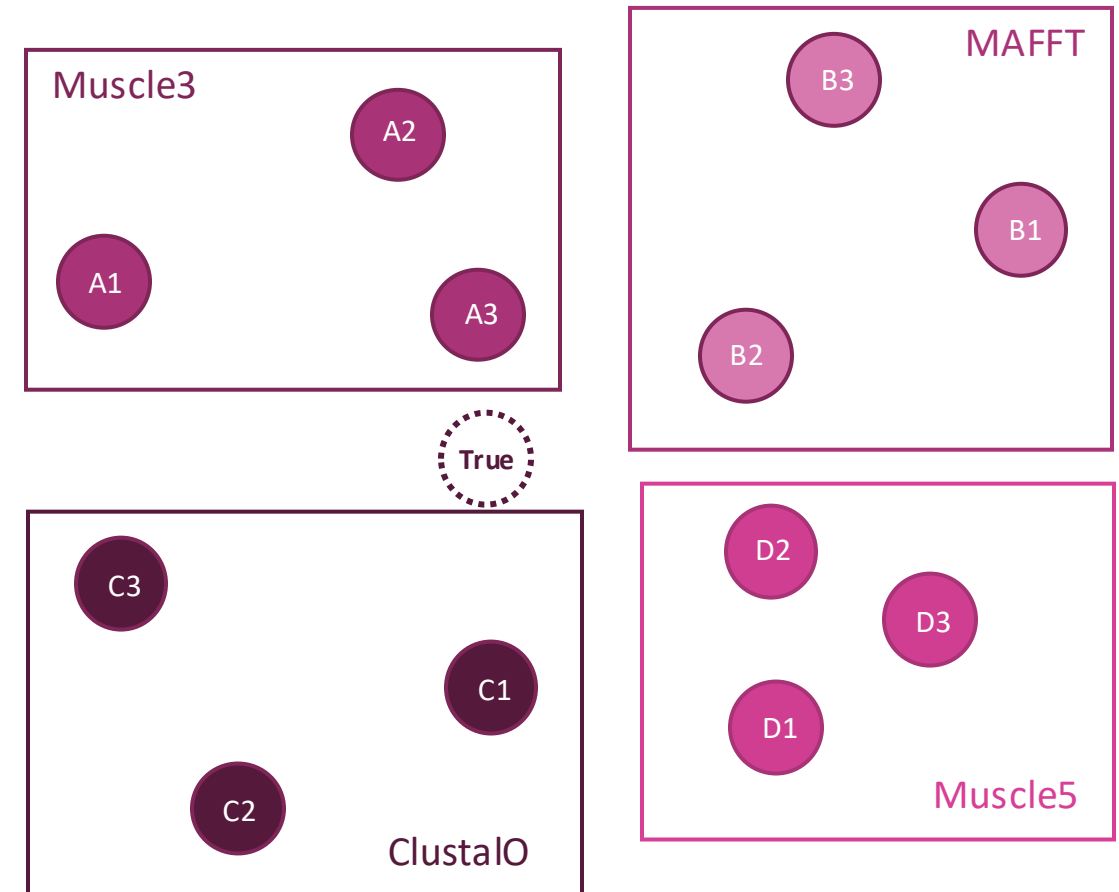
# MSA Uncertainty

# MSA Uncertainty

# MSA Uncertainty

**So far, we know that…**

o Multiple Sequence Alignment is an NP-Hard problem

o Different alignment tools with different parameters may generate different solutions

**MSA Uncertainty ~ MSA Difficulty**

- How difficult is to generate a stable MSA?
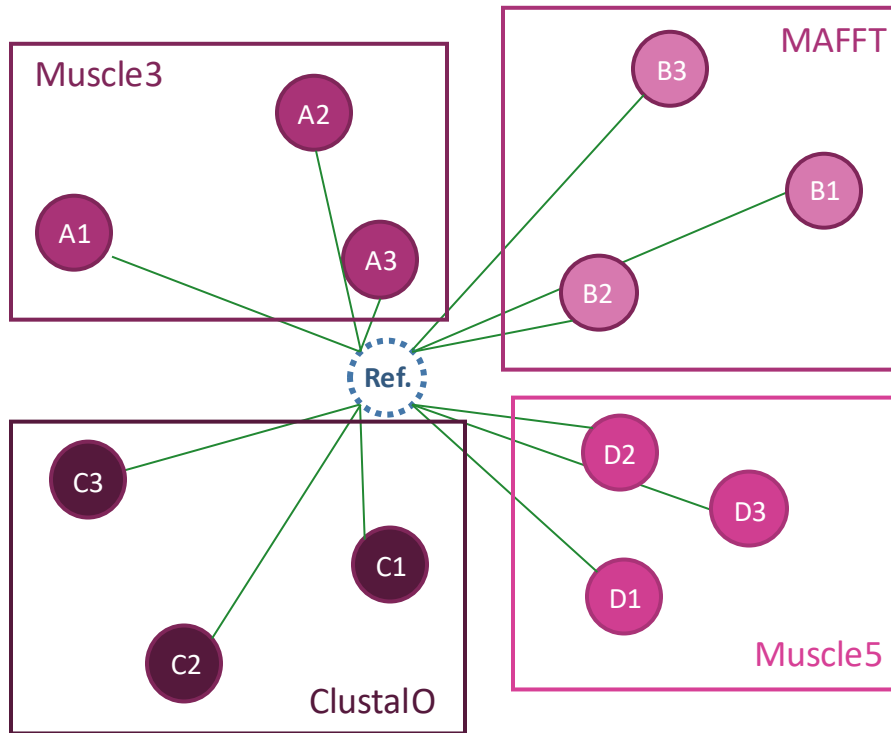- How much uncertainty exists in how accurately the MSA captures true homology?
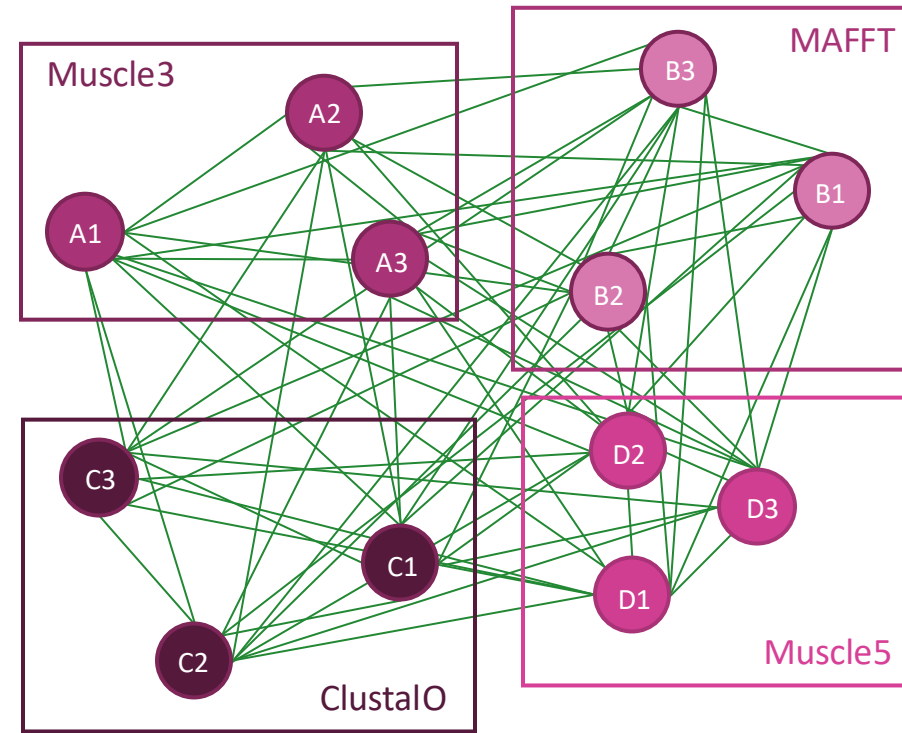


**ENSEMBLE**

Different MSAs generated from one set of sequences using various algorithms and parameters

# MSA Uncertainty

**MSA Uncertainty score**


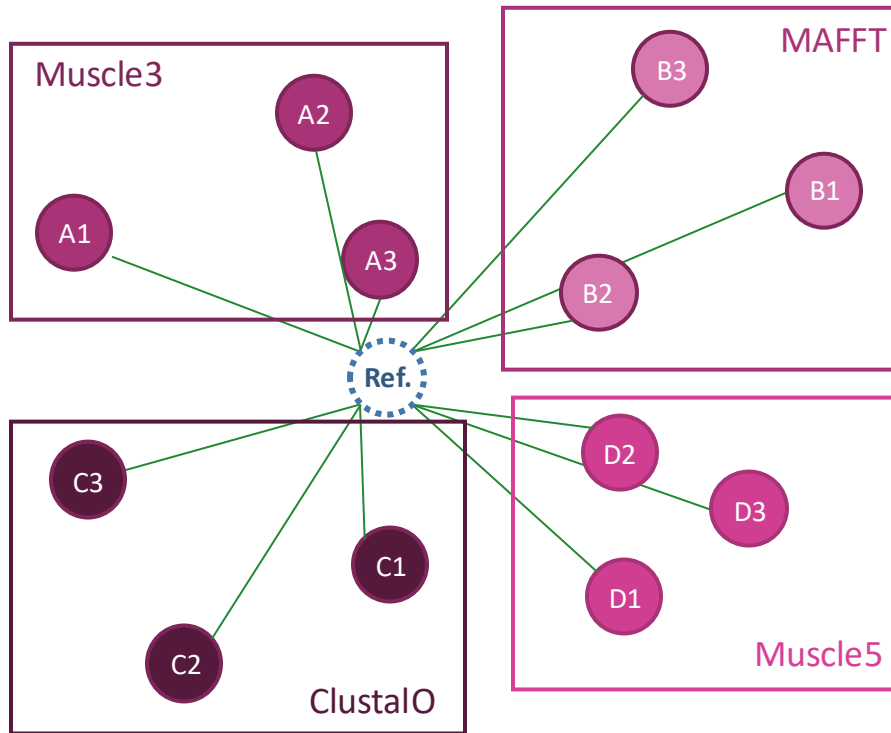
**Reference-Based score**

**Reference-free score**
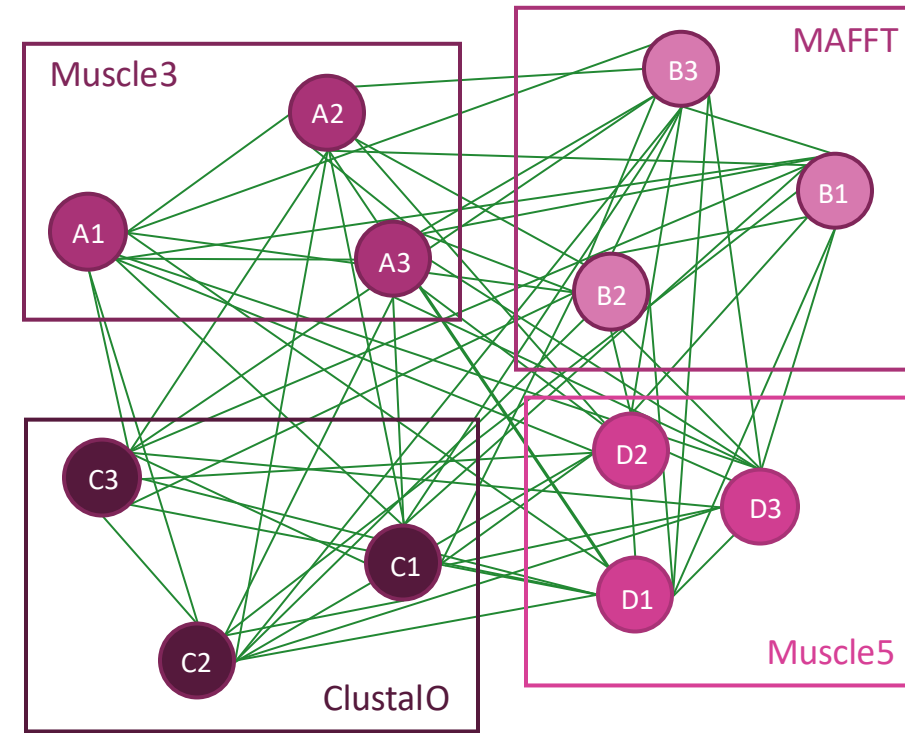
**BAliBASE:** structural benchmark database

# MSA Uncertainty

**MSA Uncertainty score**

We compared different **distance metrics** using **BAliBASE** (structural benchmark database) as reference.



**Reference-Based score**

**Reference-free score**

$d_{pos}$
- Uses positional homology sets of the alignments
- Incorporates positional information about gaps
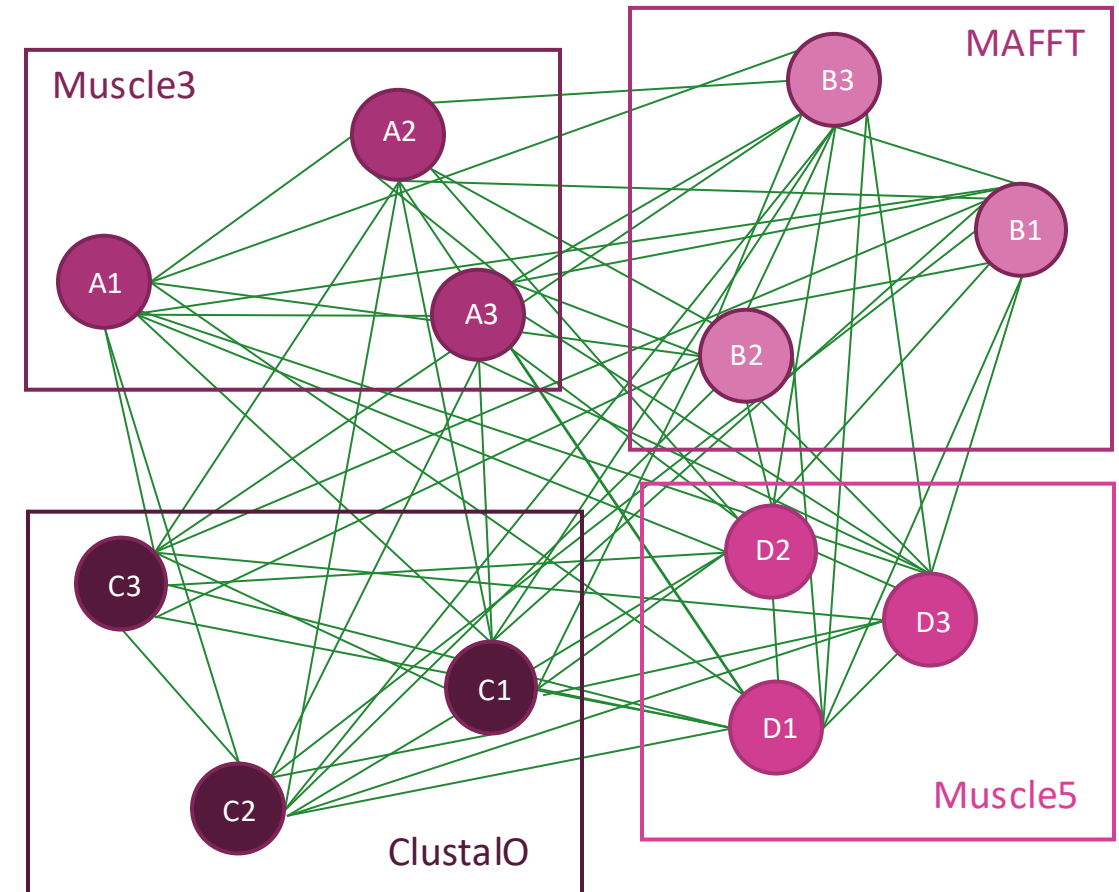
# MSA Uncertainty

## MSA Uncertainty score

We measure how difficult it is to align a set of sequences:



0            1

*easy*            *difficult*

Quantifies how much alignments differ within an ensemble

~ average norm. pairwise distance ($d_{pos}$) between all MSAs.

---

**MSA Uncertainty ~ MSA Difficulty**

- How difficult is to generate a stable MSA?
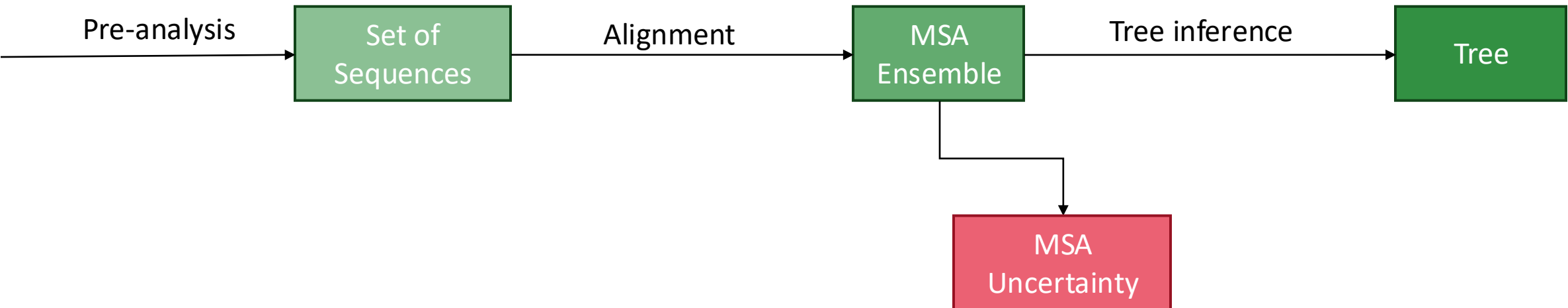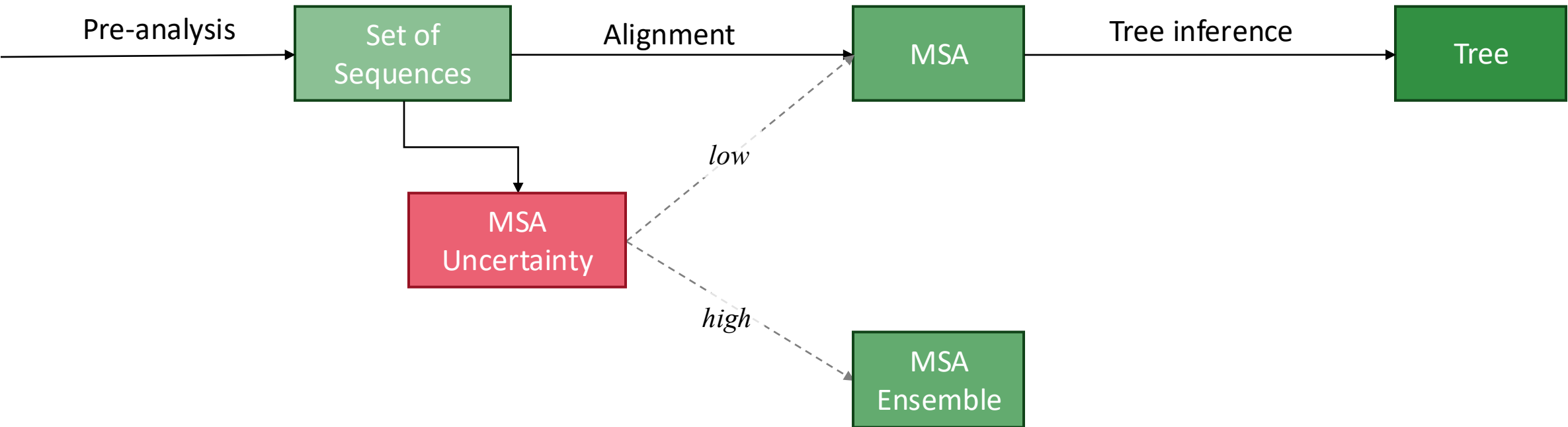- How much uncertainty exists in how accurately the MSA captures true homology?

---



**ENSEMBLE**

Different MSAs generated from one set of sequences using various algorithms and parameters

# MSA Uncertainty

# MSA Uncertainty

# Predicting MSA uncertainty

1. **Data collection**

2. **Label Generation**

3. **Feature Generation**

4. **Training the model**

5. **Results**

# Predicting MSA uncertainty

## 1. Data collection

| | |
|---|---|
| HOMSTRAD<br>SABmark<br>BAliBASE v3<br>PREFAB v2<br>OXBENCH<br>ArthropodsP450 | AA |
| TreeBASE | AA / DNA |
| BRAliBASE<br>BAliBASE v2 | RNA → DNA<br>AA → DNA |

11.432 sequence sets

# Predicting MSA uncertainty

1. **Data collection**

   11.432 sequence sets

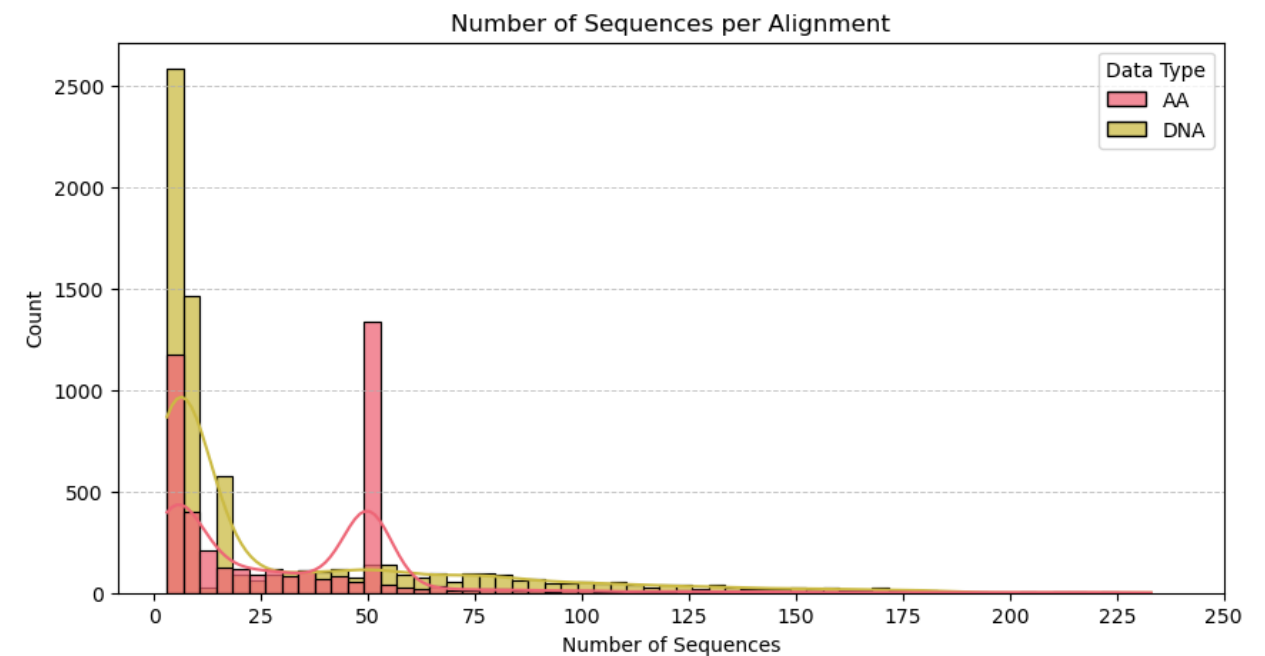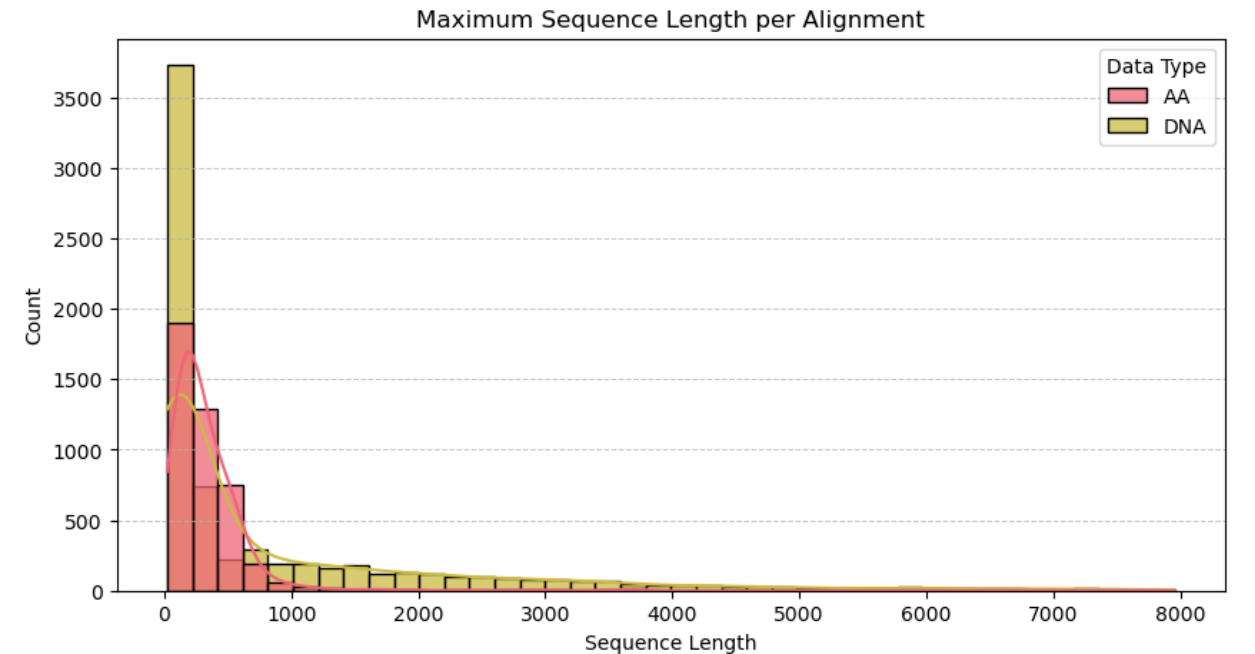2. **Label Generation**

3. **Feature Generation**

4. **Training the Model**

5. **Results**

# Predicting MSA uncertainty

1. **Data collection**

   11.432 sequence sets

2. **Label Generation**

   We calculated the MSA uncertainty score heuristically for the collected data.

   We generate ensembles of 48 alignments per sequence set.

3. **Feature Generation**

4. **Training the Model**

5. **Results**

# Predicting MSA uncertainty

1.  **Data collection**

    11.432 sequence sets

2.  **Label Generation**

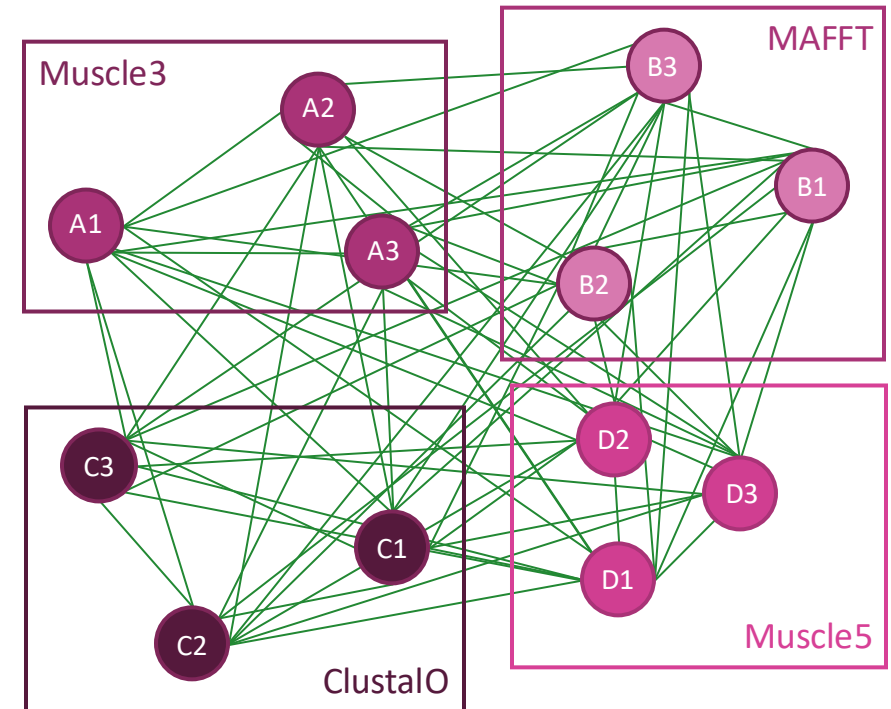    We calculated the MSA uncertainty score heuristically for the collected data.

    We generate ensembles of 48 alignments per sequence set.

3.  **Feature Generation**

    We define inexpensive to compute features on the unaligned sequences.
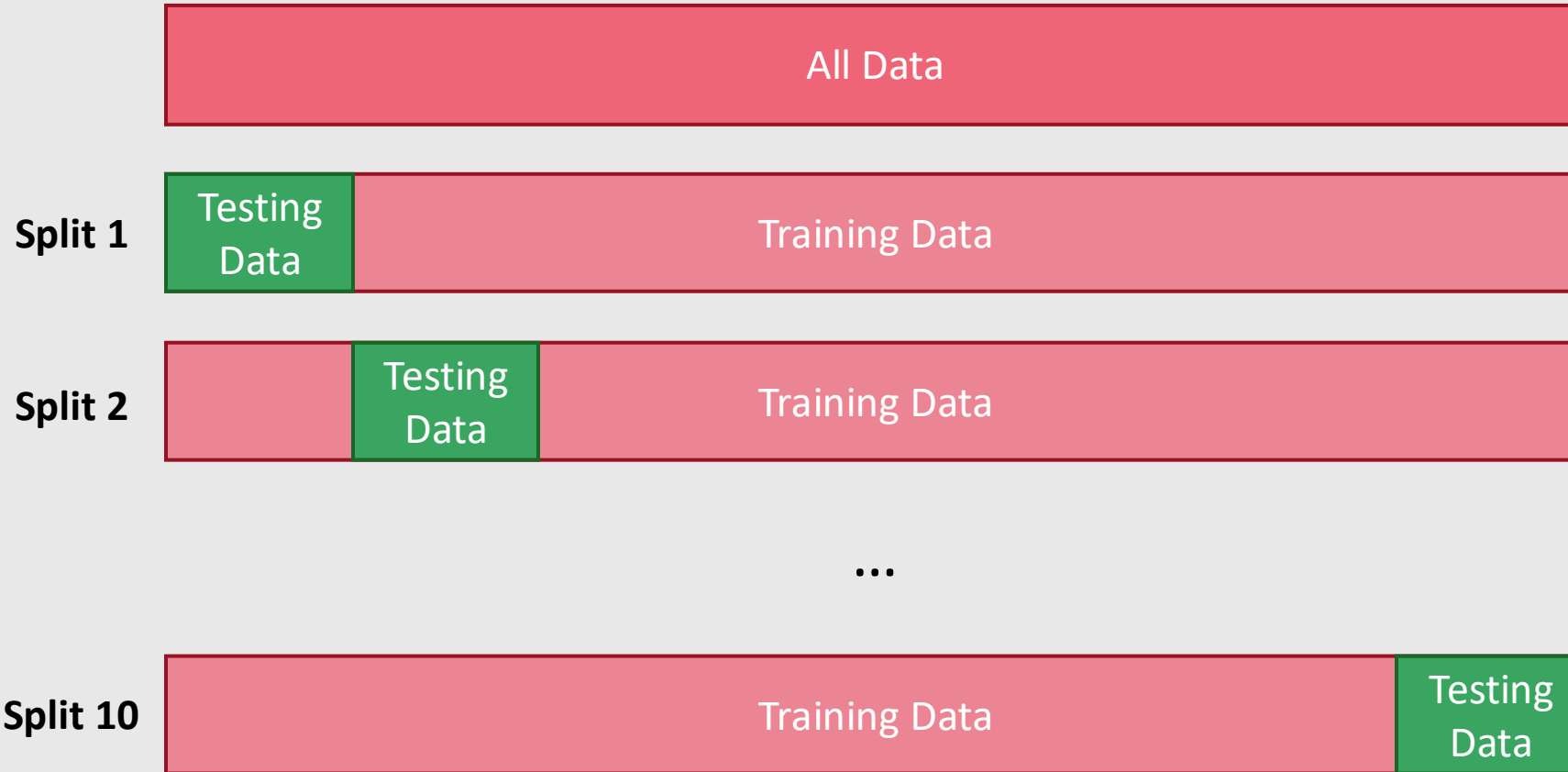
    The majority of our features are stochastic because they subsample sequences.

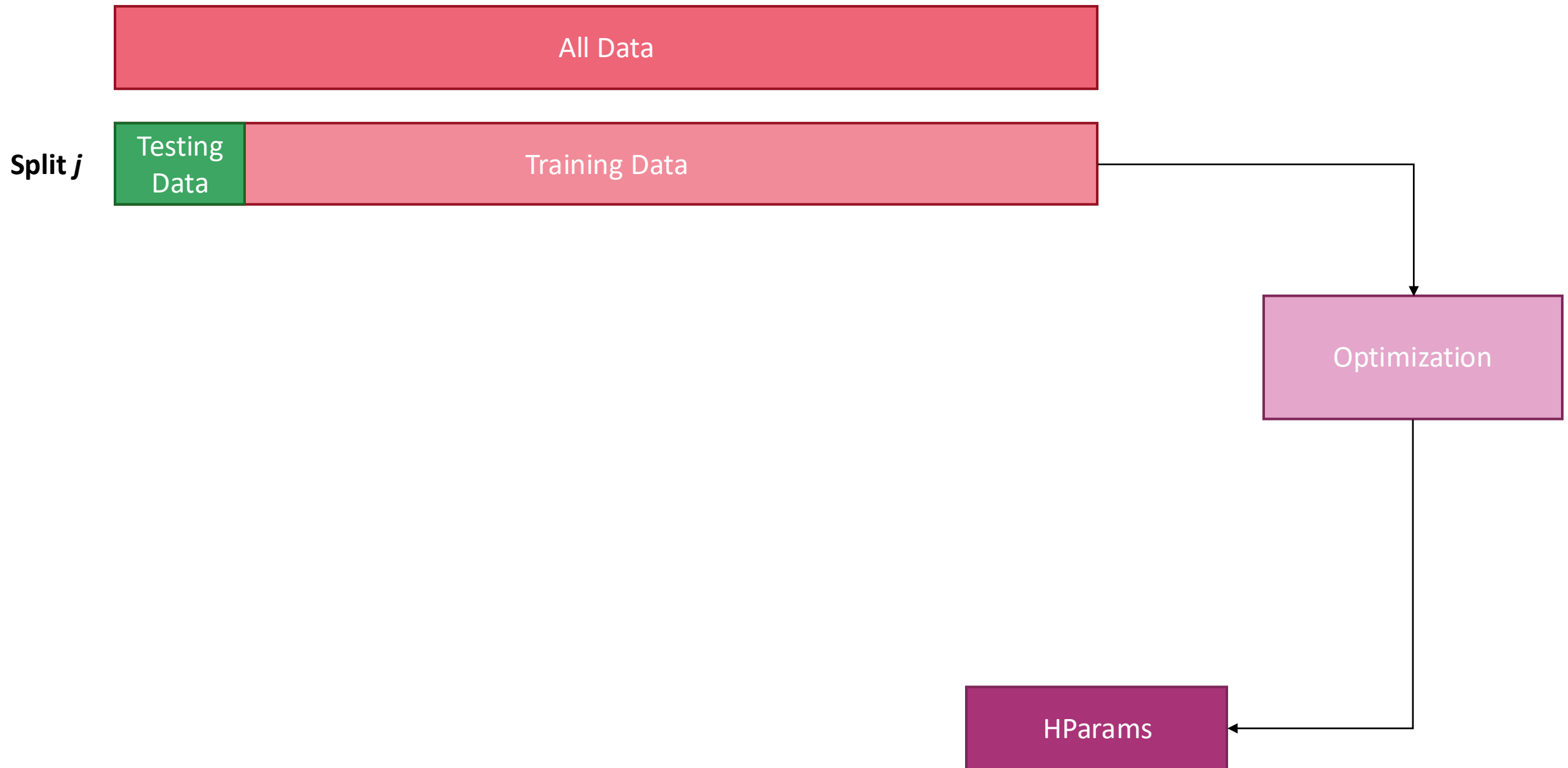4.  **Training the Model**

5.  **Results**

# Training the Model



X10 sampling seeds for stochastic Feature Generation

**10 splits x 10 sampling seeds = 100 Folds**

# Training the Model

# Training the Model

# Training the Model

# Training the Model

# Training the Model



Split 1 | Testing Data | Training Data | HParams | RMSE, $R^2$

Split 2 | Testing Data | Training Data | HParams | RMSE, $R^2$

...

Split 10 | Training Data | Testing Data | HParams | RMSE, $R^2$

**10 splits X 10 sampling seeds = 100 Folds**

**Final Model Training**

**Min(RMSE)**

**Summary($R^2$)**

Final LightGBM regressor

Best HParams

Estimated Performance

13

# Final Model

**Best HParams**

The set of optimized hyperparameters across all 100 folds with the lowest RMSE

**Estimated Performance**

**$R^2$ = 0.945 [0.939–0.951]**
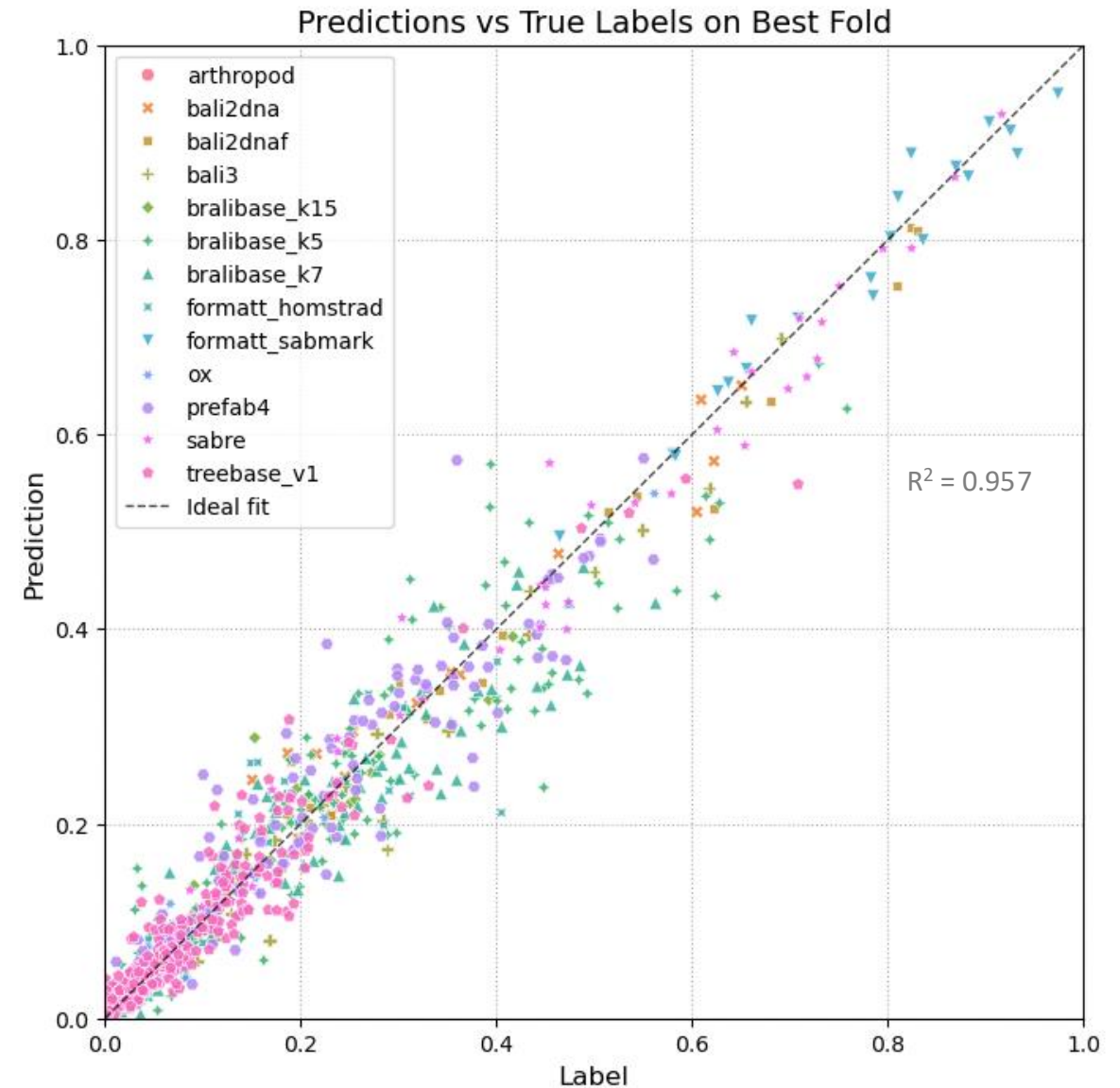$R^2$ median and its 10% and 90% percentiles across all 100 folds.

**Final LightGBM regressor**

Trained using all available data and the best hyperparameter set



Predictions vs True Labels on Best Fold

Legend:
- arthropod
- bali2dna
- bali2dnaf
- bali3
- bralibase_k15
- bralibase_k5
- bralibase_k7
- formatt_homstrad
- formatt_sabmark
- ox
- prefab4
- sabre
- treebase_v1
- ---- Ideal fit

$R^2$ = 0.957

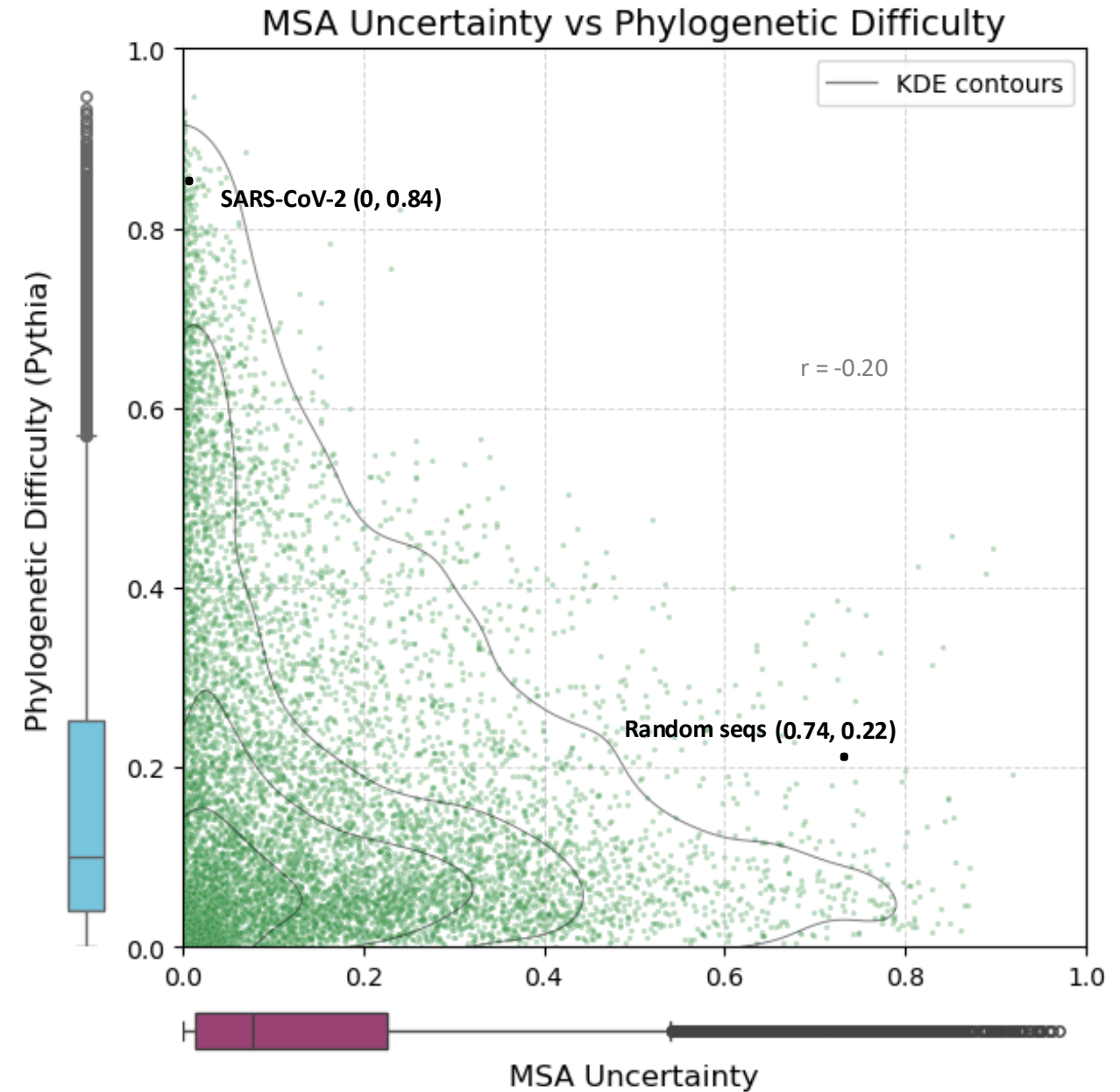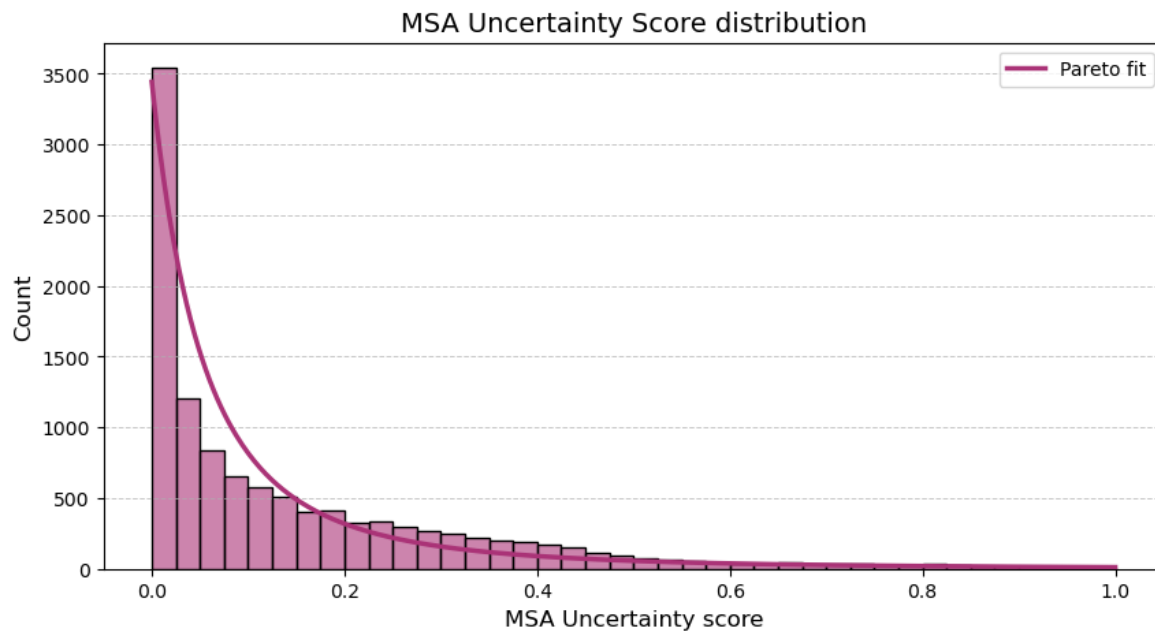Axis labels: Prediction (y-axis), Label (x-axis)

# MSA Uncertainty vs Phylogenetic Difficulty

**SARS-Cov-2.** 4869 sequences

- MSA Uncertainty = 0
- Phylogenetic Difficulty = 0.84

20 **Random DNA seqs** (500 - 525bp)

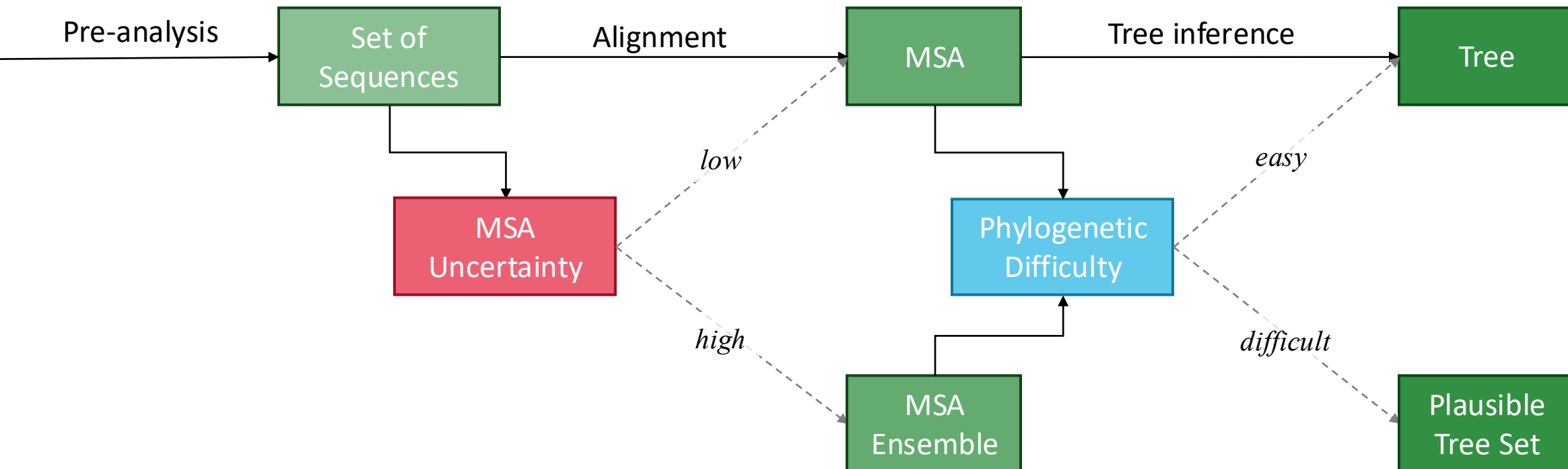- MSA Uncertainty = 0.74
- Phylogenetic Difficulty = 0.22

# Summary

## Conclusion

- Heuristic method to calculate MSA Uncertainty
  ~ average norm. pairwise distance ($d_{pos}$) between all 48 MSAs in an ensemble

- **LightGBM Regressor** to predict MSA Uncertainty given a set of sequences **(R²=0.945)**

- **Inverse correlation** between Phylogenetic Difficulty and MSA uncertainty.

# Availability

**Paper Loading . . .**



https://github.com/MaBody/aldiscore

📖 README     ⚖️ GPL-3.0 license       ✏️ ☰

## AlDiScore - Alignment Difficulty Score

AlDiScore provides two approaches for quantifying multiple sequence alignment (MSA) difficulty:

1. **Heuristic Scoring**: Compute dispersion within an ensemble of alternative alignments
2. **Predictive Scoring**: Predict alignment difficulty from unaligned sequences using ML

## Features

- Command-line interface for heuristics and prediction
- Multiple scoring methods for ensemble analysis
- Pre-trained models for difficulty prediction
- Supports DNA and amino acid sequences

**FORTH**
INSTITUTE OF COMPUTER SCIENCE

## Biodiversity Computing Group

- **Ben Bettisworth**
- **Lucía Martín Fernández**
- Georgios Koutsovoulos

- Panos Ioannidis
- Franziska Reden
- Noah Wahl

**HITS**
Heidelberg Institute for
Theoretical Studies

## Exelixis Lab

- **Mattis Bodynek**
- **Julia Haag**
- Alexey Kozlov
- Benoit Morel
- Christoph Stelz

- Lukas Hübner
- Anastasis Togkousidis
- Dimitri Höhler
- Luise Häuser
- Johannes Hengstler

Professor **Alexandros Stamatakis**