

Function annotated from sequence alignment (FAFSA)

ChemE 545/546: UW DIRECT Courses for Data Scientists

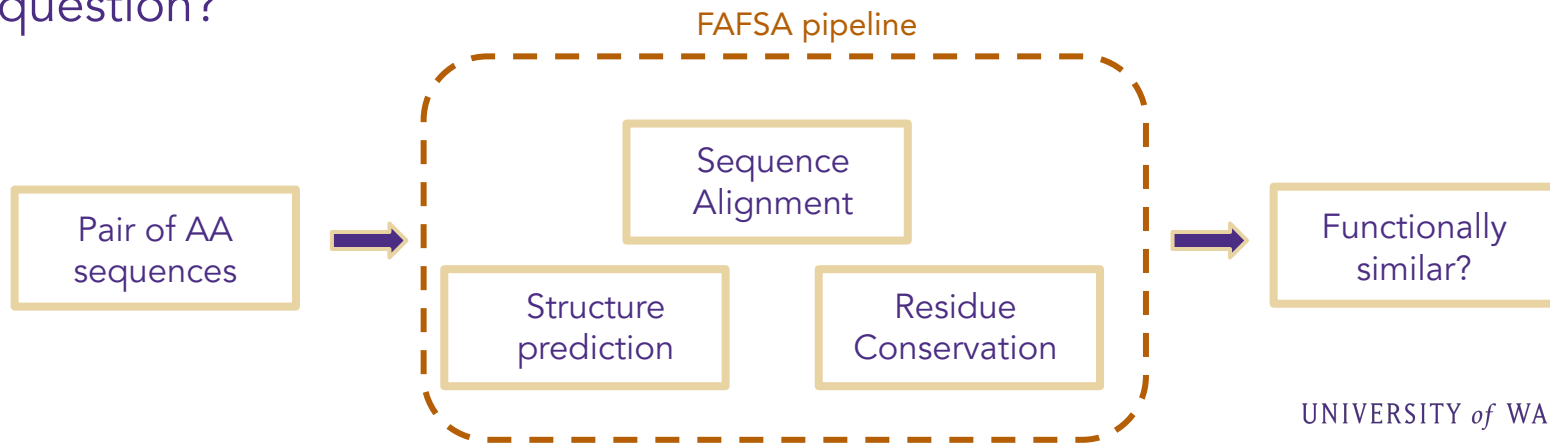
Group members: Humood, Logan, Ryan, Amin, Chau
15/03/23

Why is this useful?

Protein pair validation is problematic in the current state:

- It's time consuming and resource intensive
- Proteins are related through many non-trivial parameters
- Many software applications relate proteins through some of these parameters

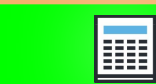
What if one tool could combine these resources to answer a simple question?



Sample preparation

Model development

Deployment



Data retrieval



Data format &
Feature selection

X data



Pfam - Family
identification

Y data



Model
development

Y data



Protein Data Bank
- Structural info

On-going and finish
in Spring

Sample preparation

Model development

Deployment

Winter quarter milestones

Data retrieval

Data format &
Feature selection

X data

Pfam - Family
identification

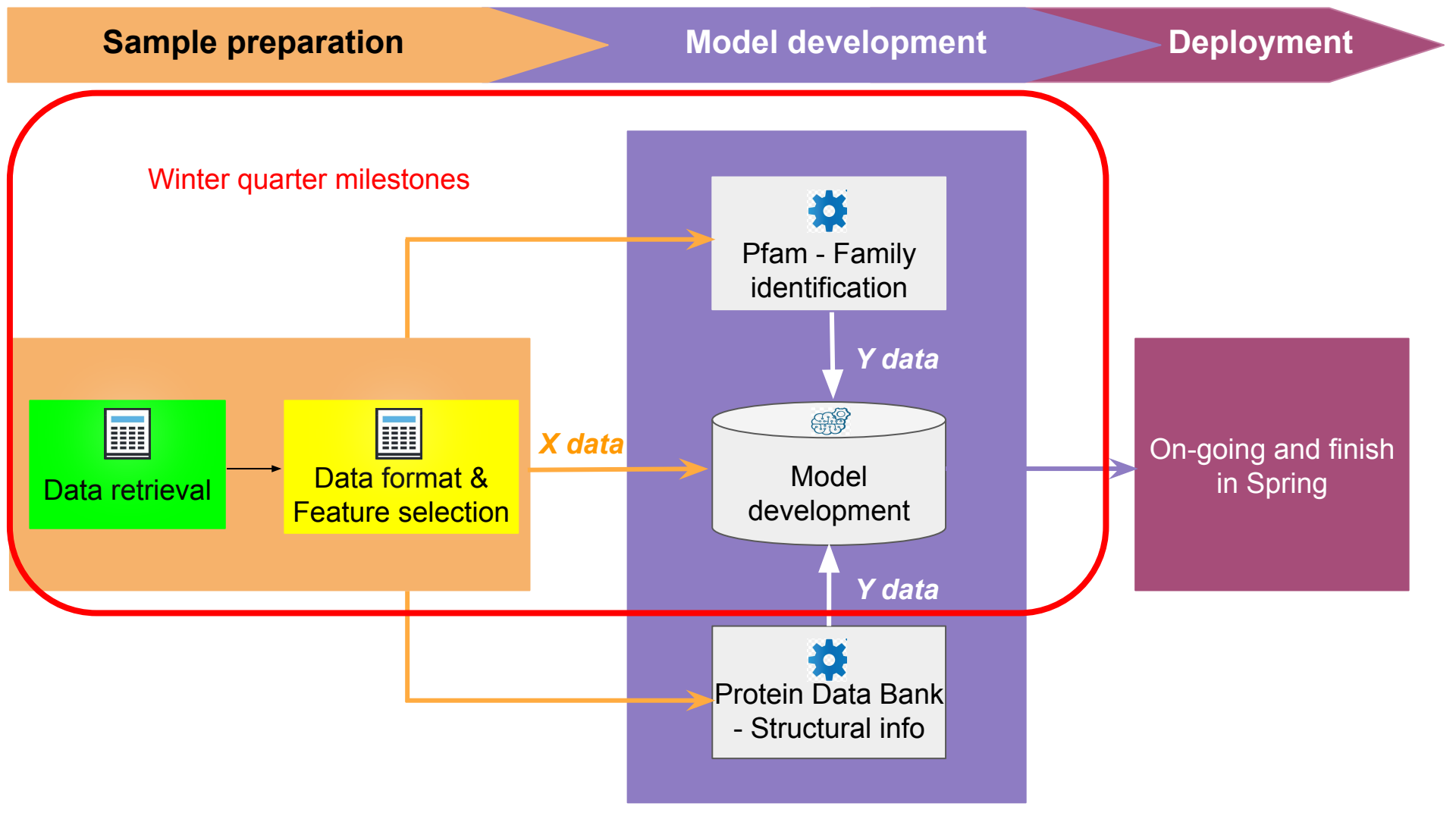
Y data

Model
development

Y data

Protein Data Bank
- Structural info

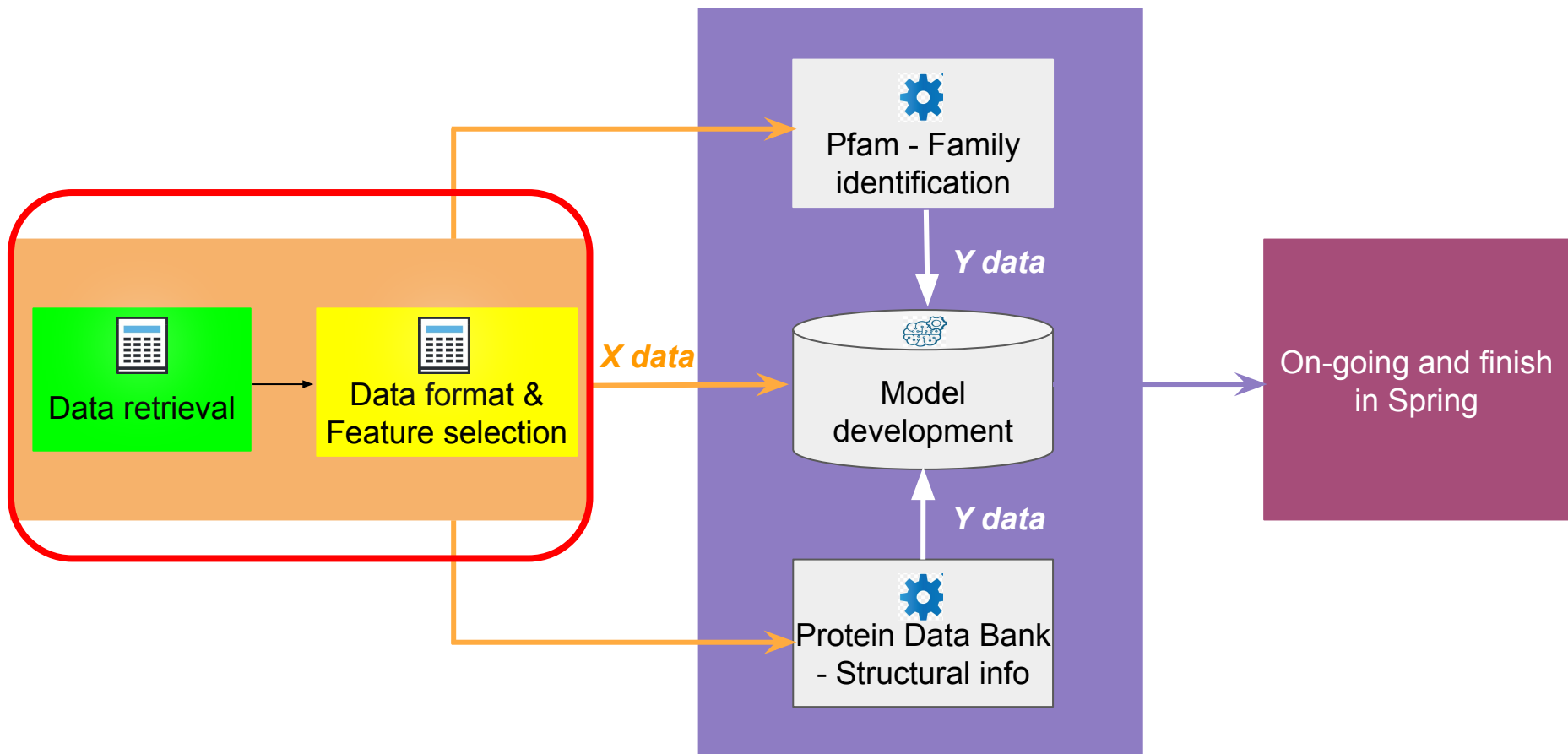
On-going and finish
in Spring



Sample preparation

Model development

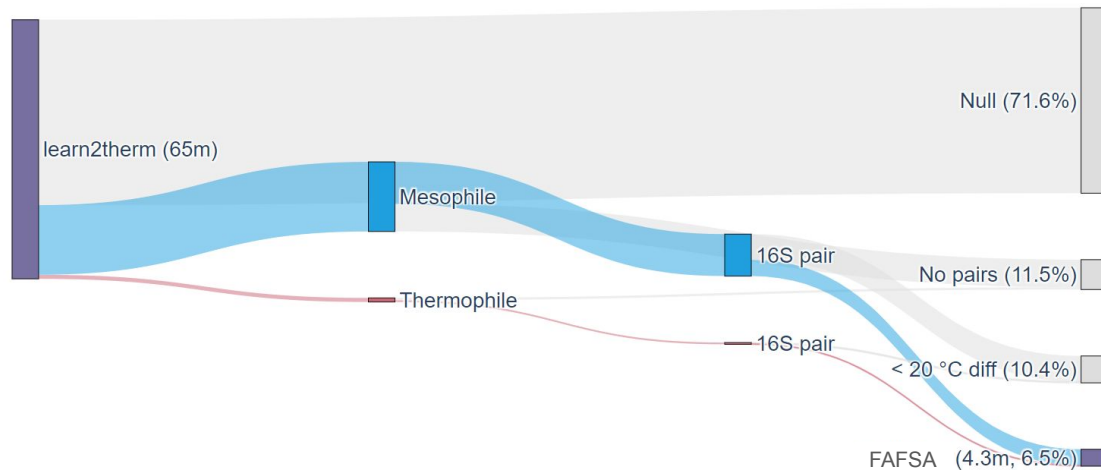
Deployment



Data retrieval from DuckDB

-Protein pairs derived from meso/thermophilic organisms with high 16S alignment.

Protein Representation



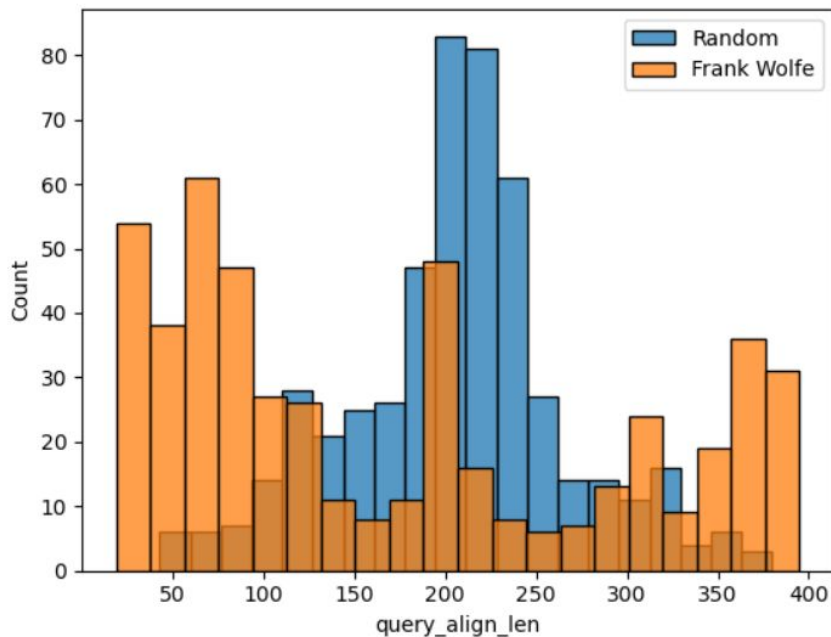
FAFSA

- > 20°C growth temp difference
- 4.3 million proteins
- 54 million pairs
- 1.8k species

Maximizing sample information

- Many features have large tailed distributions.
- D-optimal design (Frank-Wolfe) creates a training set with more information about fringe data points.
- This is computationally expensive because there are multiple features.

Sampling Strategies



Sample preparation

Model development

Deployment

Data retrieval

Data format &
Feature selection

X data

Pfam - Family
identification

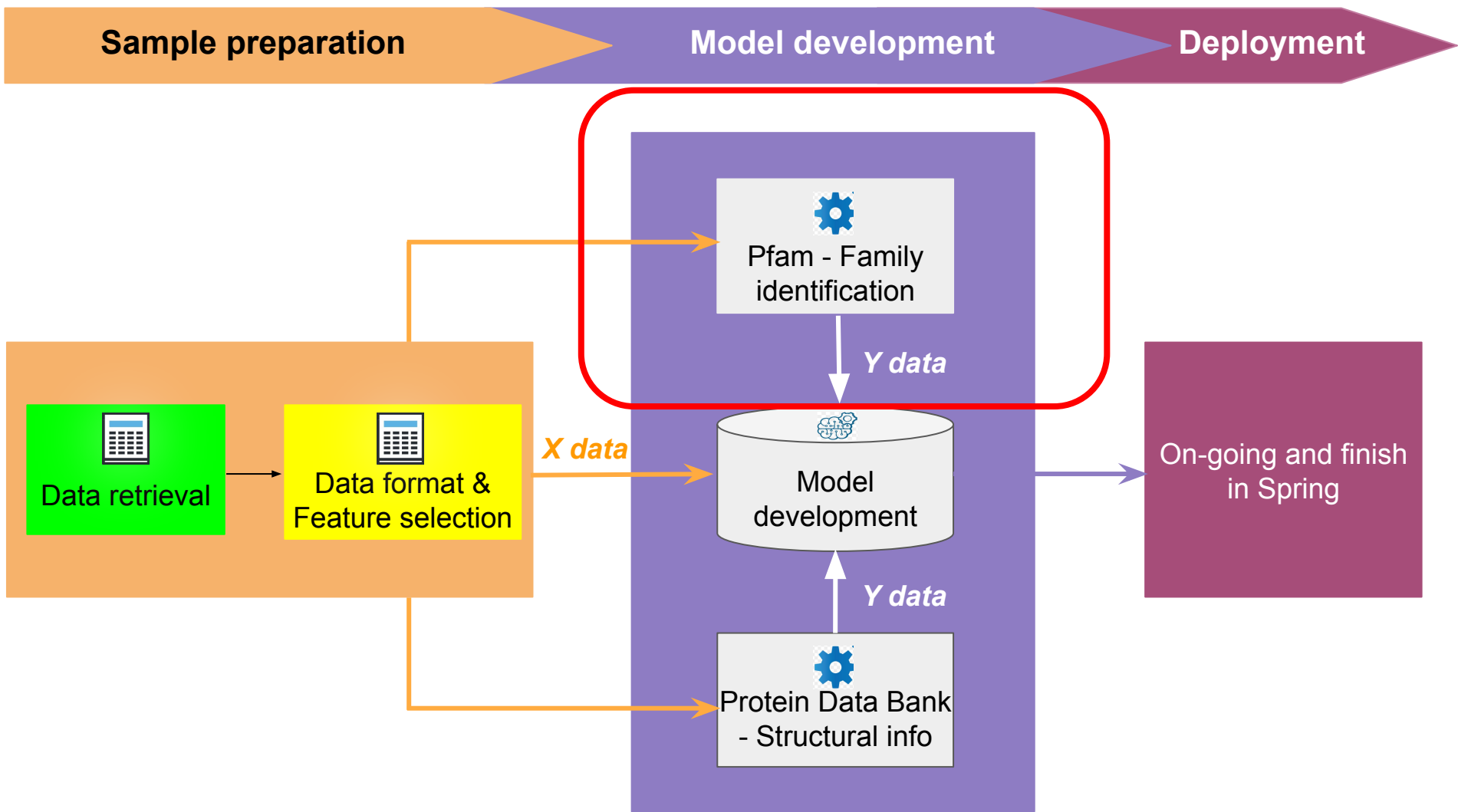
Y data

Model
development

Y data

Protein Data Bank
- Structural info

On-going and finish
in Spring



Finding functional information

Goal:

- 1) Use HMMER and pfam to compute family information for protein pairs
- 2) Parse and filter the output generated to make dataframe if a pair is functional or not

Upstream: Sampling component

Downstream: ML function prediction

MGPSENDPNLFVALYDF ← → MGSFLVRESESSPGQRSISLR

?

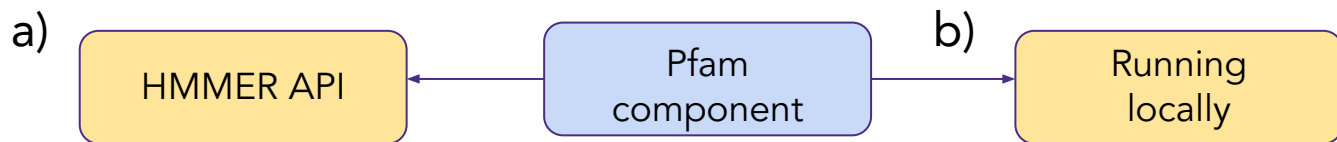


Analysis

- 1) Yes, they're functional pairs.
Appends domain HMMER information to dataframe
- 2) No, they're not functional pairs.
Appends a false boolean to functional pair

Computing functional information from protein pairs

Given two sequences, are we able to see if they're functional or not?

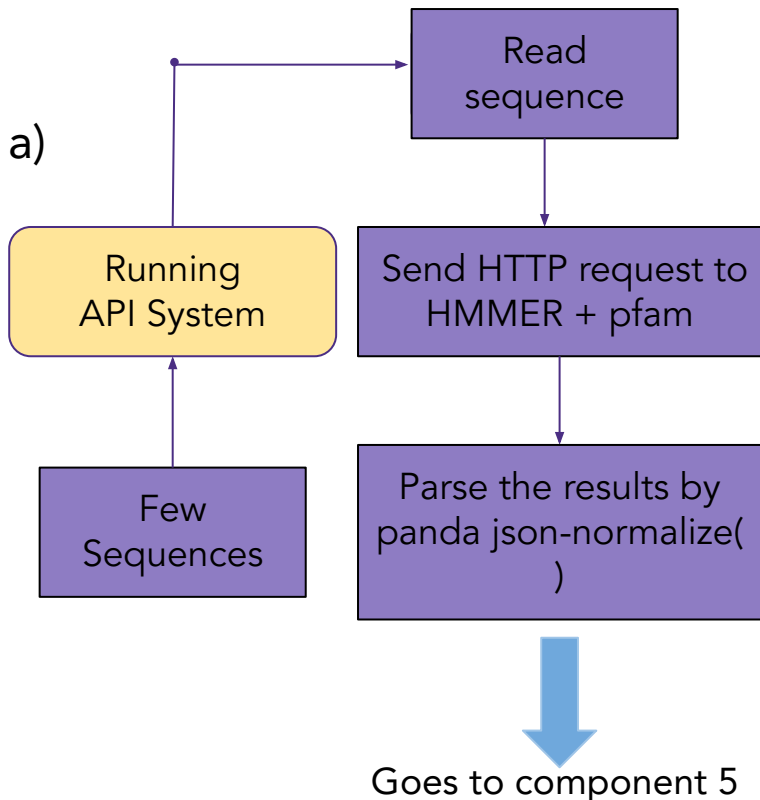


		Run with API	Run locally
1	User has a few sequences	✓	✗
2	Users want few dependencies	✓	✗
3	User doesn't want server variability	✗	✓

This choice is dependent on the user!

Computing functional information from protein pairs

a)



```
def hmmerscanner(df: pd.DataFrame, k: int):  
  
    if k > 300:  
        print("Use local function for the number of sequences more than 300.")  
        return pd.DataFrame()  
  
    # Create an empty DataFrame to store the results.  
    results_df = pd.DataFrame()  
  
    # Loop through the sequences to check them.  
    for i in range(k):  
        # This is for meso protein sequences; we can change that in the future according to our request.  
        sequence = df['m_protein_seq'][i]  
  
        # Send an HTTP request to the HMMER API to get information for the current sequence.  
        url = 'https://www.ebi.ac.uk/Tools/hmmer/search/hmmscan'  
        headers = {'Content-Type': 'application/x-www-form-urlencoded',  
                  'Accept': 'application/json'}  
        data = {'hmddb': 'pfam', 'seq': f">seq\n{sequence}"}  
        data = urllib.parse.urlencode(data).encode('ascii')  
        response = requests.post(url, headers=headers,  
                                data=data, allow_redirects=False)  
        redirect_url = response.headers.get('Location')  
  
        if redirect_url is None:  
            # If the server doesn't work, show this error.  
            print("Error: No redirect URL found in response.")  
        elif redirect_url == 'late':  
            # Raises an exception if the status is pending for too long.  
            response.raise_for_status()  
            time.sleep(180)  
            raise IOError("Error notice after 3 minutes.")  
        else:  
            response2 = requests.get(redirect_url, headers=headers)  
  
            # Put the results in the empty DataFrame.  
            results = response2.json()  
            hits = results['results']['hits']  
            dfff = pd.json_normalize(  
                hits, 'domains', ['acc', 'name', 'score', 'evalue', 'pvalue', 'desc'])  
            dfff.insert(0, 'sequence', sequence)  
            dfff = dfff.set_index('sequence')  
            results_df = pd.concat([results_df, dfff])  
            if redirect_url == 'late':  
                # Raises an exception if the status is pending for too long.  
                response2.raise_for_status()  
                time.sleep(180)  
                raise IOError("Error notice after 3 minutes.")  
  
    return results_df
```

Results from API System

sequence	alisqacc	alildCount	alirflne	is_included	alihmmname	bitscore	display	ievalue	alisqto	aliSim	...	pvalue
INFWLWPTS		35		1	DCD	28.693352	0.0	6.0e-07	111	0.660550	...	-36.480112
INFWLWPTS		23		0	DCD	15.748980	NaN	0.0051	175	0.672131	...	-36.480112
INFWLWPTS		30		1	dUTPase	40.322865	1.0	0.0	178	0.663636	...	-32.277306
FMARREGR		41		1	AAA	51.984688	1.0	1.0e-13	175	0.690476	...	-40.228785
FMARREGR		41		1	AAA_5	43.017456	0.0	0.0	171	0.669118	...	-34.13238

⋮

823 rows × 221 columns

Results from API System

```
1 # Show the first row and the information of the columns
2 pd.set_option('display.max_rows', None) # Show all details of the first row
3 print(Result_test.iloc[0])
```

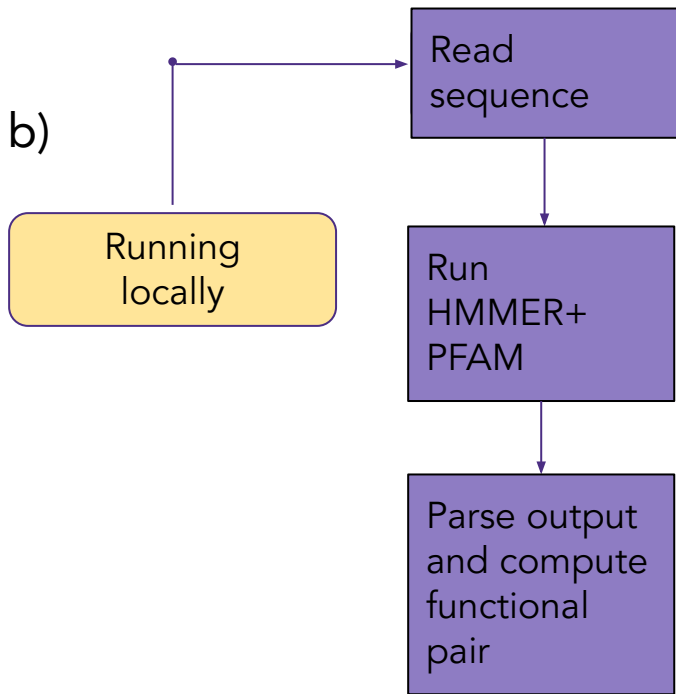
```
alisqacc
aliIdCount      27
alirflne
is_included      1
alihmmname      Sigma70_r2
bitscore      67.940552
display      1.0
ievalue      0.0
alisqto      113
aliSim      0.884058
jali      113
bias      1.66
ienv      45
cevalue      0.0
significant      1.0
alimline      1++ + p+v++l rr+l++ta ae+++Qe+++ +wr+++rfdp++g...
alihmmfrom      1
clan      CL0123
aliL      206
alihindex      16681
is_reported      1
alintseq
jenv      114
alimline
alihmmacc      PF04542.17
oasc      0.96
aliaseq      LYDLLAPRVYGLIRRVLDPALAEVTTQEVLEVWRRARFDPAQG...
alihmto      70
aliId      0.391304
alippline      688899*****
alimodel      lverylplvkrllarrllgsgadaeDlvQegflrlwaverfdperg...
aliM      71
iali      45
alicline      HHHHHHHHHHHHHCTCHHHHHHHHHHHHHHHHHHHHHGCCTTTC...
aliSimCount      61
alihmndesc      Sigma-70 region 2
alisqdesc
```

```
outcompeted      0.0
alisqname      >seq
alisqfrom      45
uniq      1.0
aliN      70
acc      PF04542.17
name      Sigma70_r2
score      68.9
evaluate      0.0
pvalue      -52.678535
desc      Sigma-70 region 2
```

```
Name: MAESGTSRRADHLVPVPGDAEPPAVALDELLRAVGRGDEQAFGRLYDLLAPRVYGLIRRVLDPALAEVTTQEVLEVWRRARFDPAQGSANAWVFTIAHRRRAVDVRAE
QKAADRTVRAGAAALDSPYSDVADEVSGRLERRQVRHCLDALTGQLQREVVTLAYYQGHSTPQVAELLKTPGTVKTRMRDGLIRLRDCLGVEATA, dtype: object
```

Computing functional information from protein pairs

Given two sequences, are we able to see if they're functional or not?



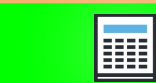
#	#																						
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	
#	#	target name	accession	len	query name	accession	qlen	E-value	score	bias	# of	c-Evalue	t-Evalue	score	bias	from	to	from	to	acc	description of target		
4	Sigma70_v2	PF08542.17	71	875	-	202	4e-19	68.3	0.1	1	1	4.3e-22	8.5e-19	67.2	0.1	5	70	24	87	22	88	8.64 Sigma-70 region 2	
5	Sigma70_v4_2	PF08201.15	54	875	-	202	2.4e-17	62.4	0.6	1	3	1.6	3.2e+03	-2.0	0.1	6	22	20	36	13	37	8.01 Sigma-70, region 4	
6	Sigma70_v4_2	PF08201.15	54	875	-	202	2.4e-17	62.4	0.6	2	3	0.89	1.7e+03	-1.1	0.0	6	13	119	126	116	130	8.64 Sigma-70, region 4	
7	Sigma70_v4_2	PF08201.15	54	875	-	202	2.4e-17	62.4	0.6	3	3	7.2e-10	1.4e-16	60.0	0.0	4	54	128	139	126	139	8.65 Sigma-70, region 4	
8	Sigma70_v4	PF08545.19	50	875	-	202	1.2e-11	44.1	0.1	1	2	2	3.9e+03	-2.4	0.0	38	48	114	124	114	126	8.69 Sigma-70, region 4	
9	Sigma70_v4	PF08545.19	50	875	-	202	1.2e-11	44.1	0.1	2	2	2.3e-24	4.6e-11	42.2	0.0	1	50	132	181	132	191	8.69 Sigma-70, region 4	
10	Sigma70_ECF	PF07638.14	105	875	-	202	0.0003	20.9	0.1	1	1	2.5e-07	0.0005	20.1	0.1	51	101	44	101	24	105	8.77 ECF sigma factor	
11	Htlr_23	PF13384.9	50	875	-	202	0.0036	17.3	0.0	1	1	5.4e-06	0.011	15.8	0.0	10	40	143	174	141	179	8.66 Homodimer-like domain	
12	Xre-Like-Htlr	PF20452.1	63	875	-	202	0.0043	17.6	0.0	1	1	9e-06	0.018	15.6	0.0	10	50	146	179	140	191	8.68 Antitoxin Xre-like helix-turn-helix domain	
13	DUF1350	PF07002.14	250	875	-	202	0.0121	14.2	0.0	1	1	1.4e-05	0.020	13.0	0.0	119	174	19	63	11	109	8.63 Protein of unknown function (DUF1350)	
14	DUF2002	PF10075.11	60	875	-	202	0.009	13.5	0.5	1	2	0.00015	0.3	11.4	0.0	30	53	26	49	9	54	8.61 Protein of unknown function (DUF2002)	
15	DUF2002	PF10075.11	60	875	-	202	0.009	13.5	0.5	2	2	0.77	1.5e+03	-0.4	0.1	11	11	131	131	83	174	8.65 Protein of unknown function (DUF2002)	
16	DUF134	PF08011.19	90	875	-	202	0.12	13.1	0.0	1	1	0.00012	0.23	12.2	0.0	42	83	143	184	111	202	8.63 Protein of unknown function (DUF134)	
17	UOCC2_CBP6	PF20106.1	90	875	-	202	0.34	11.6	2.7	1	3	1.4	2.8e+03	-1.0	0.0	57	69	43	55	29	67	8.70 Complex III assembly factor UOCC2_CBP6	
18	UOCC2_CBP6	PF20106.1	90	875	-	202	0.34	11.6	2.7	2	3	0.00073	1.4	9.6	0.0	56	81	70	93	44	102	8.62 Complex III assembly factor UOCC2_CBP6	
19	UOCC2_CBP6	PF20106.1	90	875	-	202	0.34	11.6	2.7	3	3	0.6	1.2e+03	0.2	0.0	42	63	107	120	95	141	8.69 Complex III assembly factor UOCC2_CBP6	

Goes to component 5

Sample preparation

Model development

Deployment



Data retrieval



Data format &
Feature selection

X data



Pfam - Family
identification

Y data



Model
development

Y data



Protein Data Bank
- Structural info

On-going and finish
in Spring

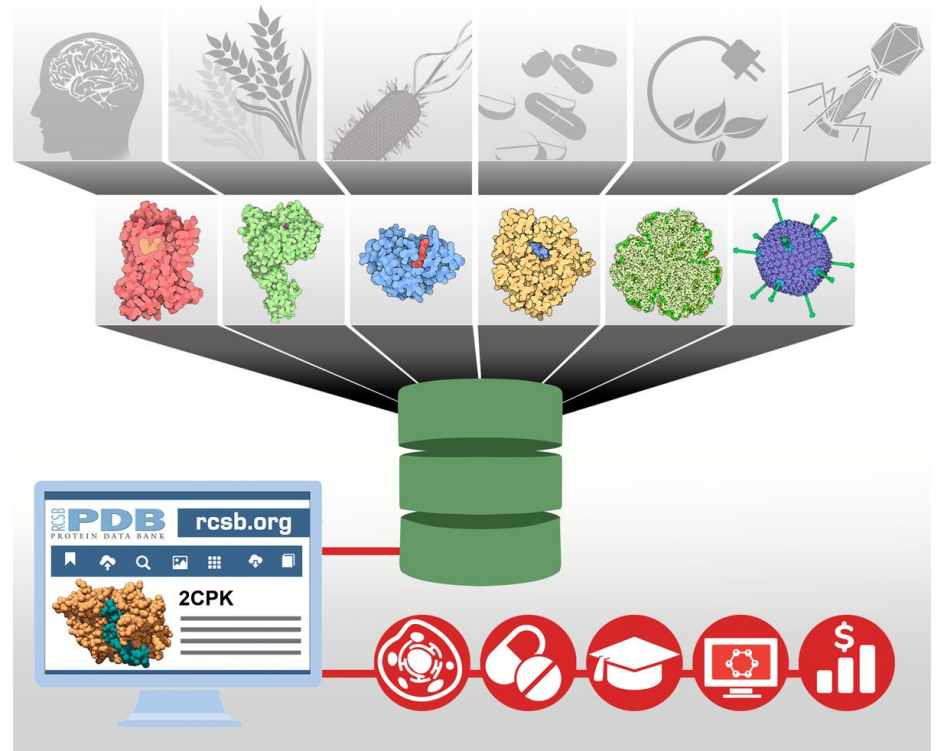
Finding molecular structure information

Goal:

- 1) Query into Protein Data Bank (PDB)
- 2) Obtain lists of proteins with structural similarity

Upstream: Sampling component

Downstream: ML function prediction



Finding molecular structure information

MLLSDRDLVSEIKSGDLSLEPFEPALLQPS...



Progress:

- 1) Successfully access PDB
- 2) Obtain lists of proteins with sequence similarity

Next:

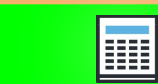
- 1) Modify the search and return to obtain output based on structures.

	identifier	score	sequence_identity	evalue	bitscore
0	2QXX_1	1.0	0.790	9.055000e-94	304
1	4A6A_1	0.992308	0.784	3.197000e-93	302
2	2QLP_1	0.842308	0.788	1.157000e-79	263
3	1XS1_1	0.365385	0.431	1.379000e-36	139
4	2V9X_1	0.361538	0.426	3.537000e-36	138
5	1XS4_1	0.357692	0.426	9.071000e-36	137
6	1XS6_1	0.357692	0.426	9.071000e-36	137
7	2J4Q_1	0.357692	0.426	9.071000e-36	137
8	2J4H_1	0.35	0.426	3.184000e-35	135
9	4XJC_1	0.2	0.338	6.627000e-22	96

Sample preparation

Model development

Deployment



Data retrieval



Data format &
Feature selection

X data



Pfam - Family
identification

Y data



Model
development

Y data



Protein Data Bank
- Structural info

On-going and finish
in Spring

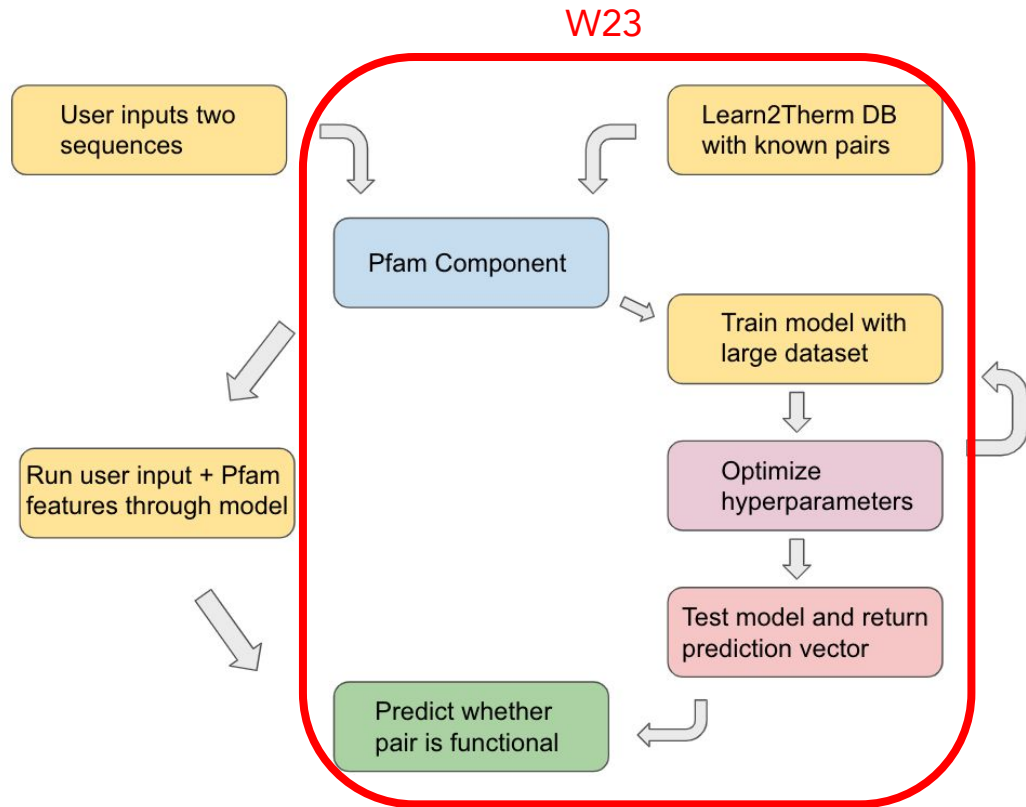
Training and validating a model to predict protein function

Goals:

1. Train a RandomForestClassifier on Learn2Therm DB in order to predict whether a protein pair is functional
2. Create robust pipeline that can handle expanded feature space

Upstream: Pfam + Data Retrieval components

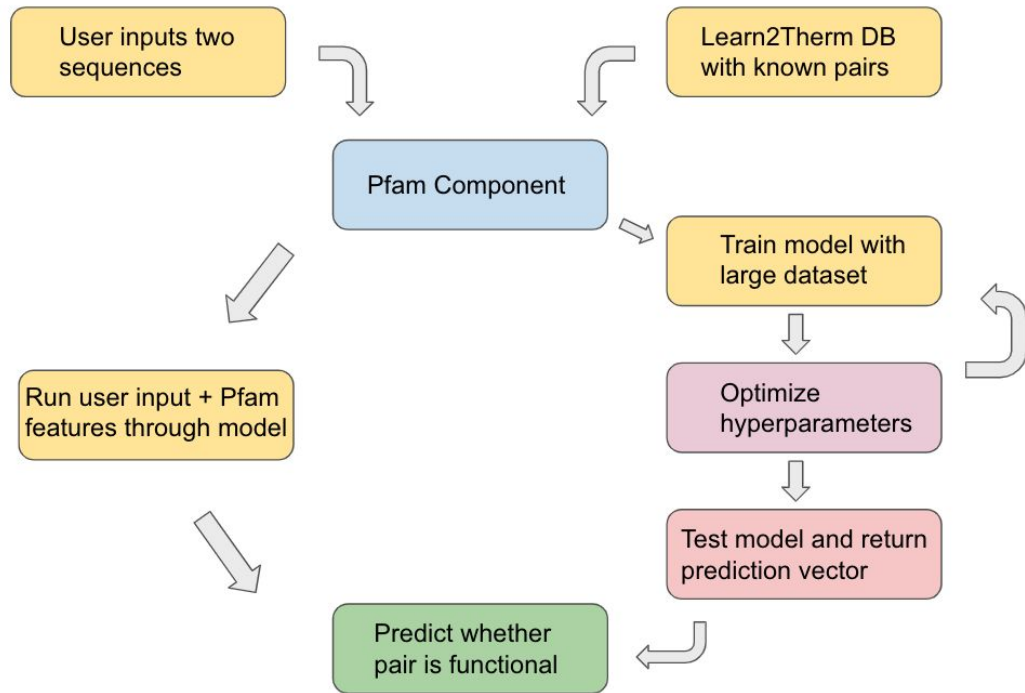
Downstream: Functionality analysis (future work)



Training and validating a model to predict protein function

Ongoing work:

1. Combine output of MSA and structural scraping to improve strength of model
2. Use model to predict functionality from a user input of two AA sequences
3. Develop metrics of functionality to be reported in addition to Boolean classification



Summary

- ❖ Obtained a training dataset from an existing SQL database.
- ❖ Ran the HMMER algorithm against Pfam database.
- ❖ Queried sequence and got output from PDB.
- ❖ Trained and validated a RandomForestClassifier to predict protein pair functionality.



Future work

- Make a technical validation pipeline that explores MSAs (evolutionary history) and structure (e.g. SWISS-model) purely from sequence pair inputs from users
 - Component 0-2: Run BLAST locally to obtain local alignment from users' protein pairs + Sample analytics
 - Component 3: More sophisticated parsing to get functional information for various pairs
 - Component 4: Integration with downstream and upstream components
 - Component 5: Increase feature space of pair functionality prediction model

Acknowledgement

Consultation and Code Review:

David Beck (PhD)
Nels Schimek
Orion Dollar

Development of Learn2Therm
DB & EXTENSIVE project
consultation:

Evan Komp

Questions?
