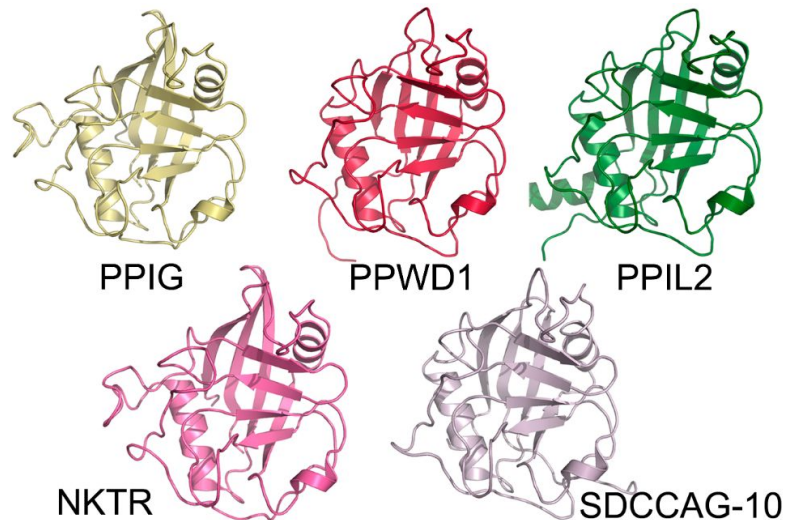# Protein Family Classification Using Profile HMMs

By: Levi Cavagnetto, Daniel Chang, and Garrett Abou-Zeid
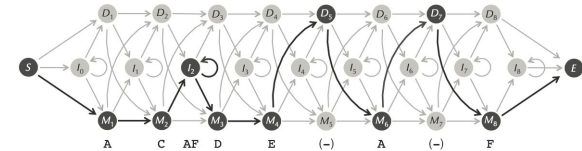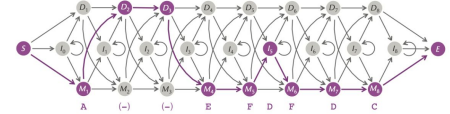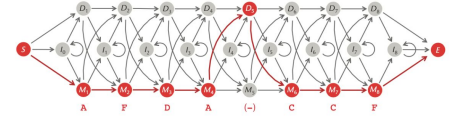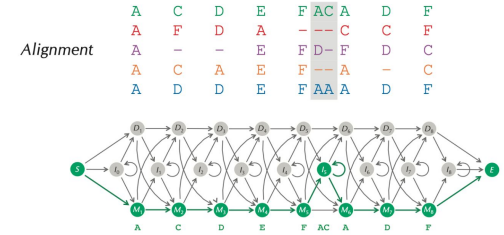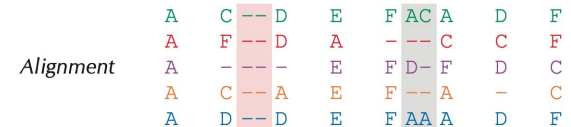
# **Motivation**

- Why is protein family classification important?
    - If you know the family of a novel protein, you know its function!

- Why are profile HMMs useful for protein family classification?
    - A sequence that is a "distance cousin" of a family may **lack strong pairwise similarities** with **any singular sequence** within the family, but may have **subtle/weak pairwise similarities** with **many sequences** in the family
        - Thus, for diverse protein families, pairwise similarity is not a sufficient criteria for determining if a novel sequence is a family member

- Profile HMMs are more efficient to use than constructing an entirely new MSA whenever testing a novel sequence for family membership

PPIG          PPWD1          PPIL2

NKTR                              SDCCAG-10

# Profile HMM

- Profile HMMs generate a probabilistic representation of a **given multiple sequence alignment**

- Profile HMMs are **"trained"** by using each aligned sequence and computing their path through the HMM.
  - After all paths have been computed, transmission and emission probabilities are estimated

- Profile HMMs can then estimate **("testing")** how likely a novel sequence belongs to the family described by the MSA used for training
  - Using the viterbi algorithm, the most likely path (and subsequent log likelihood) is estimated

# Profile HMM
# (our implementation)



"training"

"testing"

```python
sequences = np.array([
    ["A", "C", "D", "E", "F", "A", "C", "A", "D", "F",],
    ["A", "F", "D", "A", "-", "-", "-", "C", "C", "F",],
    ["A", "-", "-", "E", "F", "D", "-", "F", "D", "C",],
    ["A", "C", "A", "E", "F", "-", "-", "A", "-", "C",],
    ["A", "D", "D", "E", "F", "A", "A", "A", "D", "F",],
])

THETA = 0.35
match_state = (sequences == "-").mean(axis=0) < THETA

transition_counts = defaultdict(int)
emission_counts = defaultdict(int)

for seq in sequences:
    last_state = "start"
    for char in seq:

        emission_counts[(new_state, char)] += 1
        new_state = get_new_state(last_state, char)
        transition_counts[(last_state, new_state)] += 1
        ...
```

# Profile HMM
# (our implementation)

```python
def score(query_sequence):
    """
    Top down viterbi algorithm
    """

    def score_helper(seq, state):
        # base case
        if state is start_state:
            return 0

        # recursive case
        if state is match_state:
            return emission_prob[(state, seq[-1])] + max(
                score_helper(seq[:-1], prev_match_state(state)) + transition_prob[(prev_match_state(state), state)],
                score_helper(seq[:-1], prev_insert_state(state)) + transition_prob[(prev_insert_state(state), state)],
                score_helper(seq[:-1], prev_delete_state(state)) + transition_prob[(prev_delete_state, state)],
            )
        elif state is insert_state:
            ...
        elif state is delete_state:
            ...
    return score_helper(query_sequence, "end_state")
```



"training"

"testing"

Alignment

Text

# Methodology

- We focused on building profile HMMs for two genes: **alcohol dehydrogenase (ADH)** and **acetaldehyde dehydrogenase (ACDH)**.

- ~300 Fasta files of the two genes were downloaded from Uniprot
  - To avoid outliers, sequences were limited to those of length 200-400
  - We created a hold out set for testing purposes

- SeaView MSA tool was used to generate MSAs
  - Obviously mismatching sequences were discarded

- Created a profile HMM implementation using python, and used an online implementation of profile HMMs for comparison

- Generated log likelihoods of the trained HMMs when run on the holdout sequences



SeaView MSA

# Results: our implementation

- Using our profile HMM implementation, we trained a profile HMM for both genes

- Log likelihood scores of the test set were generated using the two trained profile HMMs

- The pHMM was able to stratify the two genes fairly well

  - Alcohol dehydrogenase seems to have more heterogeneity

- Performance was highly dependent on smoothing factor (LaPlace) and the maximum proportion of insertions per column for match states



Log Likelihood of alcohol dehydrogenase Profile HMM on holdout genes



Log Likelihood of acetaldehyde dehydrogenase Profile HMM on holdout genes

# Results: online implementation

- A profile HMM was trained using **acetaldehyde** dehydrogenase (ACDH) MSAs
- Log odds scores of the hold out genes are displayed



ADH Test



ACDH Test

# Results: online implementation

- A profile HMM was trained using **alcohol** dehydrogenase (ADH) MSAs
- Log odds scores of the hold out genes are displayed



ADH Test



ACDH Test

# **Questions?**

Our implementation: https://github.com/danielchang2002/5481_final