

SPARCLE Curation

predicting protein architecture names

2024 NCBI Codeathon - Team Gwadz-Yang

Marc Gwadz (Team Leader), NCBI, NLM, NIH

Mingzhang Yang (Technical Lead), NCBI, NLM, NIH

Christopher Meyer (Writer), University of Chicago - Center for Translational Data Science

Franziska Ahrend, ORISE Fellow NIDDK, NIH

Yixiang Deng, Post-doc MIT, Harvard

Shaojun Xie, Advanced Biomedical and Computational Science, Frederick National Laboratory for Cancer Research

Goal

Input/processing

Model

Evaluation

Inference

Future Work

SPARCLE, the Subfamily Protein Architecture Labeling Engine, is a resource for the functional characterization and labeling of protein sequences that have been grouped by their **characteristic conserved domain architecture**.

A domain architecture is defined as the sequential order of conserved domains in a protein sequence.

Conserved domains can be clustered into **superfamilies** that generate over-lapping annotations on the same protein

Goal: Develop Machine Learning methods to use curated architectures to apply names (and other information) to un-curated architectures

SPARCLE > Architecture Viewer

non-ribosomal peptide synthetase

non-ribosomal peptide synthetase is a modular multidomain enzyme that acts as an assembly line to catalyze the biosynthesis of complex natural products

Arch. ID	Version	Date Published	Review Level
13483227	2	2023-05-09	curated

PRK12467 NRPS-para261 PRK12316

Sequences with this architecture

This architecture currently does not link to any protein sequence records.

↑ Do Table of Contents

Name, label and taxonomic scope		
Taxonomic scope	Name	Label
All organisms	non-ribosomal peptide synthetase	non-ribosomal peptide synthetase is a modular multidomain enzyme that acts as an assembly line to catalyze the biosynthesis of complex natural products

Supporting evidence

Protein Accession	Evidence
• PRK12316	non-ribosomal peptide synthetase, partial (Pseudomonas aeruginosa)
• ESQ64119	fusaricidin synthetase (Pseudomonas aeruginosa HB15)

CDO

- TIGR01720: NRPS para261: non-ribosomal peptide synthetase domain TIGR01720
- smr00022: PKS_PP: Phosphopantetheine attachment site
- pfam00501: AMP-binding: AMP-binding enzyme
- TIGR01732: AA-adenyl-dom: amino acid adenylation domain
- pfam00668: Condensation: Condensation domain
- COG1020: EntF: Non-ribosomal peptide synthetase component F [Secondary metabolites biosynthesis, transport and catabolism]
- cd17649: A_NRPS_PvdJ-like: non-ribosomal peptide synthetase
- PRK12316: peptide synthase, Provisional
- PRK12467: peptide synthase, Provisional

Conserved domains on [gi|557216254|gb|ESQ64119|]

fusaricidin synthetase [Pseudomonas aeruginosa HB15]

View Standard Results

Graphical summary Zoom to residue level show extra options

Query seq. 1 250 500 750 1000 1250 1500 1750 2000 2250 2500 2750 3000 3250

Specific hits

Non-specific hits

Superfamilies

Search for similar domain architectures Refine search

#	Name	Accession	Description	Interval	E-value
[H]	PRK12316	PRK12316	peptide synthase, Provisional	599-3200	0e+00
[H]	A_NRPS_PvdJ-like	cd17649	non-ribosomal peptide synthetase; This family of the adenylation (A) domain of nonribosomal ...	1074-1564	0e+00
[H]	E_NRPS	cd19534	Epimerization domain of nonribosomal peptide synthetases (NRPS); belongs to the ...	161-580	1.97e-169
[H]	EntF	COG1020	Non-ribosomal peptide synthetase component F [Secondary metabolites biosynthesis, transport ...	830-1482	1.83e-166
[H]	AA-adenyl-dom	TIGR01733	amino acid adenylation domain; This model represents a domain responsible for the specific ...	2606-3017	3.17e-149
[H]	AMP-binding	pfam00501	AMP-binding enzyme;	2585-2993	9.11e-96
[H]	PRK12467	PRK12467	peptide synthase, Provisional	6-415	6.33e-67
[H]	Condensation	pfam00668	Condensation domain; This domain is found in many multi-domain enzymes which synthesize ...	160-582	2.78e-65
[H]	NRPS-para261	TIGR01720	non-ribosomal peptide synthetase domain TIGR01720; This domain appears to be located immediately ...	451-606	8.83e-57
[H]	entF	PRK10252	enterobactin non-ribosomal peptide synthetase EntF;	5-134	9.14e-31
[H]	A_NRPS_PvdJ-like	cd17649	non-ribosomal peptide synthetase; This family of the adenylation (A) domain of nonribosomal ...	6-63	9.70e-19
[H]	alpha_am_amid	TIGR03443	L-aminoacylate-semialdehyde dehydrogenase; Members of this protein family are ...	7-118	1.08e-15
[H]	PP-binding	pfam00550	Phosphopantetheine attachment site; A 4-phosphopantetheine prosthetic group is attached ...	83-141	8.46e-13
[H]	EntF	COG1020	Non-ribosomal peptide synthetase component F [Secondary metabolites biosynthesis, transport ...	372-579	4.63e-06
[H]	PKS_PP	smart00823	Phosphopantetheine attachment site; Phosphopantetheine (or pantetheine 4' phosphate) is the ...	3098-3168	2.99e-04
[H]	PKS_PP	smart00823	Phosphopantetheine attachment site; Phosphopantetheine (or pantetheine 4' phosphate) is the ...	70-142	4.43e-04

Reset search parameters

Goal

Input Features: SuperFamilies, Specific Architectures, Title Strings
Output : Curated Names

Input/processing

Target output :
Curated Names

Features to train Models: Superfamily and specific domain models that make up the architecture. Descriptive text associated with specific domains

Model

Evaluation

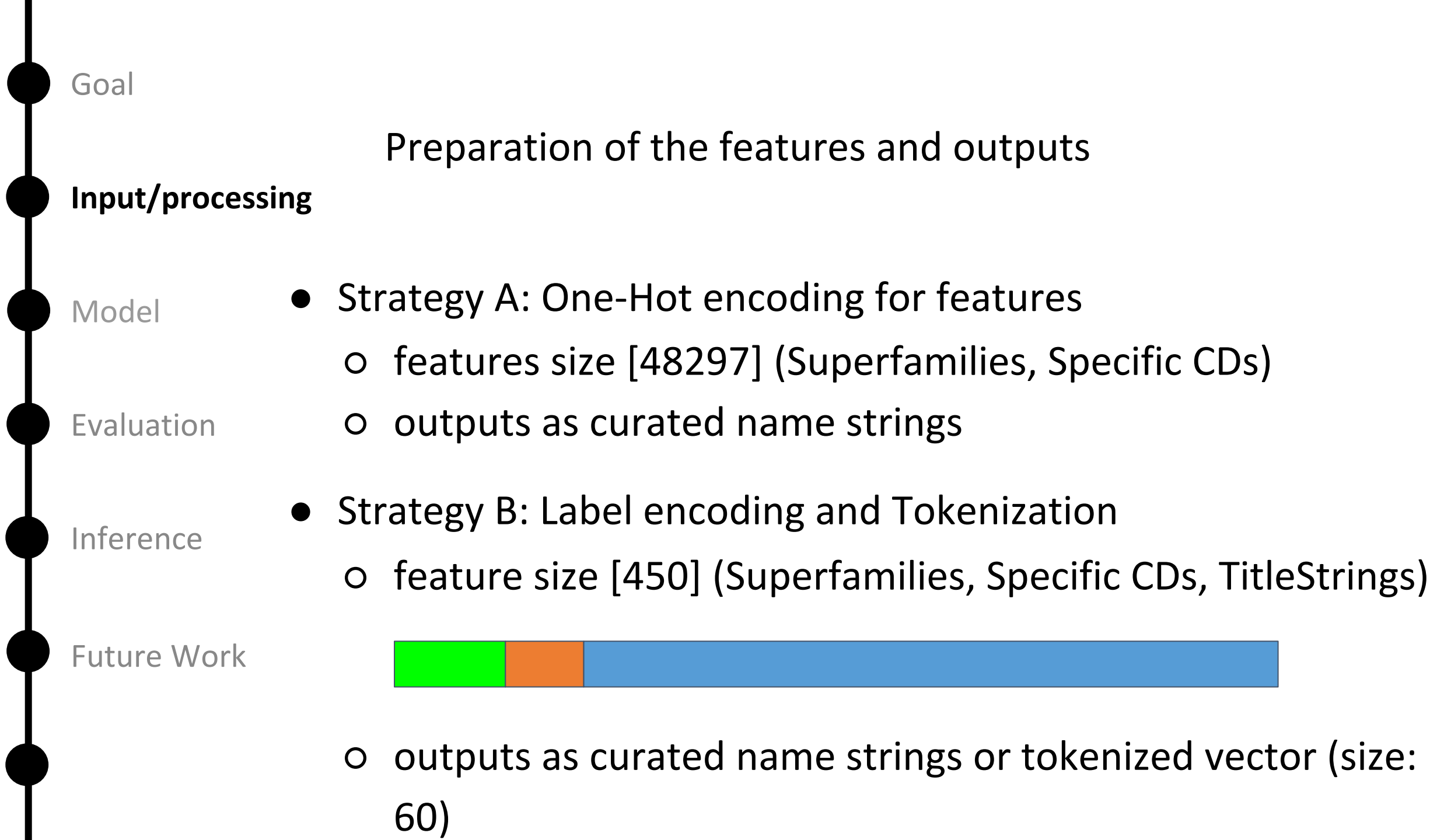
Inference

Future Work

Curated Name	LabelState	SpecificArch	superfamilyarch	TitleStrings
non-ribosomal peptide synthetase	curated	cl36129 TIGR01720 cl36106	NRPS-para261 PRK12316 PRK12467	PRK12467 non-ribosomal peptide synthase domain TIGR01720 PRK12316
NRPS-para261 and AFD_class_I domain-containing protein	namedByDomain	TIGR01720 cl36106	NRPS-para261 PRK12316	non-ribosomal peptide synthase domain TIGR01720 PRK12316
NRPS-para261 and A_NRPS_Cytc1-like domain-containing protein	namedByDomain	cl36129 cl36106 TIGR01720 cl36106 cl36129 COG1020	NRPS-para261 PRK12316 PRK12467 EntF	PRK12467 PRK12316 non-ribosomal peptide synthase domain TIGR01720 PRK12316 PRK12467 Non- ribosomal peptide synthetase component F [Secondary metabolites biosynthesis, transport and catabolism]
A_NRPS_Srf_like and A_NRPS domain-containing protein	namedByDomain	cl36129 COG3321 cl36129 cl11771	NRPS-para261 PRK12467 PksD	PRK12467 Acyl transferase domain in polyketide synthase (PKS) enzymes [Secondary metabolites biosynthesis, transport and catabolism] PRK12467 NRPS-para261

Pre-processing :

- Curated Names were manually reduced to decrease # categories and help in testing.
- Common / uninformative words in TitleStrings (articles, prepositions, “protein” etc.) were eliminated prior to vectorization.
- Input/Output Data organized into One-Hot encoded matrices and vectors (SentencePiece)



Goal

Input/processing

Model

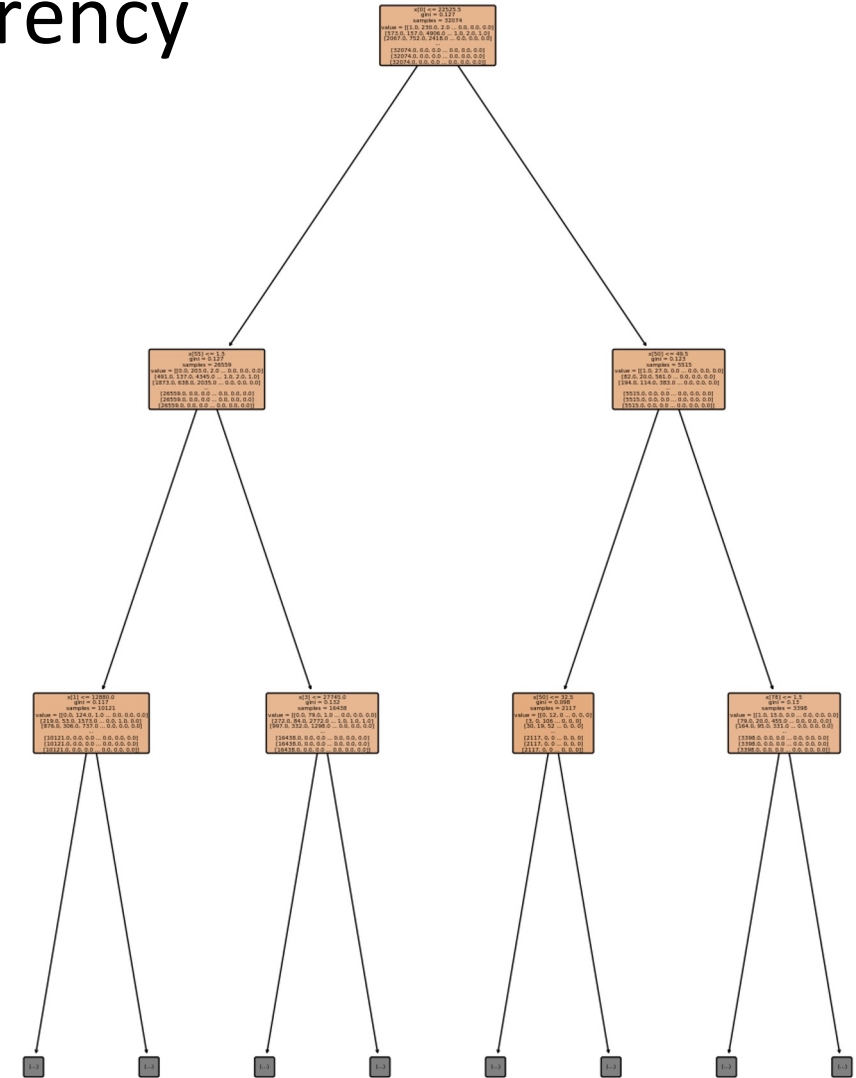
Evaluation

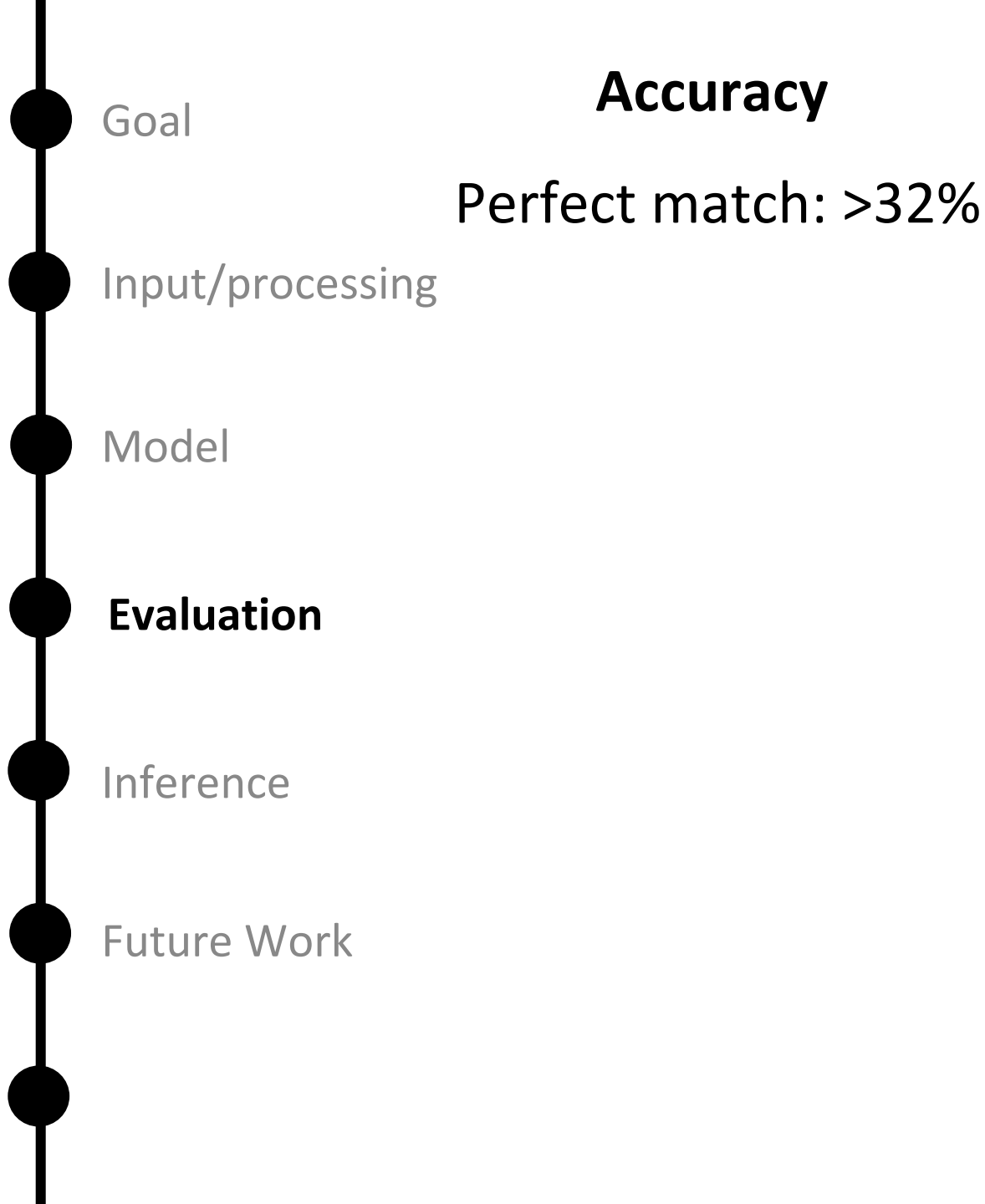
Inference

Future Work

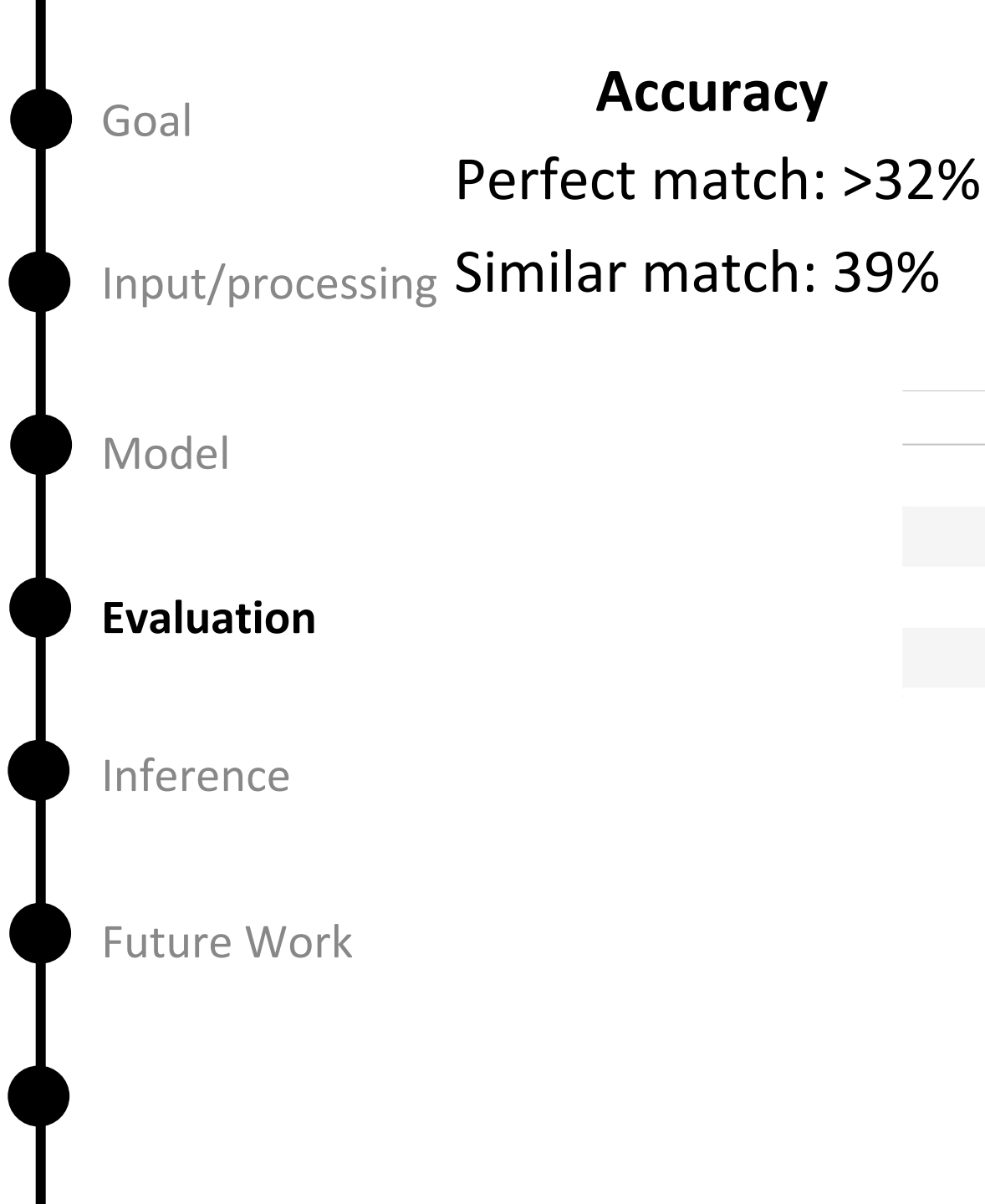
Decision Tree

- Interpretability and Transparency
- Categorical Data Handling
- Reveals Feature Importance
- Models Non-Linear Relationships





Actual	Predicted	Equal
T-complex protein 1 subunit	T-complex protein 1 subunit	True
mannose-6-phosphate isomerase	mannose-6-phosphate isomerase	True
PUA domain-containing protein	PUA domain-containing protein	True
hybrid sensor histidine kinase/response regulator	hybrid sensor histidine kinase/response regulator	True
hybrid sensor histidine kinase/response regulator	hybrid sensor histidine kinase/response regulator	True
lysine decarboxylase	lysine decarboxylase	True
SANT/Myb-like DNA-binding domain-containing pr...	SANT/Myb-like DNA-binding domain-containing pr...	True
MFS transporter	MFS transporter	True
enoyl-CoA hydratase	enoyl-CoA hydratase	True
LCP family protein	LCP family protein	True
cell division protein	cell division protein	True
methyl-accepting chemotaxis protein	methyl-accepting chemotaxis protein	True
E3 ubiquitin-protein ligase	E3 ubiquitin-protein ligase	True
helix-hairpin-helix domain-containing protein	helix-hairpin-helix domain-containing protein	True
adenylate/guanylate cyclase domain-containing ...	adenylate/guanylate cyclase domain-containing ...	True
cadherin repeat domain-containing protein	cadherin repeat domain-containing protein	True
LacI family DNA-binding transcriptional regulator	LacI family DNA-binding transcriptional regulator	True
BTB/POZ domain and ankyrin repeat-containing p...	BTB/POZ domain and ankyrin repeat-containing p...	True
dTDP-glucose 4,6-dehydratase	dTDP-glucose 4,6-dehydratase	True
class I SAM-dependent methyltransferase	class I SAM-dependent methyltransferase	True
FRMD7 family protein	FRMD7 family protein	True
non-ribosomal peptide synthetase	non-ribosomal peptide synthetase	True
BTB/POZ domain-containing protein	BTB/POZ domain-containing protein	True
cation diffusion facilitator family transporter	cation diffusion facilitator family transporter	True
transposase	transposase	True



Actual	Predicted
peroxisomal biogenesis factor	peroxisome biogenesis factor
methylmalonyl-CoA mutase subunit	methylmalonyl-CoA mutase
ubiquitin carboxyl-terminal hydrolase family p...	ubiquitin carboxyl-terminal hydrolase
C2HC-type zinc finger protein	C2H2-type zinc finger protein
AP-4 complex subunit	AP-3 complex subunit

Actual	Predicted
DUF4299 domain-containing protein	DUF4278 domain-containing protein
DUF4345 domain-containing protein	DUF4346 domain-containing protein
DUF3015 domain-containing protein	DUF2515 domain-containing protein
DUF5043 domain-containing protein	DUF5035 domain-containing protein

Actual	Predicted
ParB/RepB/Spo0J family partition protein	ROK family transcriptional regulator
flagellar brake protein	leucine--tRNA ligase
transposase	heat shock 70 family protein
GATOR complex protein	Myb family transcription factor
2-aminoethylphosphonate ABC transporter substr...	low molecular weight protein-tyrosine-phosphat...
PTS phosphocarrier protein NPr	PTS galactitol transporter subunit
DUF3083 family protein	YusG family protein
MspA family porin	class I SAM-dependent methyltransferase
CPBP family intramembrane glutamic endopeptidase	DUF4430 domain-containing protein
coproporphyrinogen III oxidase	NADH dehydrogenase subunit

Goal

Challenges for training/testing :

- Sparse data (many architectures only have rare superfamilies)
- Large number of categories and some inconsistencies in Curated names

Input/processing

While most super superfamilies occur in very few Architectures, others can occur many times in both curated and uncurated Architectures.

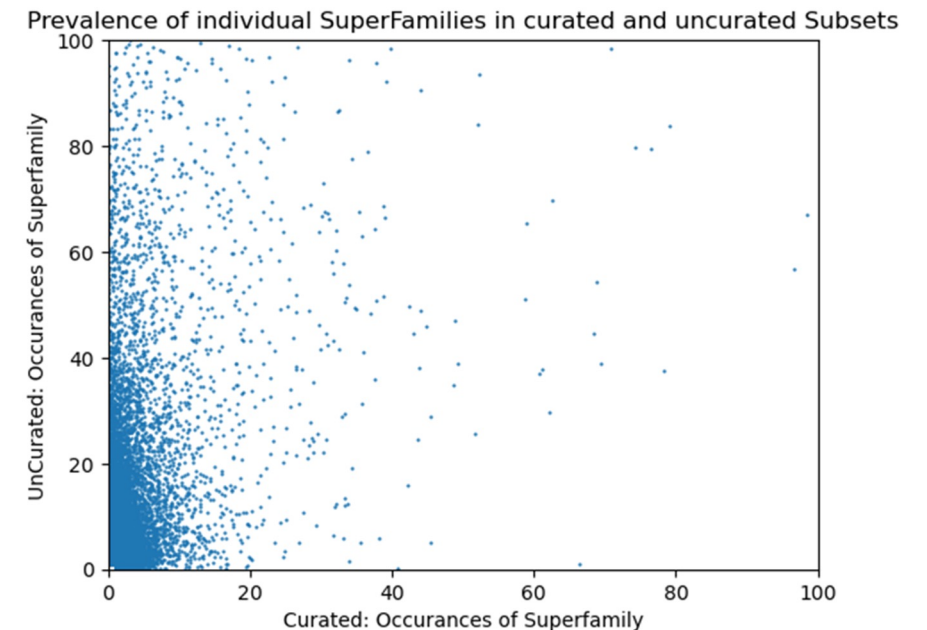
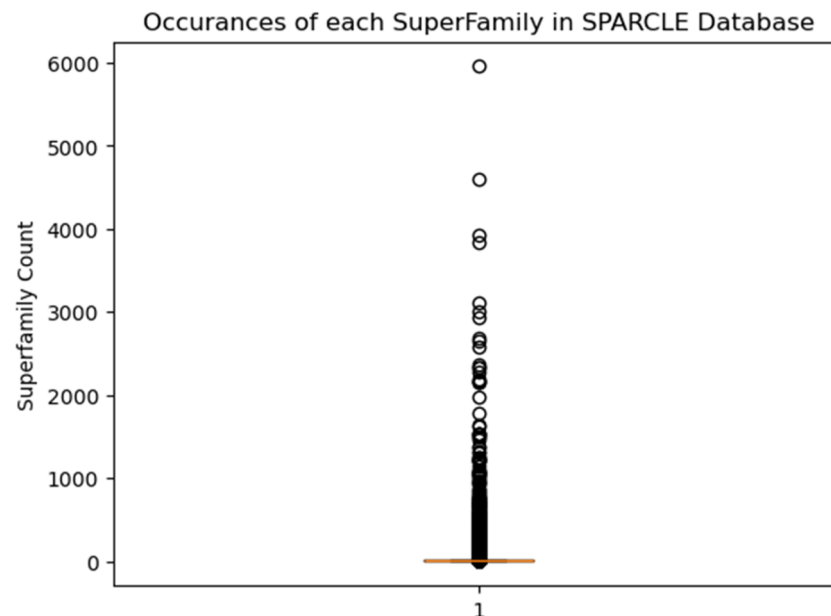
Model

Architectures with frequently occurring superfamilies can hopefully be leveraged to auto-name numerous related architectures.

Evaluation

Inference

Future Work





Goal

Input/processing

Model

Evaluation

Inference

Future Work

- Standardized/generic target names
- Collect all the information / knowledge about every domain model to create a corpus.
 - With this corpus, we can generate word embedding for every single domain model.
 - Then we can create biologically meaningful vectors for the specificArch and superfamilyArch.
- Try other ML methods and employ better validation methods
- Parameter optimization: change “opt” variable to use different regressors.
 - We used “Ridge regression” (opt=ridge), a linear model with performance score $R^2 = 0.25$.
 - This is fastest, but would love to evaluate other methods.
- Evaluate choice and filtering of input features
 - e.g. limiting text to informative strings
- Curator inspection of proposed names for uncurated data

#team-gwadz-yang

Team Leader
Marc Gwadz



Technical Lead
Mingzhang Yang



Writer
Christopher Meyer



Franziska Ahrend



Yixiang Deng



Shaojun Xie

