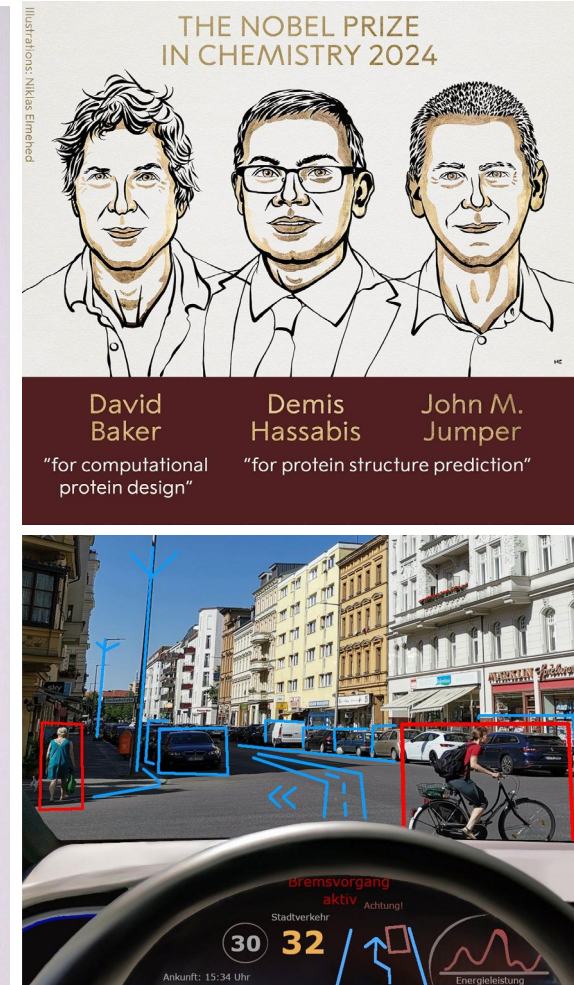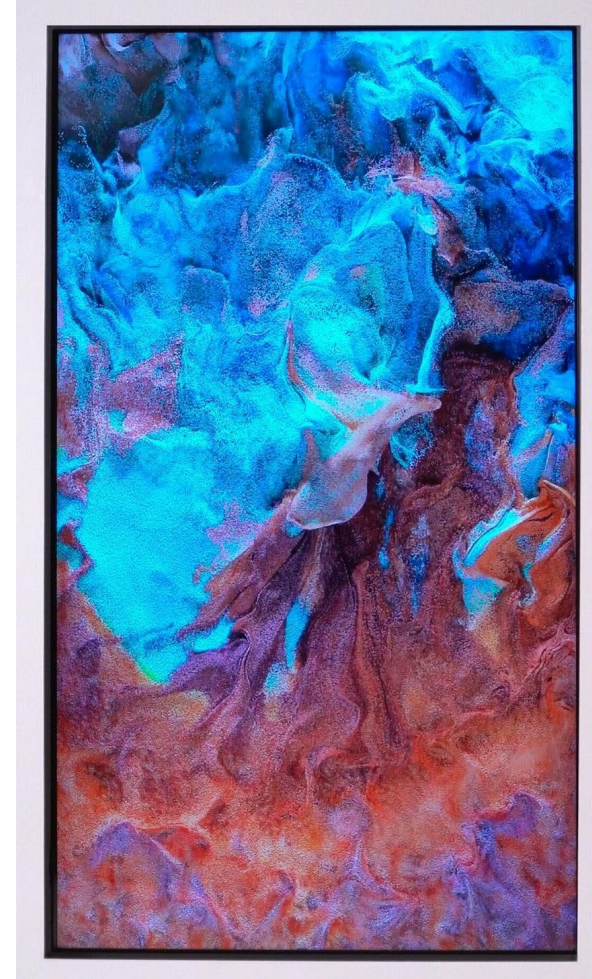# A machine-learning based alternative to phylogenetic bootstrap

Tal Pupko,
The Shmunis School of Biomedicine and Cancer Research
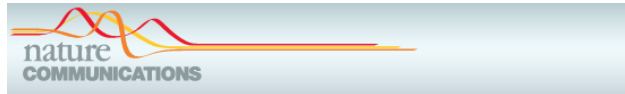
IMSI, Aug 12 2025
Chicago, USA

- Large Language Models, e.g., ChatGPT
- Autonomous systems (self-driving cars, drones)
- Creativity - first Christie's AI art sale
- Drug discovery and health care, AlphaFold2

# Our research aims to harness AI to improve phylogenetic inference

ARTICLE

https://doi.org/10.1038/s41467-021-22073-8    OPEN

## Harnessing machine learning to guide phylogenetic-tree search algorithms

Dana Azouri [1,2], Shiran Abadi [1], Yishay Mansour[3], Itay Mayrose [1] & Tal Pupko [2]

Phylogenetics

## BetaAlign: a deep learning approach for multiple sequence alignment

Edo Dotan [1,2], Elya Wygoda[1], Noa Ecker[1], Michael Alburquerque[1], Oren Avram [3], Yonatan Belinkov[2,*], Tal Pupko [1,*]

[1]The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 69978, Israel
[2]The Henry and Marilyn Taub Faculty of Computer Science, Technion—Israel Institute of Technology, Haifa 3200003, Israel
[3]The Department of Computer Science, University of California Los Angeles, Los Angeles, CA 90095, United States

## The Tree Reconstruction Game: Phylogenetic Reconstruction Using Reinforcement Learning

Dana Azouri [1,2,†], Oz Granit [3,†], Michael Alburquerque [2,†], Yishay Mansour [3,*], Tal Pupko [2,*] and Itay Mayrose [1,*]

[1]School of Plant Sciences and Food Security, Tel Aviv University, Ramat Aviv, Tel Aviv 69978, Israel
[2]The Shmunis School of Biomedicine and Cancer Research, Tel Aviv University, Ramat Aviv, Tel Aviv 69978, Israel
[3]Balvatnik School of Computer Science, Tel Aviv University, Ramat Aviv, Tel Aviv 69978, Israel
[†]These author contributed equally.

*Corresponding authors: E-mails: itaymay@tauex.tau.ac.il; talp@tauex.tau.ac.il; mansour@tauex.tau.ac.il.
Associate editor: Andrey Rzhetsky

### Abstract

The computational search for the maximum-likelihood phylogenetic tree is an NP-hard problem. As such, current tree search algorithms might result in a tree that is the local optima, not the global one. Here, we introduce a paradigm shift for predicting the maximum-likelihood tree, by approximating long-term gains of likelihood rather than …

## A LASSO-based approach to sample sites for phylogenetic tree search

Noa Ecker[1], Dana Azouri[1,2], Ben Bettisworth[3,4], Alexandros Stamatakis[3,4], Yishay Mansour[5], Itay Mayrose[2,*] and Tal Pupko [1,*]

[1]The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 69978, Israel, [2]School of Plant Sciences and Food Security, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 69978, Israel, [3]Computational Molecular Evolution Group, Heidelberg Institute for Theoretical Studies, 69118 Heidelberg, Germany, [4]Institute of Theoretical Informatics, Karlsruhe Institute of Technology, 76128 Karlsruhe, Germany and [5]The Blavatnik School of Computer Science, Raymond & Beverly Sackler Faculty of Exact Sciences, Tel Aviv University, Tel Aviv 69978, Israel

*To whom correspondence should be addressed.

Abstract

## ModelTeller: Model Selection for Optimal Phylogenetic Reconstruction Using Machine Learning

Shiran Abadi [1], Oren Avram,[2] Saharon Rosset,[3] Tal Pupko,[2] and Itay Mayrose[*,1]

[1]School of Plant Sciences and Food security, Tel-Aviv University, Tel-Aviv, Israel
[2]School of Molecular Cell Biology & Biotechnology, Tel-Aviv University, Tel-Aviv, Israel
[3]Department of Statistics and Operations Research, School of Mathematical Sciences, Tel-Aviv University, Tel-Aviv, Israel
*Corresponding author: E-mail: itaymay@tauex.tau.ac.il.
Associate editor: Li Liu

### Abstract

Statistical criteria have long been the standard for selecting the best model for phylogenetic reconstruction and downstream statistical inference. Although model selection is regarded as a fundamental step in phylogenetics, existing methods for this task consume computational resources for long processing time, they are not always feasible, and sometimes depend on preliminary assumptions which do not hold for sequence data. Moreover, although these methods are dedicated to revealing the processes that underlie the sequence data, they do not always produce the most accurate trees. Notably, phylogeny reconstruction consists of two related tasks, topology reconstruction and branch-length estimation. It was previously shown that in many cases the most complex model, GTR+I+G, leads to topologies that are as accurate as using existing model selection criteria, but overestimates branch lengths. Here, we present ModelTeller, a computational methodology for phylogenetic model selection, devised within the machine-learning framework, opti-

## A machine-learning-based alternative to phylogenetic bootstrap

Noa Ecker[1], Dorothée Huchon [2,3], Yishay Mansour [4], Itay Mayrose [5], Tal Pupko [1,*]

[1]The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel
[2]School of Zoology, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel
[3]The Steinhardt Museum of Natural History and National Research Center, Tel Aviv University, Tel Aviv 0997801, Israel
[4]The Blavatnik School of Computer Science, Raymond & Beverly Sackler Faculty of Exact Sciences, Tel Aviv University, Tel Aviv 0997801, Israel
[5]School of Plant Sciences and Food Security, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel

*Corresponding author. The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel. E-mail: talp@tauex.tau.ac.il

### Abstract

**Motivation:** Currently used methods for estimating branch support in phylogenetic analyses often rely on the classic Felsenstein's bootstrap, parametric tests, or their approximations. As these branch support scores are widely used in phylogenetic analyses, having accurate, fast, and interpretable scores is of high importance.
**Results:** Here, we employed a data-driven approach to estimate branch support values with a probabilistic interpretation. To this end, we simulated thousands of realistic phylogenetic trees and the corresponding multiple sequence alignments. Each of the obtained alignments was used to infer the phylogeny using state-of-the-art phylogenetic inference software, which was then compared to the true tree. Using these extensive data, we trained machine-learning algorithms to estimate branch support values for each bipartition within the maximum-likelihood …

# AI phylogeny

## Part 1:

OXFORD

# A machine-learning-based alternative to phylogenetic bootstrap

Noa Ecker[1], Dorothée Huchon [2,3], Yishay Mansour [4], Itay Mayrose [5], Tal Pupko [1,*]

[1]The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel
[2]School of Zoology, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel
[3]The Steinhardt Museum of Natural History and National Research Center, Tel Aviv University, Tel Aviv 0997801, Israel
[4]The Blavatnik School of Computer Science, Raymond & Beverly Sackler Faculty of Exact Sciences, Tel Aviv University, Tel Aviv 0997801, Israel
[5]School of Plant Sciences and Food Security, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 0997801, Israel

*Corresponding author. The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 6997801, Israel. E-mail: talp@tauex.tau.ac.il

## Abstract

**Motivation:** Currently used methods for estimating branch support in phylogenetic analyses often rely on the classic Felsenstein's bootstrap, parametric tests, or their approximations. As these branch support scores are widely used in phylogenetic analyses, having accurate, fast, and interpretable scores is of high importance.

**Results:** Here, we employed a data-driven approach to estimate branch support values with a probabilistic interpretation. To this end, we simulated thousands of realistic phylogenetic trees and the corresponding multiple sequence alignments. Each of the obtained alignments was used to infer the phylogeny using state-of-the-art phylogenetic inference software, which was then compared to the true tree. Using these extensive data, we trained machine-learning algorithms to estimate branch support values for each bipartition within the maximum-likelihood trees obtained by each software. Our results demonstrate that our model provides fast and more accurate probability-based branch support values
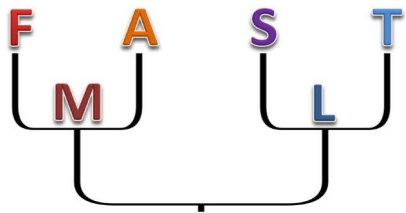
# Phylogenetic-based tools and algorithms

**The GUIDANCE2 Server**
Server for alignment confidence score
HOME     OVERVIEW     GALLERY     SOURCE CODE     CITING & CREDITS     CONTACT US

M1CR0B1AL1Z3R
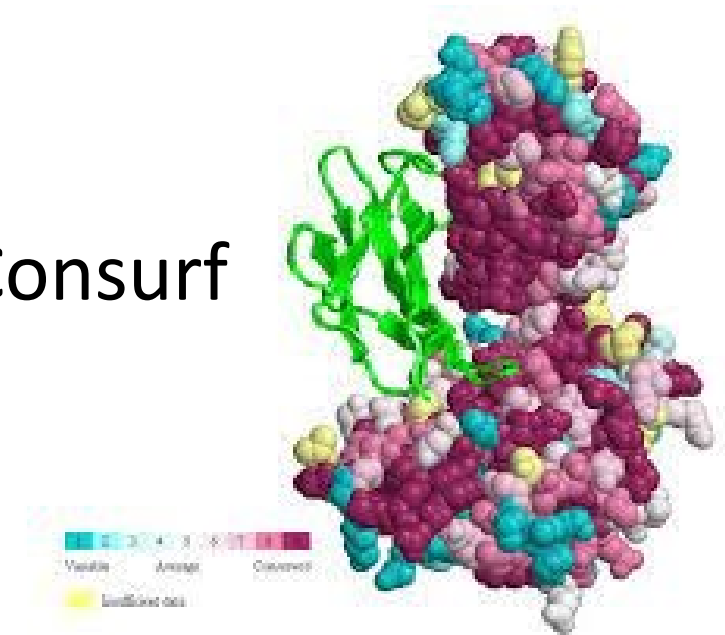A web server for analyzing bacterial genomics data. Easily.

Consurf

**The FASTML Server**
Server for computing Maximum Likelihood
ancestral sequence reconstruction
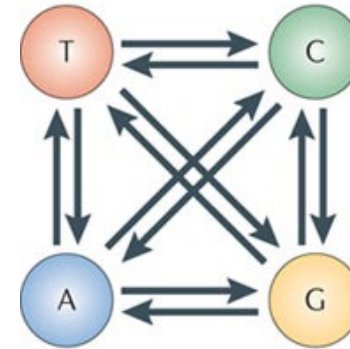HOME     OVERVIEW     GALLERY     SOURCE CODE     CITING & CREDITS

# The tree score = log-likelihood

*T = Tree*



*M = Evolutionary model*



*Tree likelihood*

$P(D|M,T)$: The conditional probability of the data, given the model, *M,* and the tree, *T*

*D = Data = Alignment*

```
Acanthocalyx_albus    tcgaaacttg cccagcagag cgaccagcga acacctccgt
Abelia_spathulata     tcgaaacctg cacagcagaa cgacccgcga acacgttcgt
Abelia_engleriana     tcgaaacctg cacagcagaa cgacccgcga acacgttcgt
Abelia_chinensis      tcgaaacctg cacagcagaa cgacccgcga acacgttcgt
Abelia_parvifolia     --------g cacagcagaa cgacccgcga acacgttcgt
Abelia_Agrandiflora   tcgaaacctg cacagcagaa cgacccgcga acacgttcgt
Abelia_mexicana       -tcgaaactg cacagcagaa cgacccgcga acacgttcgt
Abelia_occidentalis   tcgaaacctg cacagcagaa cgacccgcga acacgttcgt
```

# Finding trees keeps becoming more difficult

**Today:**

- Data: Aligned genomic sequences

- Size:  Up to GB of DNA sites, thousands of species

# The tree space is huge

| Number of Taxa | Number of rooted trees |
|---|---|
| 3 | 3 |
| 4 | 15 |
| 5 | 105 |
| 6 | 945 |
| 7 | 10,395 |
| 8 | 135,135 |
| 9 | 2,027,025 |
| 10 | 34,459,425 |
| 20 | $8.20 \times 10^{21}$ |
| 30 | $4.95 \times 10^{38}$ |
| 40 | $1.01 \times 10^{57}$ |
| 50 | $2.75 \times 10^{76}$ |

- Start the search with a good guess for a starting tree

- Examine all "neighboring" trees by making small modifications to the current tree

- Move to the neighbor with the highest (likelihood) score

- To avoid local maxima, we start from multiple starting points



likelihood

Phylogenetic tree space

Start tree

# Different (unrooted) trees



LL score = -110



LL score = -107



LL score = -117

How confidence we are in the best tree?

# Jackknife

- We create new data sets by sampling randomly half of the characters without replacement.

- We generate 100 pseudo-data sets.

- We do not change the number of sequences, just the number of positions!

# Jackknife

## Original data

```
s1   AAGTAA
s2   CAAAAC
s3   CAGGAA
s4   AAATAC
```

## Pseudo-data 1

```
s1   AGA
s2   AAA
s3   AGA
s4   AAA
```

We removed positions 1,4, and 6

# Jackknife

Pseudo-data 1

```
s1   AGA
s2   AAA
s3   AGA
s4   AAA
```

ML tree search

Best tree for pseudo-data 1

- We repeat the process 100 times and get 100 best trees

# Jackknife percentage


30 times


50 times


20 times

Our confidence in the pink tree in terms of jackknife support is 50%.

# The phylogenetic bootstrap

# Bootstrap

Original data

Pseudo-data 1

```
s1   AAGTAA
s2   CAAAAC
s3   CAGGAA
s4   AGATAC
```

```
s1   AGAGAT
s2   AAAAAA
s3   AAAGAG
s4   AGAAGT
```

Same idea as jackknife but we sample with repetition. In this example, we sample positions 2 twice and position 5 twice and zero for position 1 and 6.

# Bootstrap

- Bootstrap is used more than jackknife in phylogeny, because it has the same data-size as the original data.

# Bootstrap for splits

- Instead of getting the support for each tree, we can compute the support for a given split.

- The support for a given split is the percentage of pseudo-tree in which this split appears.

# Bootstrap for splits

# Bootstrap is very slow

# There's a need for faster estimates

**Table 1.** Several branch support methods implemented in current tree search software.

| Program | Branch-support method | References |
|---|---|---|
| RAxML-NG | Standard Felsenstein's bootstrap | (Kozlov et al., 2019) |
| RAxML-NG | Transfer bootstrap expectation | (Kozlov et al., 2019; Lemoine et al., 2018) |
| IQTREE | Ultrafast bootstrap | (Hoang et al., 2018; Minh et al., 2013) |
| IQTREE | aLRT test | (Soang et al., 2018; Anisimova and Gascuel, 2006) |
| IQTREE | aBayes test | (Hoang et al., 2018; Anisimova et al., 2011) |
| FastTree | SH-like test | (Price et al., 2010) |

# Likelihood based methods



- Likelihood-based methods will estimate the support using differences in log-likelihood

# How to evaluate different branch-support methods?

"TRUE TREE"

ML tree with BP
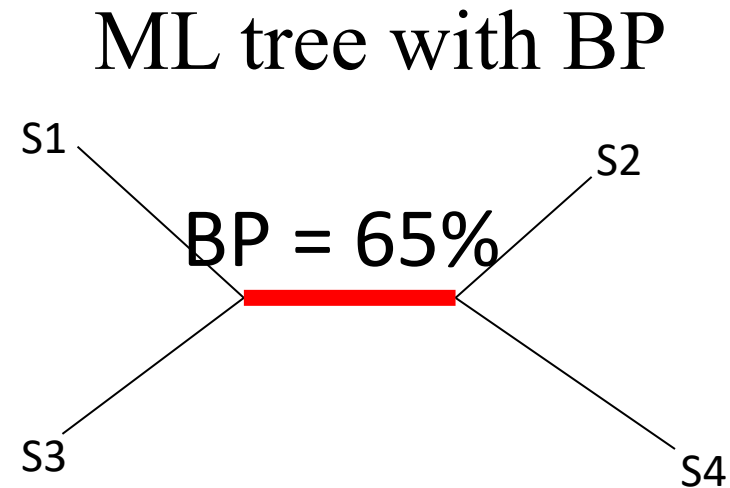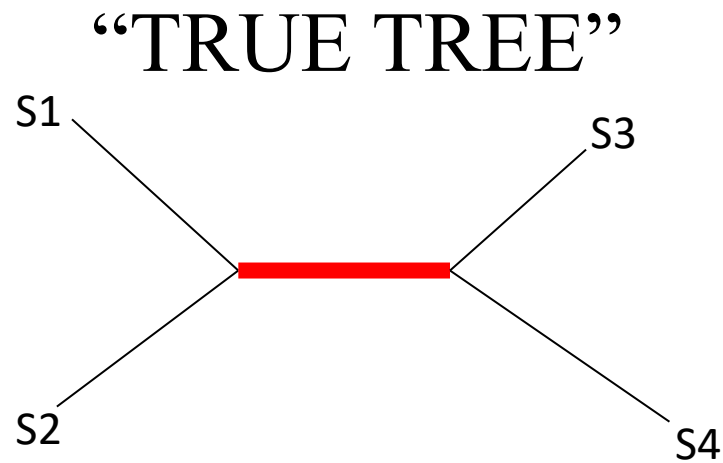
S1
S3
S2
S4

S1
BP = 65%
S3
S2
S4

- In this case we have a true positive inference (TP)
- POSITIVE = the estimate is that the split exists
- TRUE = the estimate is correct

# How to evaluate different branch-support methods?

"TRUE TREE"

ML tree with BP

BP = 35%

- In this case we have a false negative inference (FN)
- NEGATIVE = the estimate is that the split does not exist
- FALSE = the estimate is wrong

# How to evaluate different branch-support methods?

## "TRUE TREE"

S1

S3

S2

S4

## ML tree with BP

S1

S2

BP = 65%

S3

S4

- In this case we have a false positive inference (FP)

- POSITIVE = the estimate is that the split exists

- FALSE = the estimate is wrong

# We can compute confusion matrices and AUC scores

# ML (machine learning) for branch support values



**Noa Ecker**

# ML (machine learning) for branch support values



Prof. Yishay Mansour

Prof. Itay Mayrose

Prof. Dorothee Huchon

# Intuition

- I would trust the red branch more than the blue one

# Features (out of 39)

- Branch length at the partition site
- Branch length divided by the mean branch length across the tree
- Branch length divided by the mean branch length among the four neighboring branches
- Number of MSA columns
- Number of unique MSA columns
- Percentage of constant sites
- The LL of NNI neighbors around the branch

- The count and proportion of taxa on the smaller or equal side of the bipartition



Proportion of taxa on the smaller side of the bipartition = 0.25
Number of taxa on the smaller side of the bipartition = 2

- 6,000 simulated MSAs with 100 to 10,000 sites and between 30 to 1,000 taxa.
- Each MSA was simulated along a different tree topology using Alisim (Ly-Trong et al., 2022), based on the script provided in the Github repository of RAxML-grove (Höhler et al., 2022).
- Each MSA was simulated using the DNA model associated with that tree in RAxML-Grove.
- Train set: 70% of the data; test set: the remaining 30%

# Results, performance

- RAxML BP = 0.946   **Machine-learning = 0.968**

- RAxML transfer BP = 0.907

# Results, running time

- **21 times faster** without optimizing the feature extraction algorithms:

The computation time of RAxML-NG standard bootstrap exhibited a median running time of 138 min on a single CPU. On the same data, our machine-learning model had a median running time of 6.5 min.
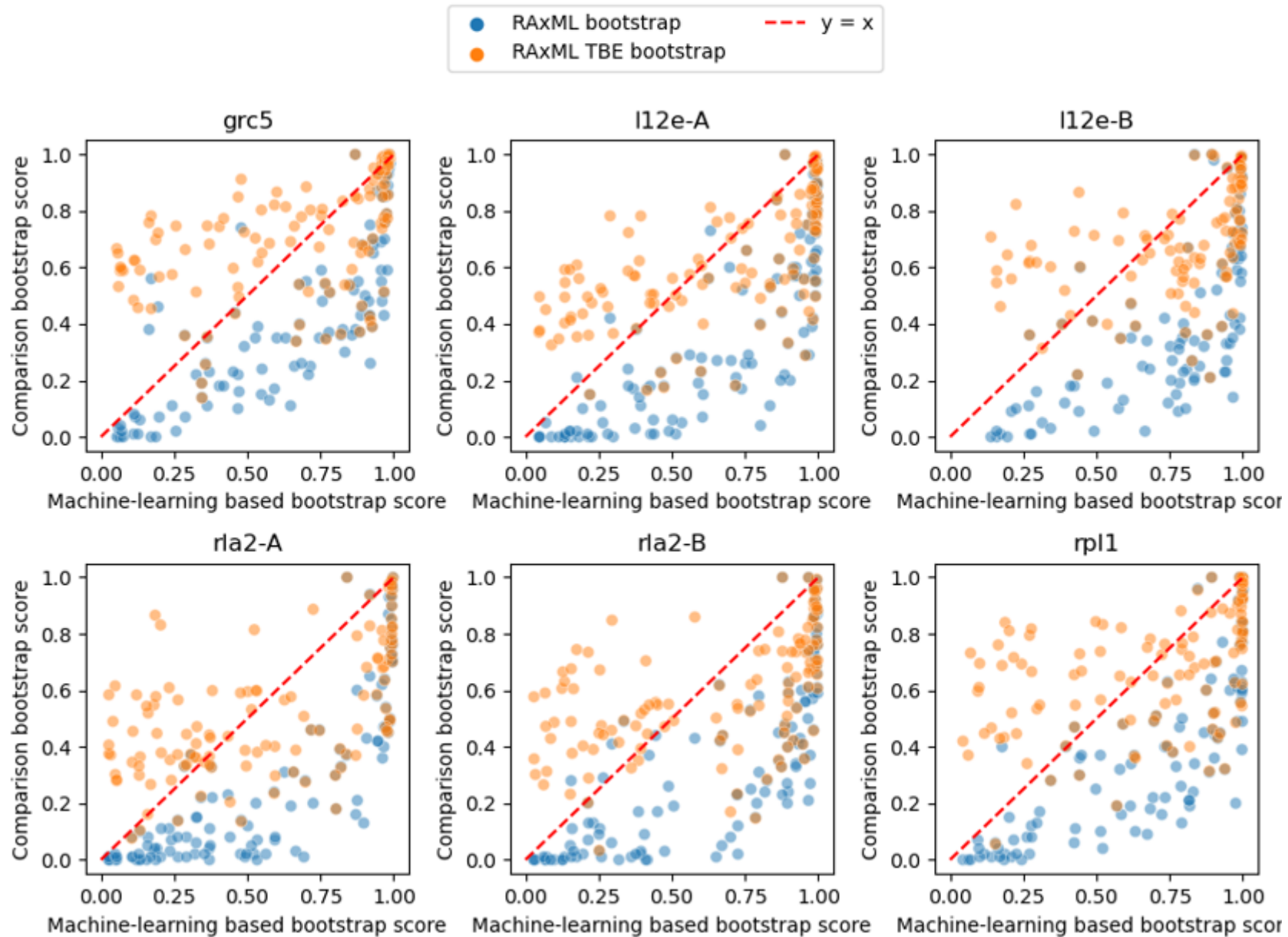
# Results, calibration



**Branch support score**

Expected Calibration Error (ECE) of machine-learning method
(IQTREE) = 0.002

ECE for ultrafast bootstrap (IQTREE) = 0.043

# Results, empirical data

# PART 2:

OXFORD

## Phylogenetics

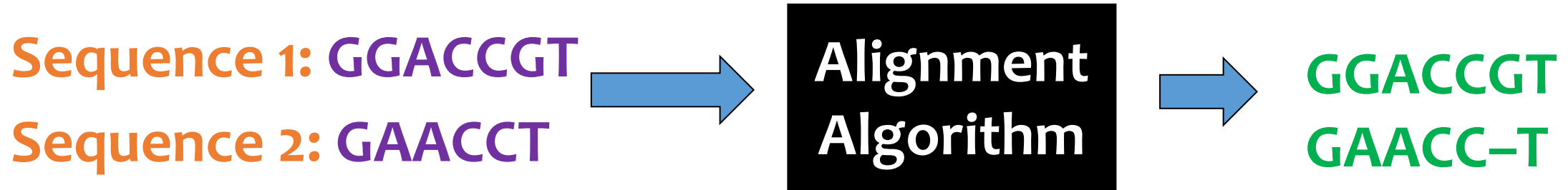# BetaAlign: a deep learning approach for multiple sequence alignment

Edo Dotan [1,2], Elya Wygoda[1], Noa Ecker[1], Michael Alburquerque[1], Oren Avram [3], Yonatan Belinkov[2,*], Tal Pupko [1,*]

[1]The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 69978, Israel
[2]The Henry and Marilyn Taub Faculty of Computer Science, Technion—Israel Institute of Technology, Haifa 3200003, Israel
[3]The Department of Computer Science, University of California Los Angeles, Los Angeles, CA 90095, United States
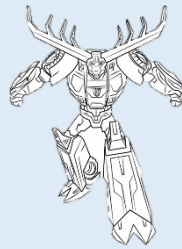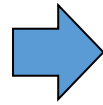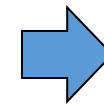
# NLP-based sequence alignment

**Sequence 1: GGACCGT**
**Sequence 2: GAACCT**

→

**Alignment Algorithm**

→

**GGACCGT**
**GAACC–T**

# The transformer

**L'ascension des robots**

Transformer

**Rise of the robots**

AAG
ACCG

Transformer

A A – G
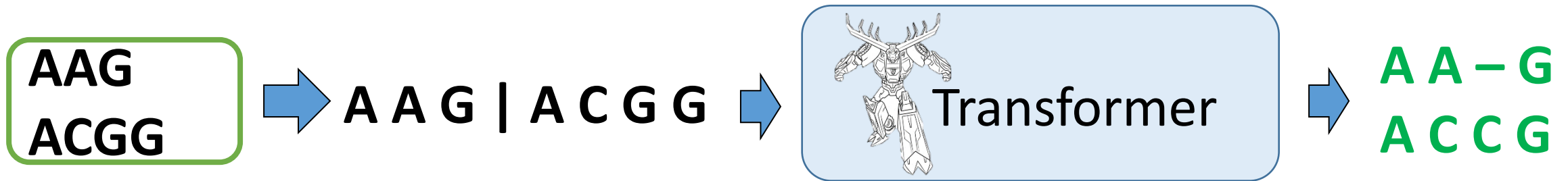A C C G

# Encoding: the "concat" language

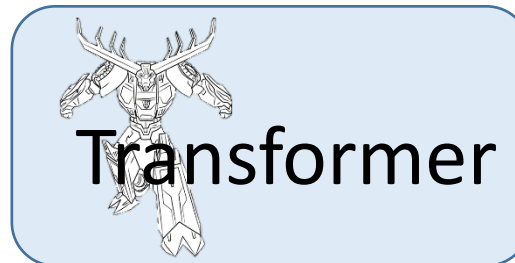- Each letter in the "concat" language is a word, and the language has 5 different words (5 tokens)

**AAG**
**ACGG**  ➡  **A A G | A C G G**  ➡  Transformer  ➡  A A – G
                                                      A C C G

# The output language

| Upper letter | Lower letter | Symbol |
|:---:|:---:|:---:|
| A | A | A |
| A | C | B |
| A | G | C |
| A | T | D |
| A | - | E |
| … | … | … |
| - | T | X |

**AAG**
**ACGG**
➡
**A A G | A C G G**
➡
Transformer
➡
**ABVM**
➡
**A A − G**
**A C C G**

# Performance

Edo Dotan[1‡], Yonatan Belinkov[2¹*], Oren Avram[3], Elya Wygoda[1], Noa Ecker[1], Michael Alburquerque[1], Omri Keren[1], Gil Loewenthal[1], and Tal Pupko[1]

Phylogenetics

**BetaAlign: a deep learning approach for multiple sequence alignment**

Edo Dotan [1,2], Elya Wygoda[1], Noa Ecker[1], Michael Alburquerque[1], Oren Avram [3], Yonatan Belinkov[2,*], Tal Pupko [1,*]

# Performance



Protein sequences

## PART 3 (work in progress):

Better than sum-of-pairs: a machine-learning-based score to evaluate multiple sequence alignments

Nimrod Serok[1,*], Ksenia Polonsky[1,*], Haim Ashkenazy[2], Dorothée Huchon[3,4], Itay Mayrose[5], Jeffrey Thorne[6,7], Tal Pupko[1†]

[1] The Shmunis School of Biomedicine and Cancer Research, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv 69978, Israel.

[2] Department of Molecular Biology, Max Planck Institute for Biology Tübingen, Tübingen, Germany.

[3] School of Zoology, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel

# An ML-based MSA objective function

- Inference of MSAs is a very difficult problem

  - The implicitly assumed indel evolutionary models are oversimplified

  - Which objective function should be optimized (the likelihood are very difficult to compute)

  - Optimizing the objective function is difficult

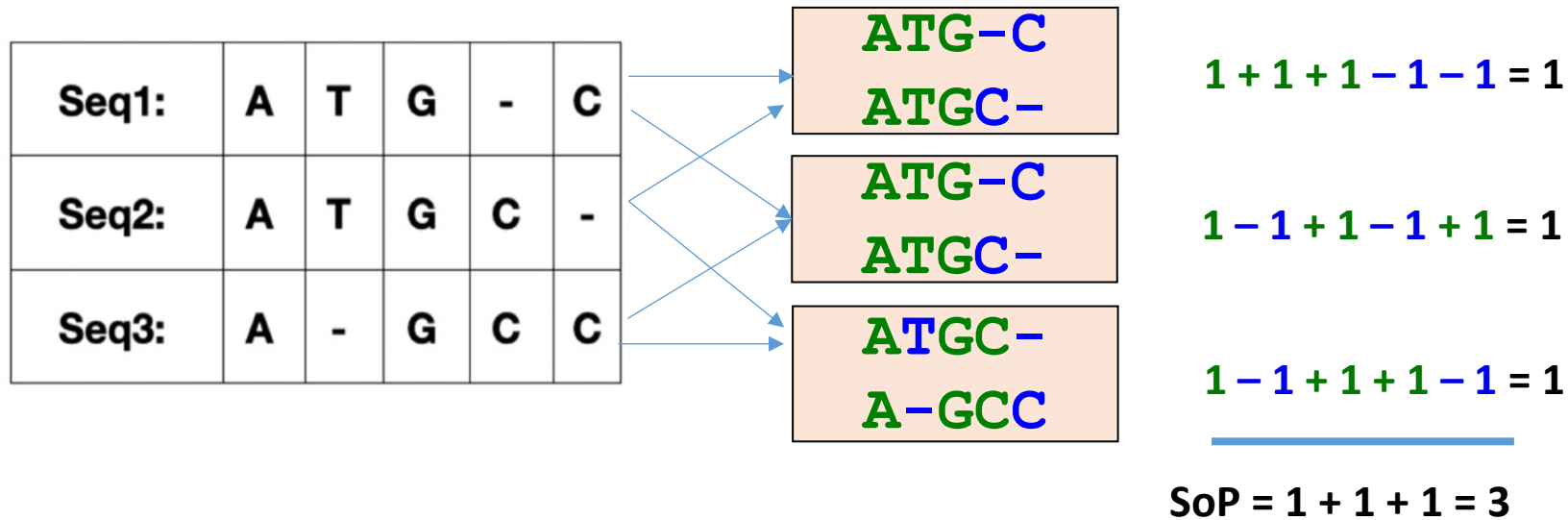  - The MSA depends on the tree and vice versa

- Inference of MSAs is a very difficult problem

  - The implicitly assumed indel evolutionary models are oversimplified

  - **Which objective function should be optimized (the likelihood are very difficult to compute)**

  - Optimizing the objective function is difficult

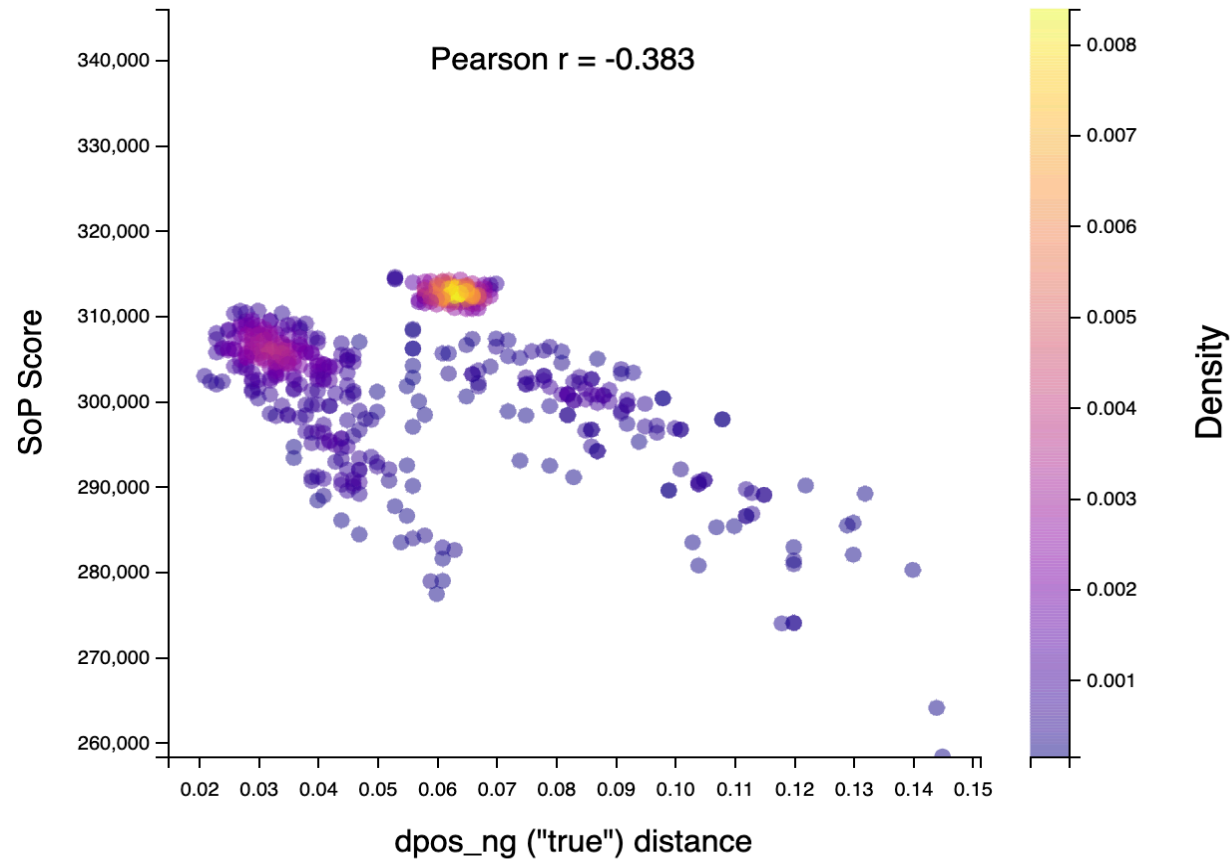  - The MSA depends on the tree and vice versa

# Sum of pairs

- The sum-of-pairs (SoP) score is widely used for MSA scoring

- Higher SoP is supposed to be an indicator of a better MSA

**Naïve scoring:**
Mismatch: -1
Indel (gap): -1
Perfect match: +1

| | | | | | |
|---|---|---|---|---|---|
| Seq1: | A | T | G | - | C |
| Seq2: | A | T | G | C | - |
| Seq3: | A | - | G | C | C |

ATG-C
ATGC-

ATG-C
ATGC-

ATGC-
A-GCC

$1 + 1 + 1 - 1 - 1 = 1$

$1 - 1 + 1 - 1 + 1 = 1$

$1 - 1 + 1 + 1 - 1 = 1$

SoP = 1 + 1 + 1 = 3

# Sum of pairs does not corelate well with accuracy

- 500 alternative MSAs of a single dataset

# Proposed solution

- Employing AI to develop novel scoring functions for MSAs that are better than the sum-of-pairs and to use our developed score to discriminate among MSAs

## Features
### 70+
We tried to use many features related to MSA and their corresponding MLE trees; 26 features were used in the final version of the model
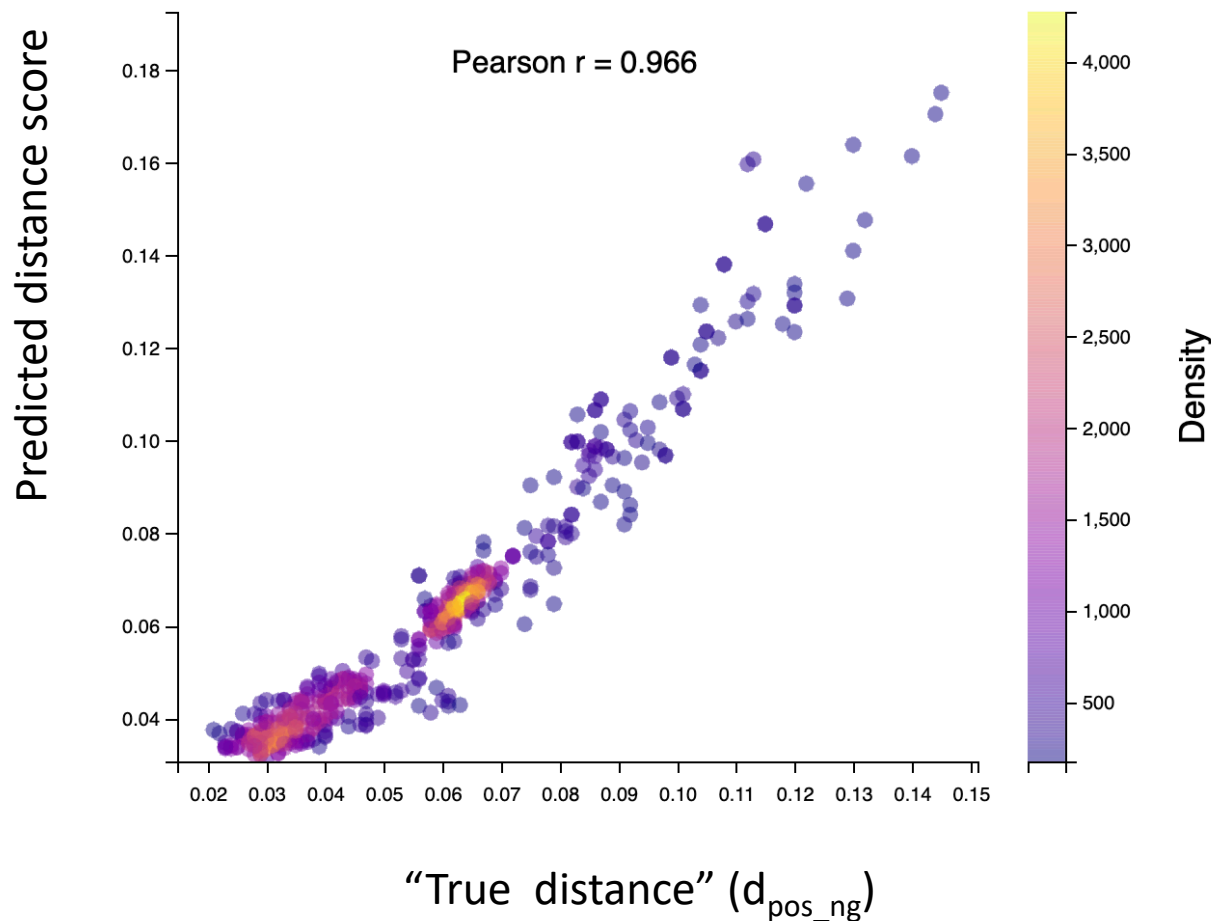
## Alternative MSAs
### >600K
Alternative MSAs were created using four different aligning programs, GUIDANCE and by refinement

## Deep Learning Network
Many model architectures, hyperparameters, scaling, and batching techniques were tested
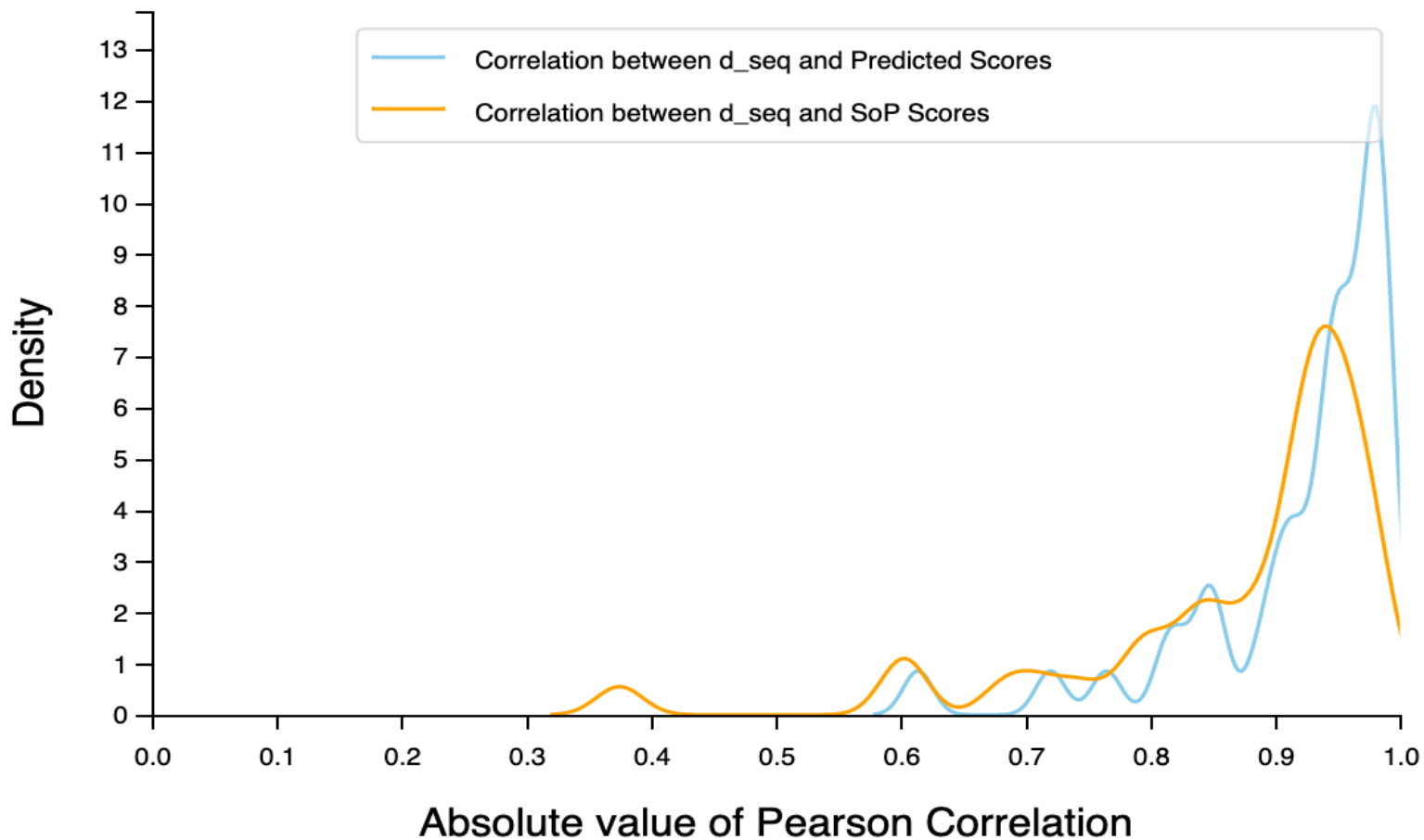
# Our score well corelates with accuracy

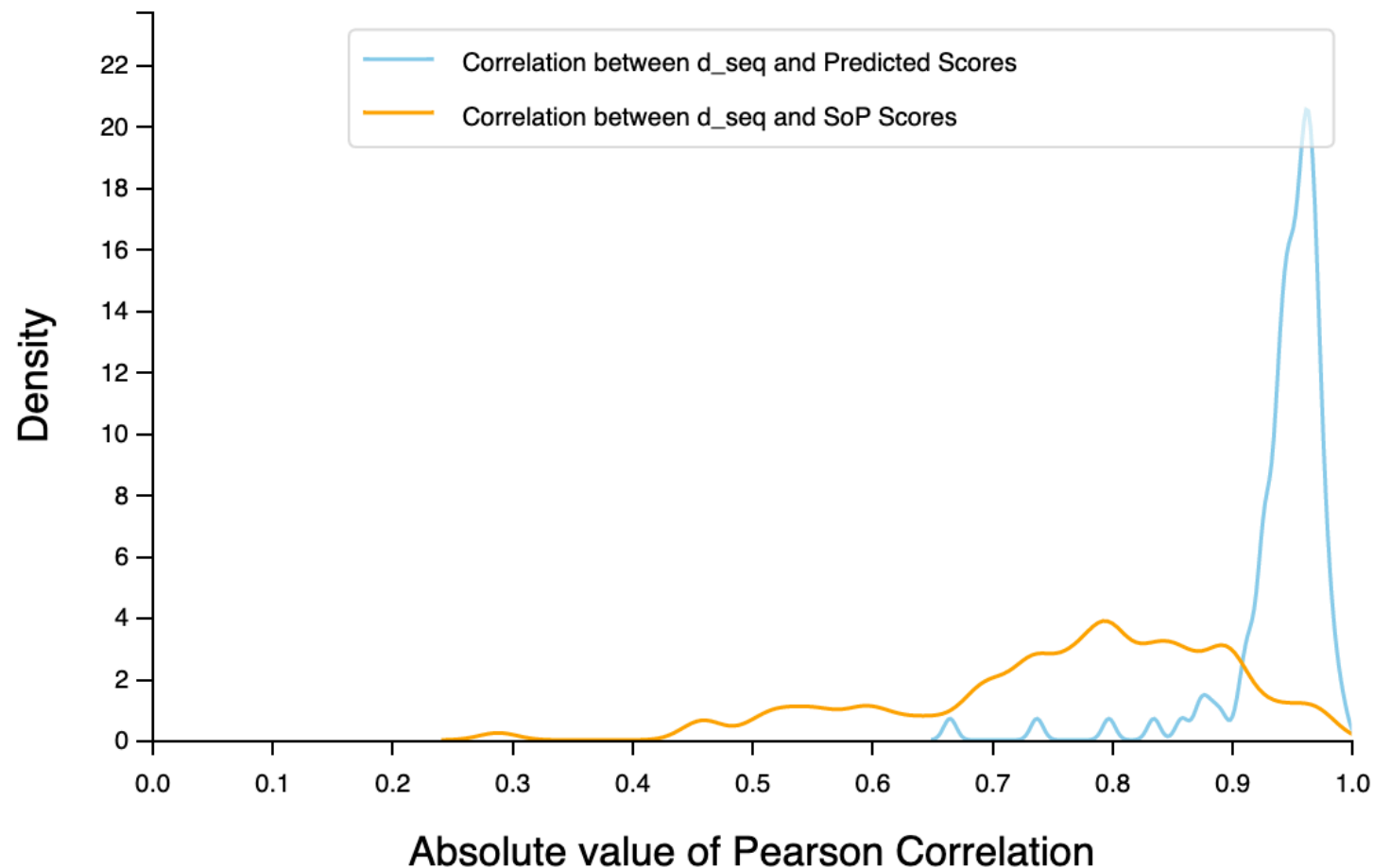- 500 alternative MSAs of a single dataset



Pearson r = 0.966

"True distance" ($d_{pos\_ng}$)

# Our score well corelates with accuracy
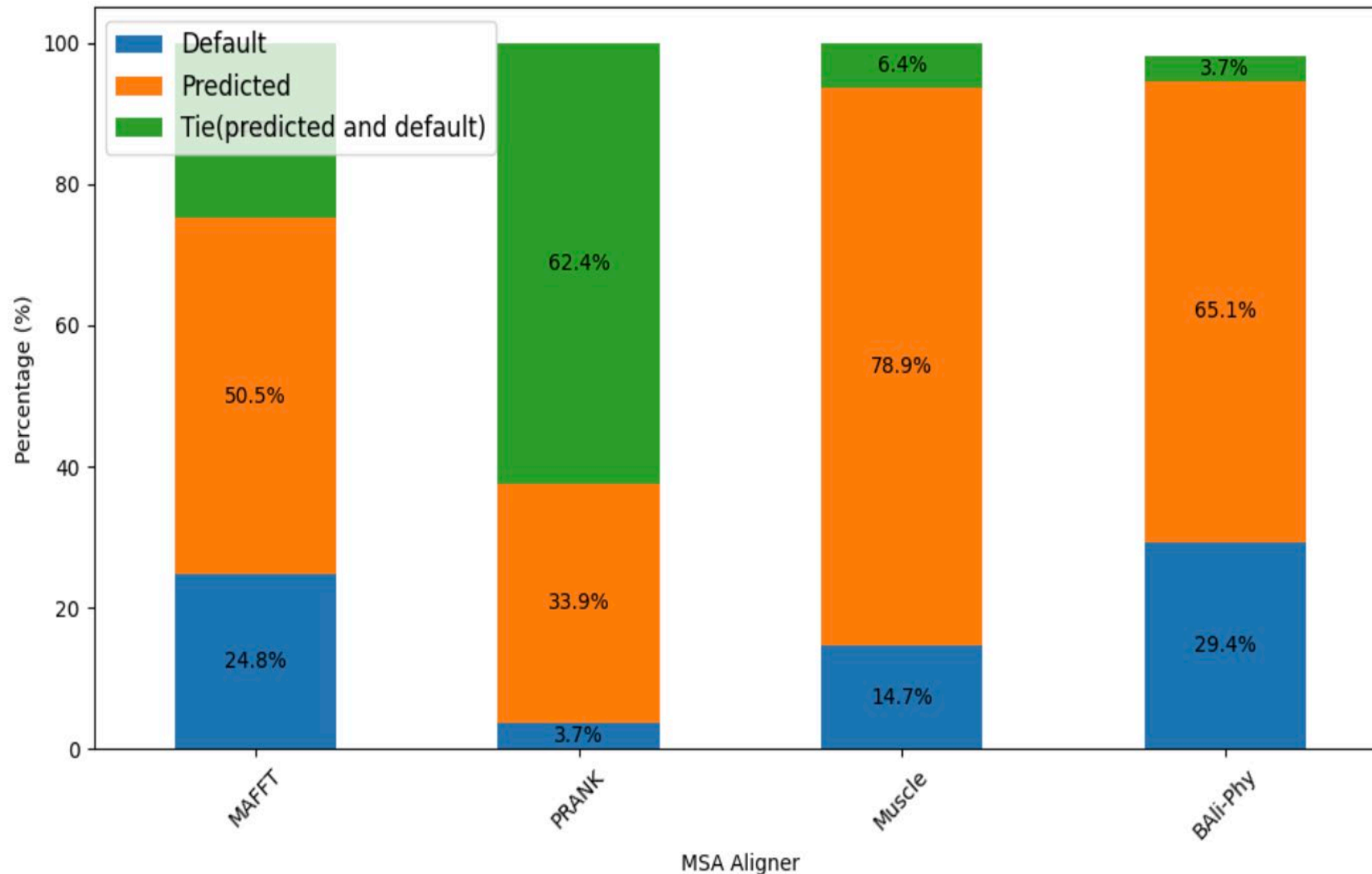
- Simulated datasets

# Our score well corelates with accuracy
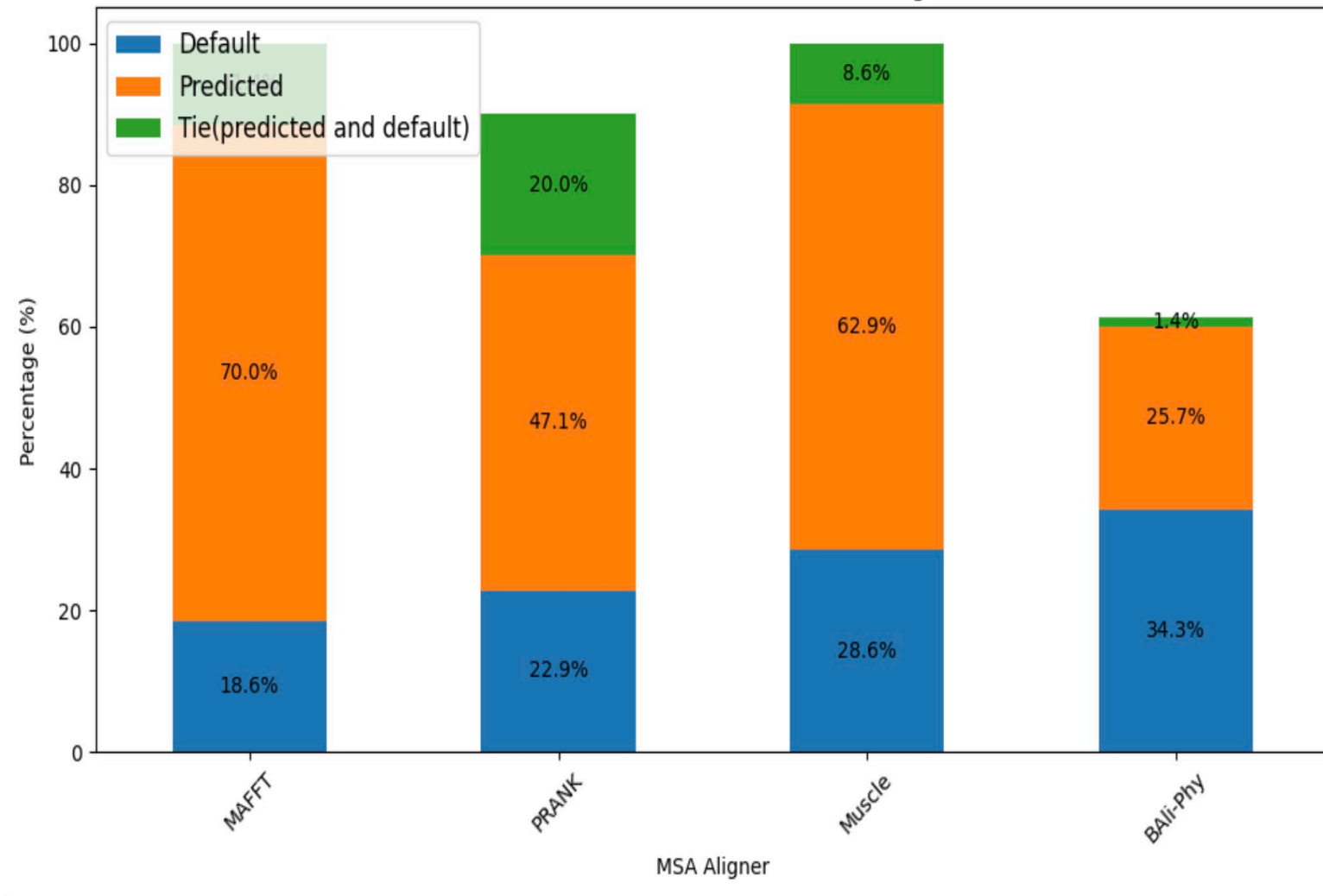
- empirical datasets

# Our score well corelates with accuracy

- Pick me game (simulations)

# Our score well corelates with accuracy

- Pick me game (empirical)

# Joint work



Prof. Yishay Mansour

Prof. Itay Mayrose

Dr. Yonathan Belinkov

Prof. Dorothee Huchon

Prof. Jeffrey Thorne

Noa Ecker

Edo Dotan

Dr. Nimrod Serok

Dr. Haim Ashkenazy

Ksenia Polonsky

# Acknowledgements

Faculty: Itay Mayrose, Dorothee Huchon, Yishay Mansur

Lab members: Naama Wagner, Nimrod Serok, Noa Ecker, Elya Wygoda, Edo Dotan, Ksenia Polonsky, Yair Shimony, Naiel Jabareen, Noam Bracha, and Boris Trayvas