



RISING STARS IN DATA SCIENCE 2024



Stanford
Data
Science



THE UNIVERSITY OF CHICAGO
DATA SCIENCE
INSTITUTE

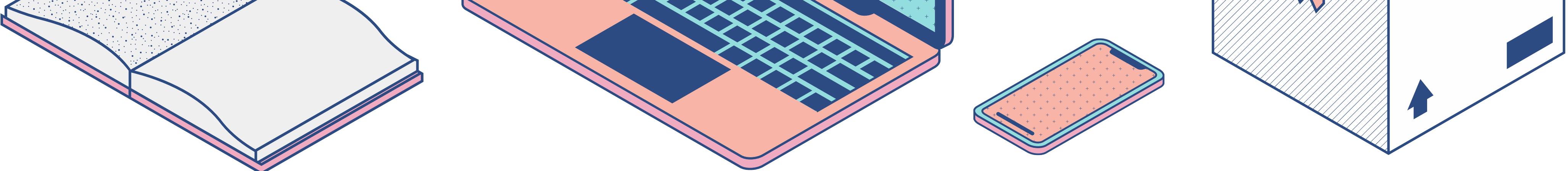


HALİCİOĞLU DATA SCIENCE INSTITUTE



The Future of Human-Centric eXplainable AI is not post-hoc explanations

VINITRA SWAMY
EPFL



Main Takeaways

EVALUATING THE EXPLAINERS
TRUSTING THE EXPLAINERS

Post-hoc explainers have serious problems,
and there is no effective way to validate them



Main Takeaways

MULTIMODN, INTERPRETCC, ILLUMINATE

With interpretable-by-design NNs,
guaranteed interpretability
does not have to come at the cost of performance
or human-understandability

Vinitra Swamy

XAI Research



5th year PhD Student at EPFL



Co-advised by
Tanja Käser (ML4ED) and
Martin Jaggi (MLO)



UC Berkeley → Microsoft AI → EPFL



AI for Education, eXplainable AI



Teaching data science

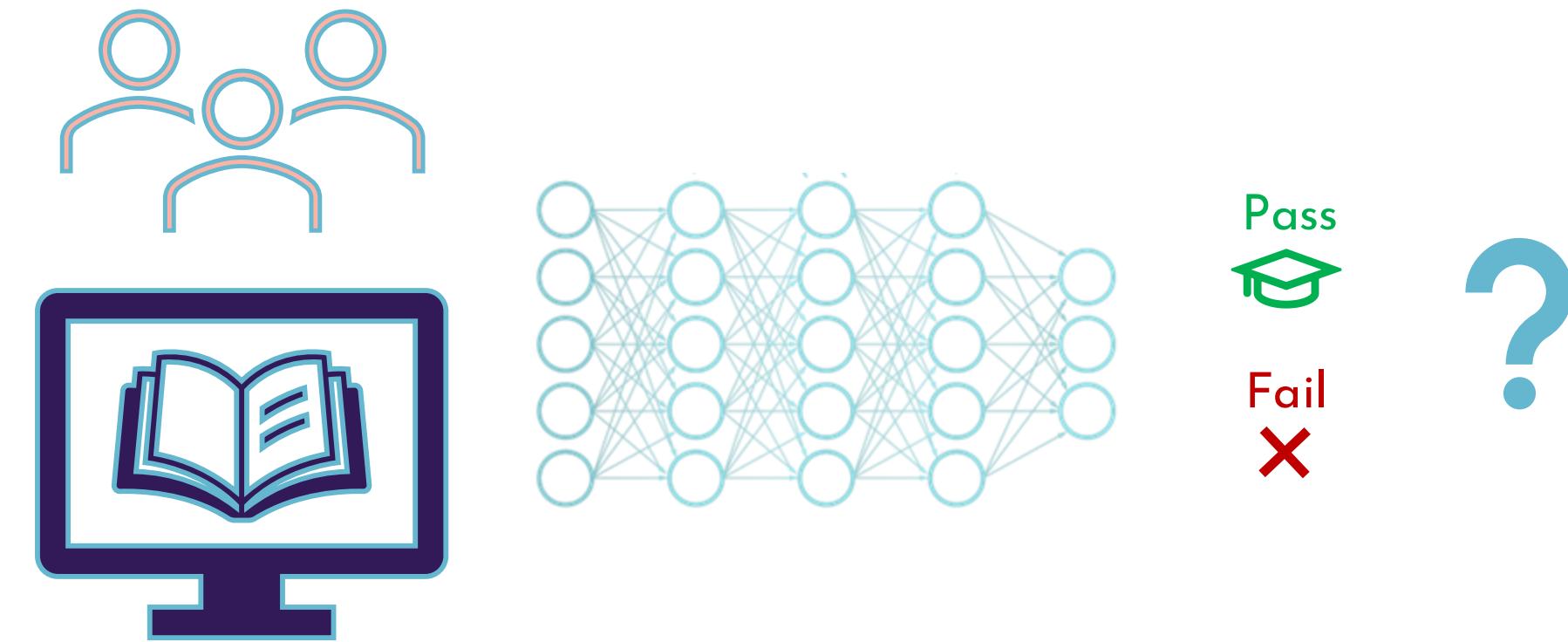
Overview

- 1) Why should we care about explainability?
- 2) Why are existing XAI methods not enough?
- 3) Three ideas towards an eXplainable AI future:
 - MultiModN
 - InterpretCC
 - iLLuMinaTE

Cost of using neural networks

AI for Education

Problem: Our current SoTA trades transparency for accuracy



6

Identifying “why” is important for effective, personalized interventions

Solution: Explainable Machine Learning

Why is eXplainable AI important?

1. Building stakeholder trust in models
2. Auditing models when they make mistakes
3. Improving models with their own reasoning

Framework of Explainability Needs for Human-Centric Computing

Consistent (Li et al., 2021)

[multiple generated explanations are the same]

Real-Time (Xu et al., 2017)

[next minute, next lesson, 1 week, after the course]

Accurate (Marx et al., 2023; Leichtmann et al., 2023)

[model is confident in the explanation]

Actionable (Joshi et al., 2019)

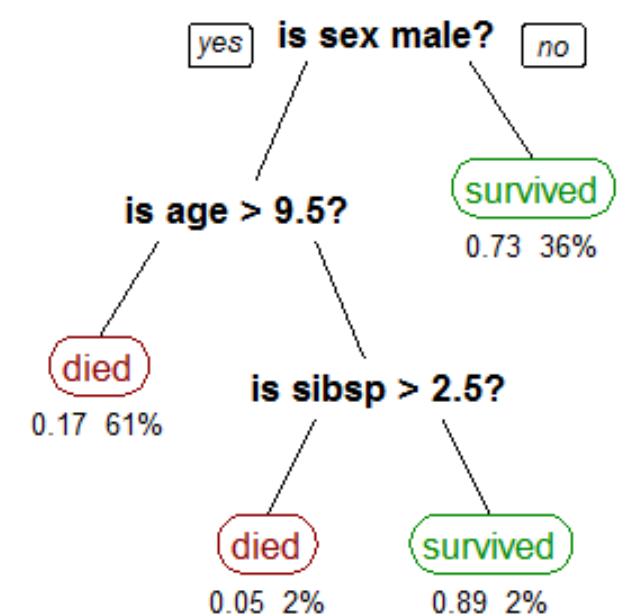
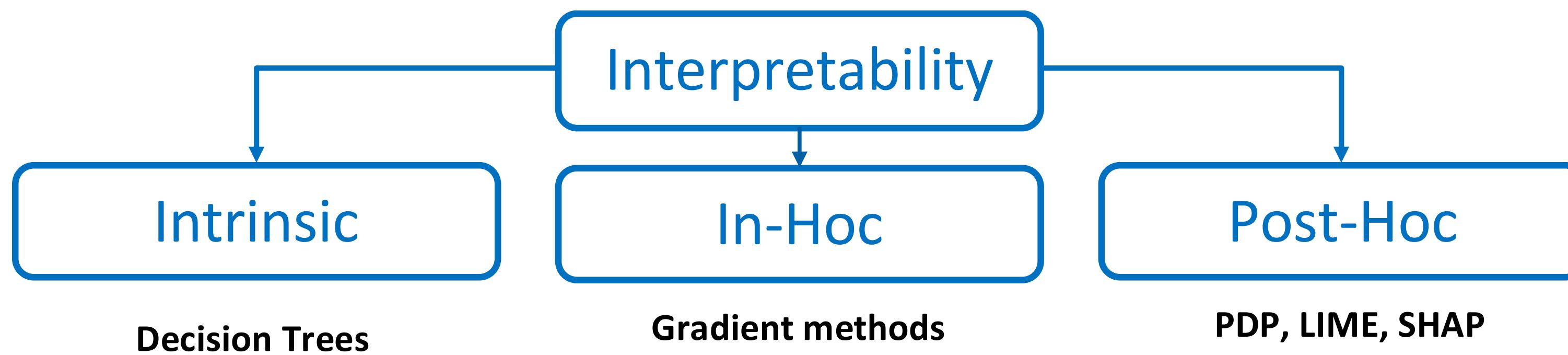
[able to take a next step based on the insight]

Human-Interpretable (Hudon et al., 2021; Haque et al., 2023)

[easy for a non-data scientist to understand]

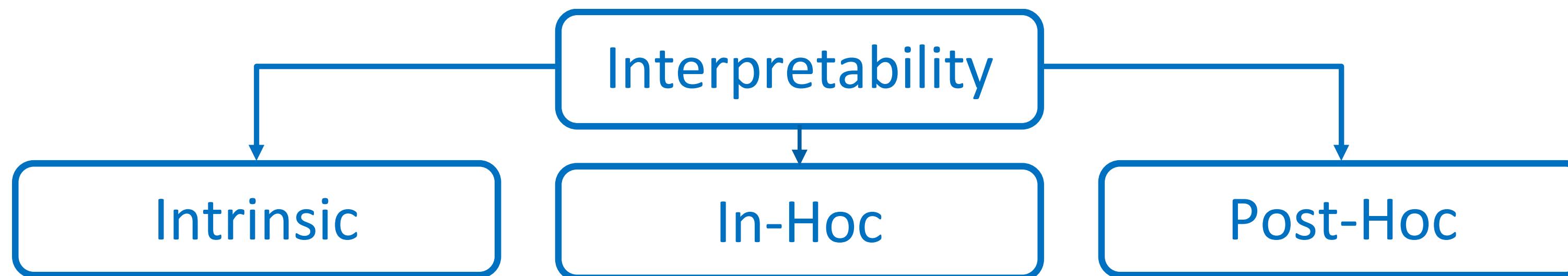
Interpretability

XAI Fundamentals

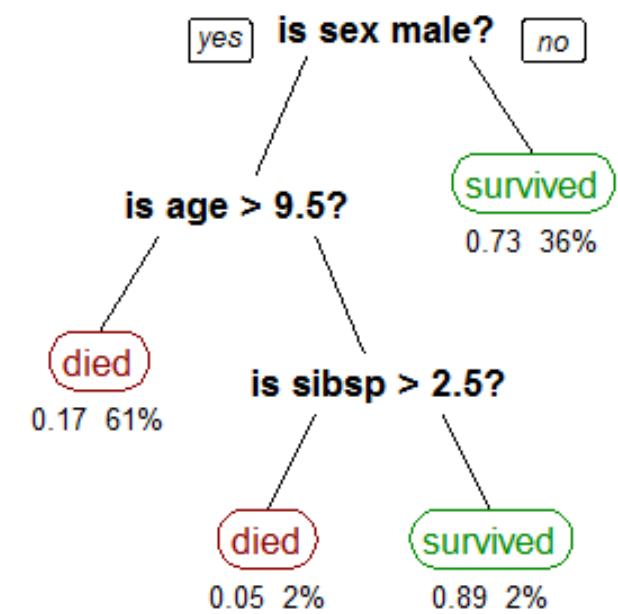


Interpretability

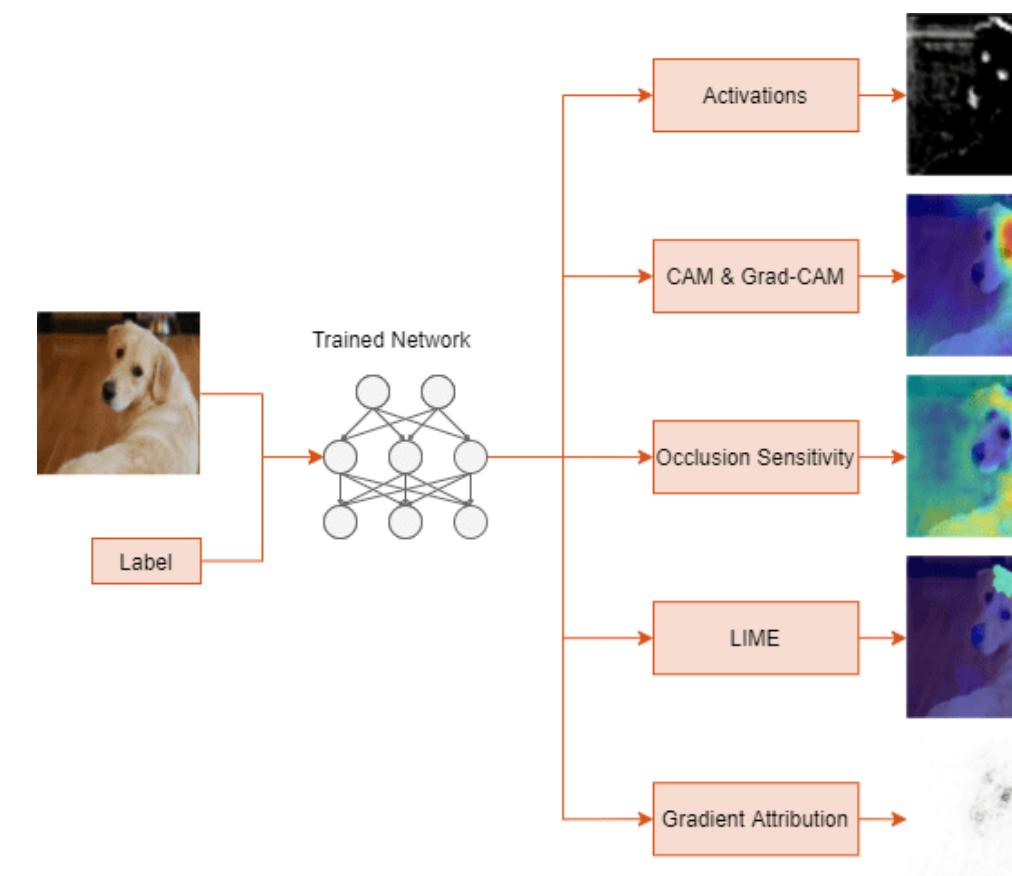
XAI Fundamentals



Decision Trees



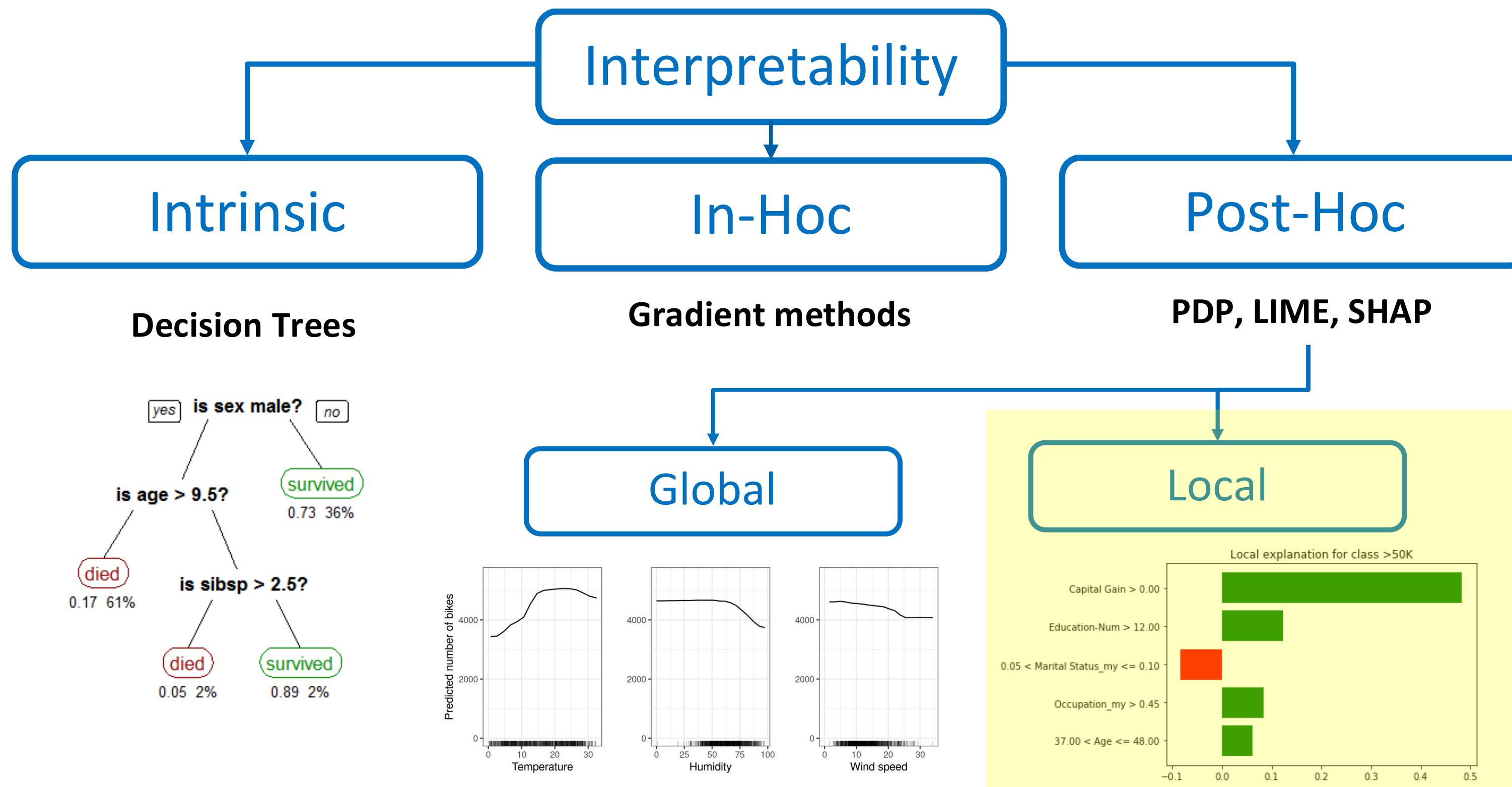
Gradient methods



PDP, LIME, SHAP

Interpretability

XAI Fundamentals



Local XAI

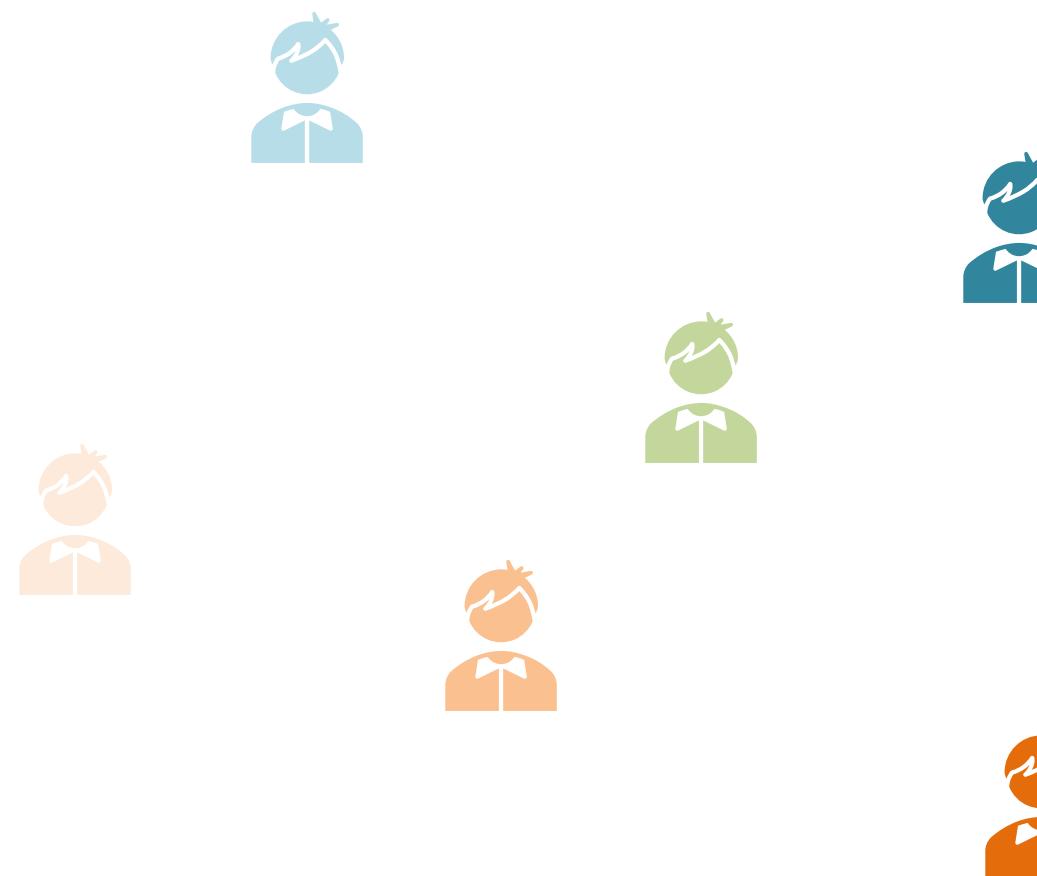
METHODOLOGY

LIME

Local Interpretable Model-Agnostic Explanations

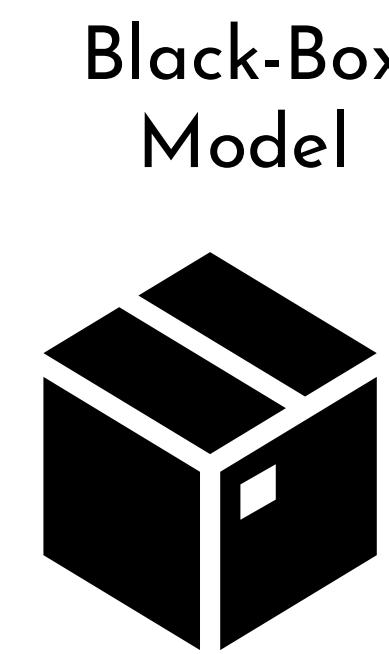
1

Select a specific point to explain: $(X_{\text{student}}, Y_{\text{student}})$

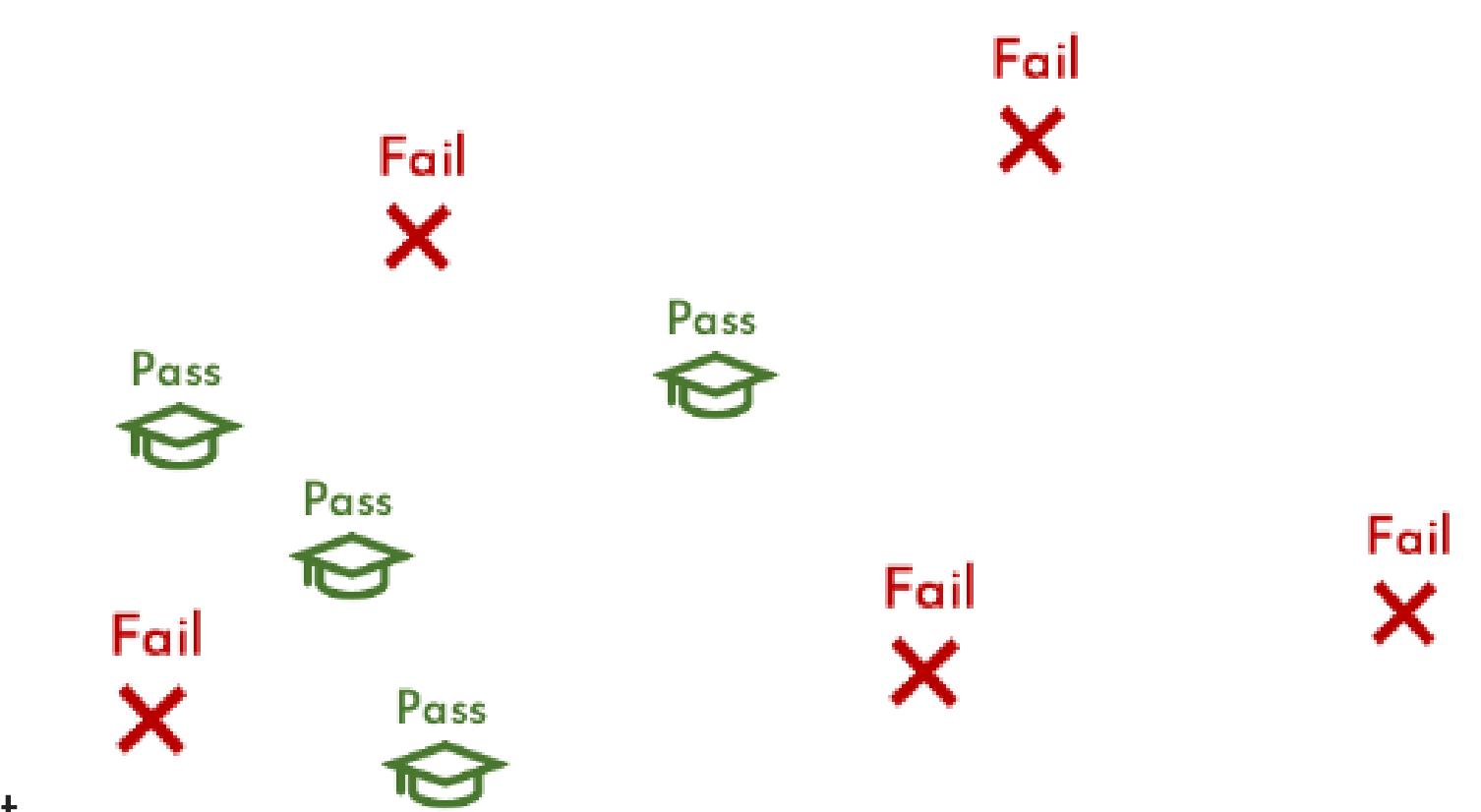


2

Perturb features of selected point to get $\{X^1_{\text{student}} \dots X^N_{\text{student}}\}$ neighbors

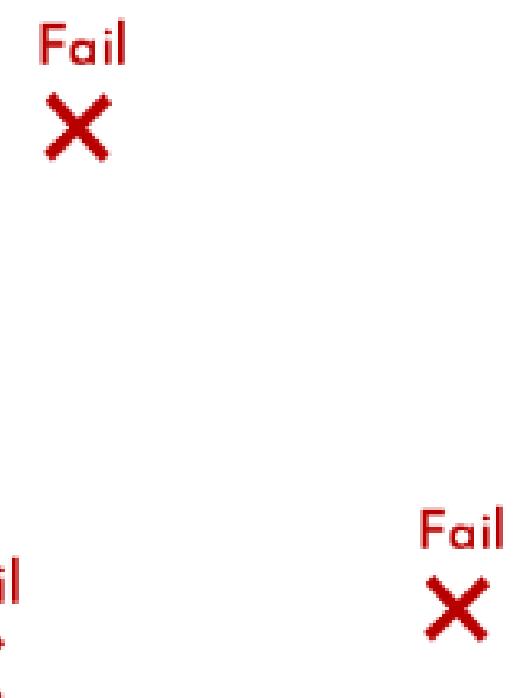


Y_{student}



3

Feed in X_{student} neighbors to the black-box model and get predictions $\{Y^1_{\text{student}} \dots Y^N_{\text{student}}\}$



Local XAI

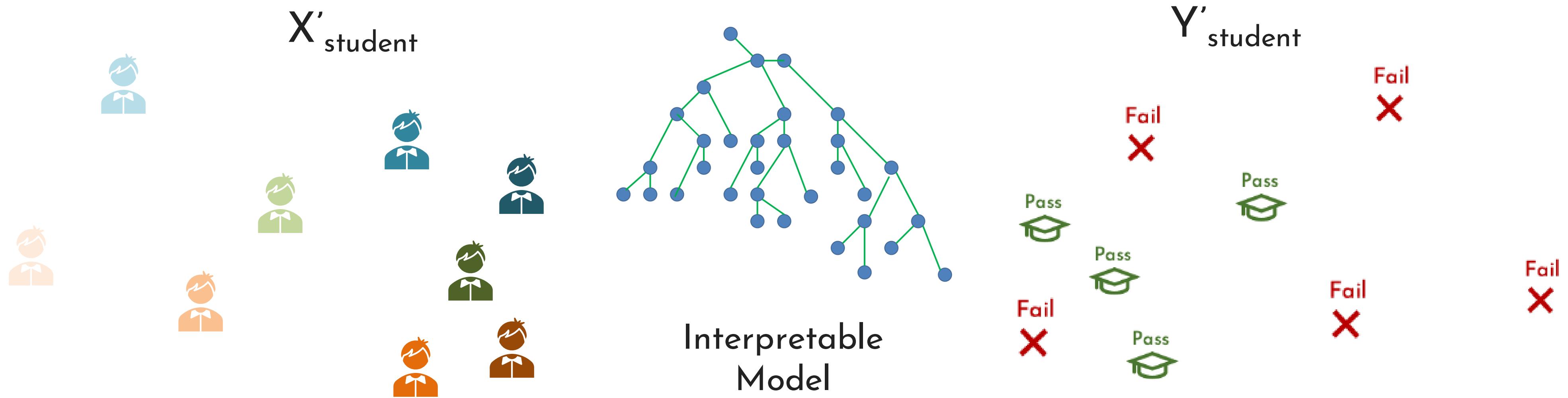
METHODOLOGY

LIME

Local Interpretable Model-Agnostic Explanations

4

Train an interpretable local model using (weighted) X'_{student} and Y'_{student}

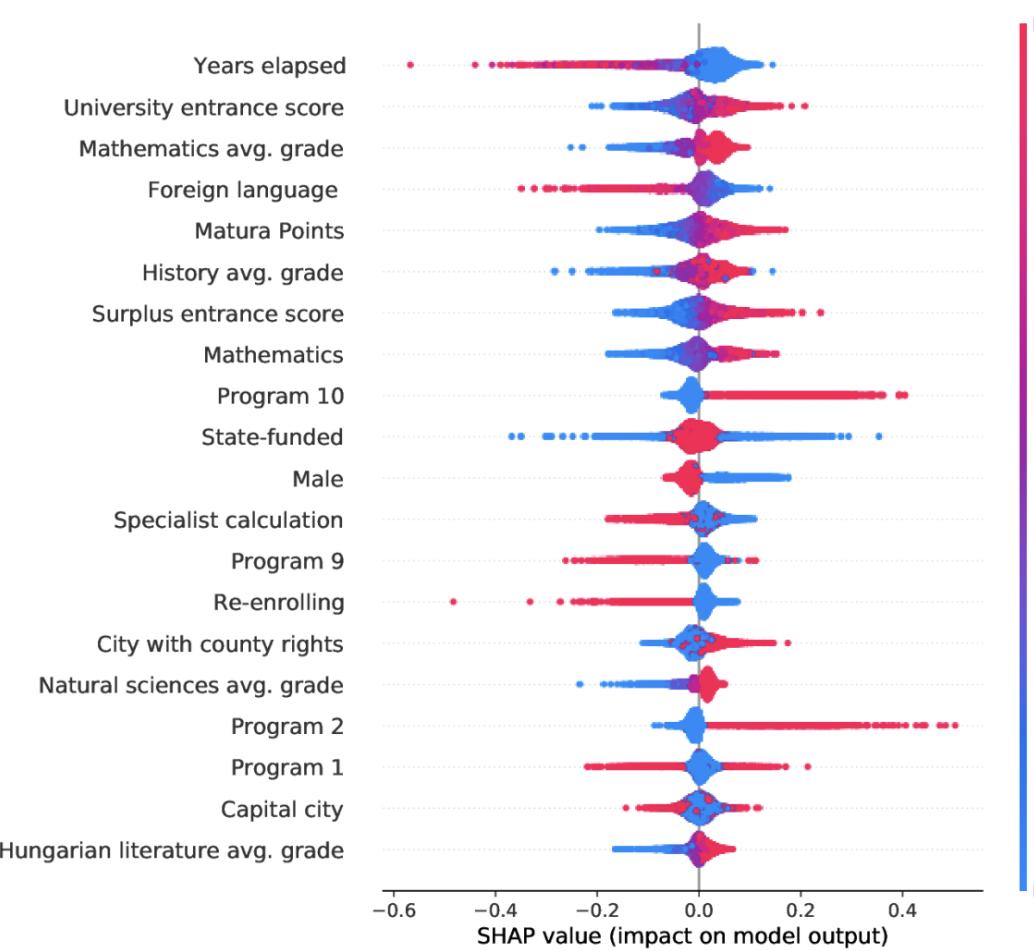


Previous Work

MOTIVATION

Previous work: In (minimal) related literature, only one explainability method is picked per application

SHAP for student dropout^[1]



LIME for student advising^[2,3]



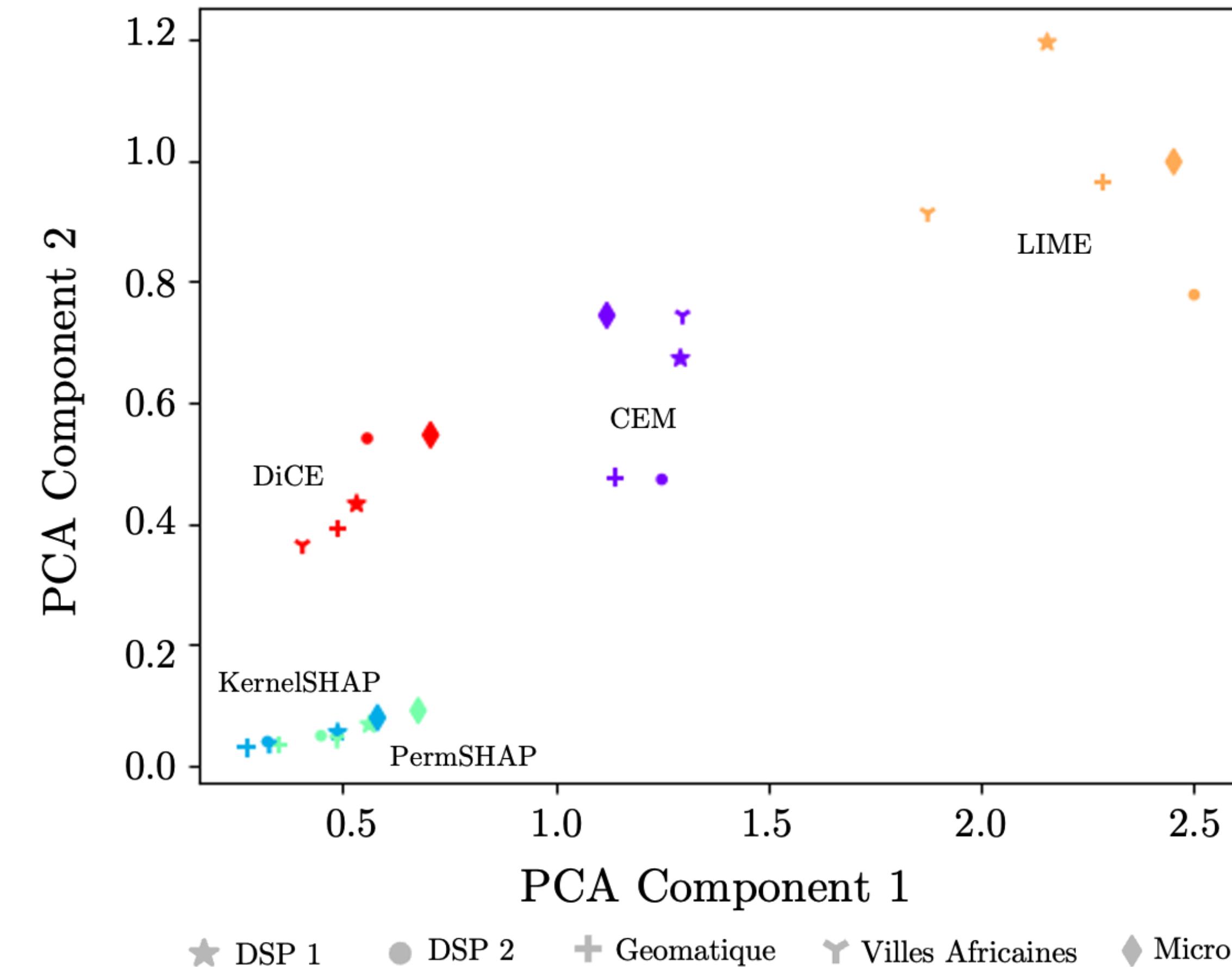
5 Courses

RESULTS

How do explanations
(quantitatively)
compare across
courses?

PCA Analysis

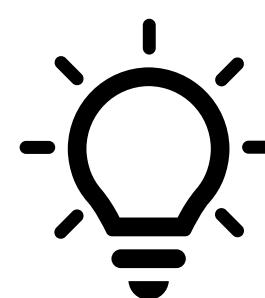
Feature importance
clusters by explainability
method, not by course





XAI methods
systematically
disagree.

How can we build
trust in them?

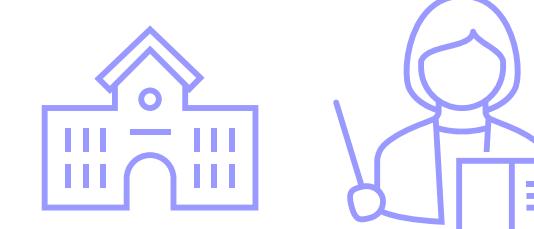


Human expert validation!

Study Participants

RESULTS

Who can we most trust to validate educational explainers?



Over 25 STEM professors
45 minute semi-structured interviews

- diverse in geographic location
- 80.95% identify as male
- age: 36.5 yrs (std: 9 years)
- strong MOOC expertise

Interview Structure

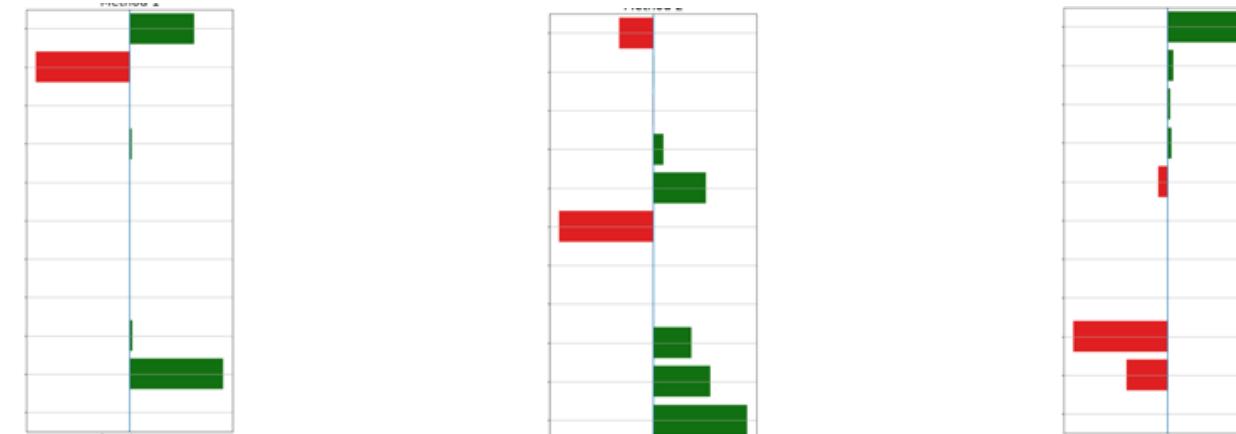
RESULTS

Course pairs that show differences in learning behavior

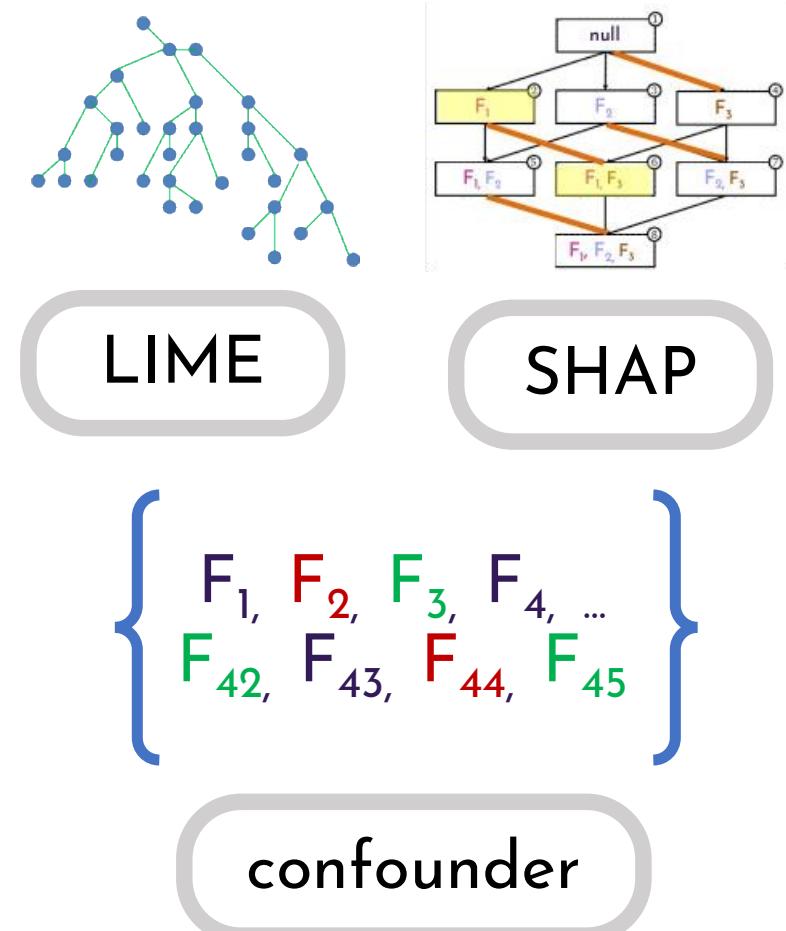
Qualitative Validation



M1 M2 M3



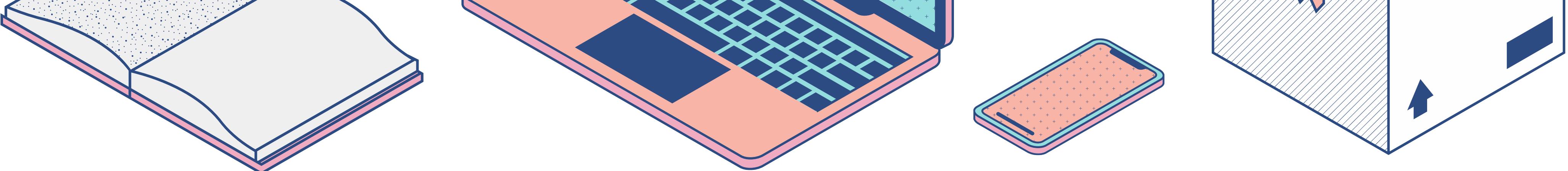
25+ expert interviews
with educators



Expert Trust

RESULTS

- diverse prior perceptions of what factors enable student success and failure
- trust explanations that aligned with their beliefs (46.8%) (which led to significant disagreement across experts)
- only 3 chose the same method for both courses in the pair



Main Takeaways

EVALUATING THE EXPLAINERS
TRUSTING THE EXPLAINERS

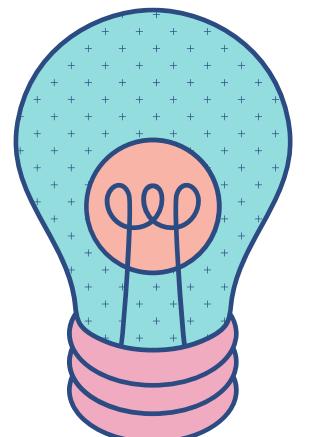
Explainability methods, systematically,
do not agree
on which features are important for predictions

.... and neither do human experts



Where do we go from here?

Interpretable-by- design neural nets!

- 
- 1) **MultiModN** - Multimodal Interpretability
 - 2) **InterpretCC** - Mixture of Experts Interpretability
 - 3) **iLLuMinaTE** - role of LLMs in explainability

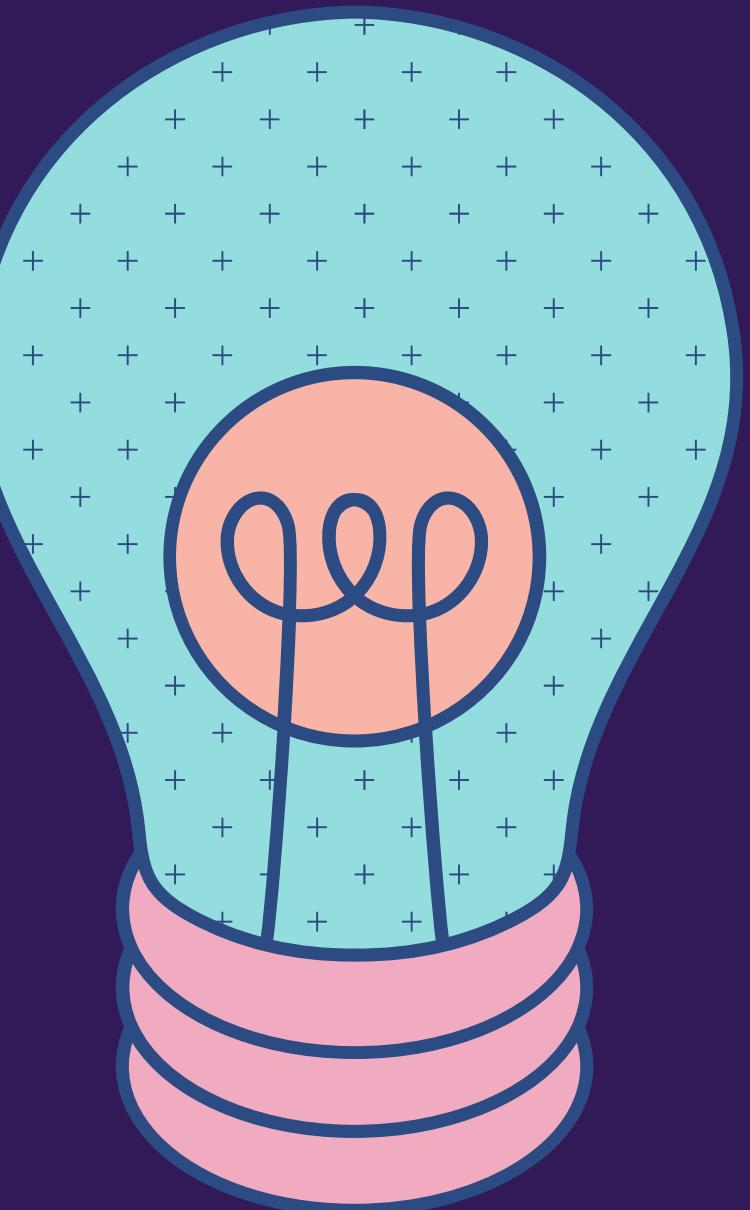
MultiModN

github.com/epfl-iglobalhealth/MultiModN

Multimodal, Multitask, Interpretable Modular Networks



EPFL Yale



Vinitra Swamy ^{1*} Malika Satayeva ^{1*} Jibril Frej ¹



Thierry Bossy ¹ Thijs Vogels ¹ Martin Jaggi ¹ Tanja Käser ^{1*} Mary-Anne Hartley ^{1,2*}

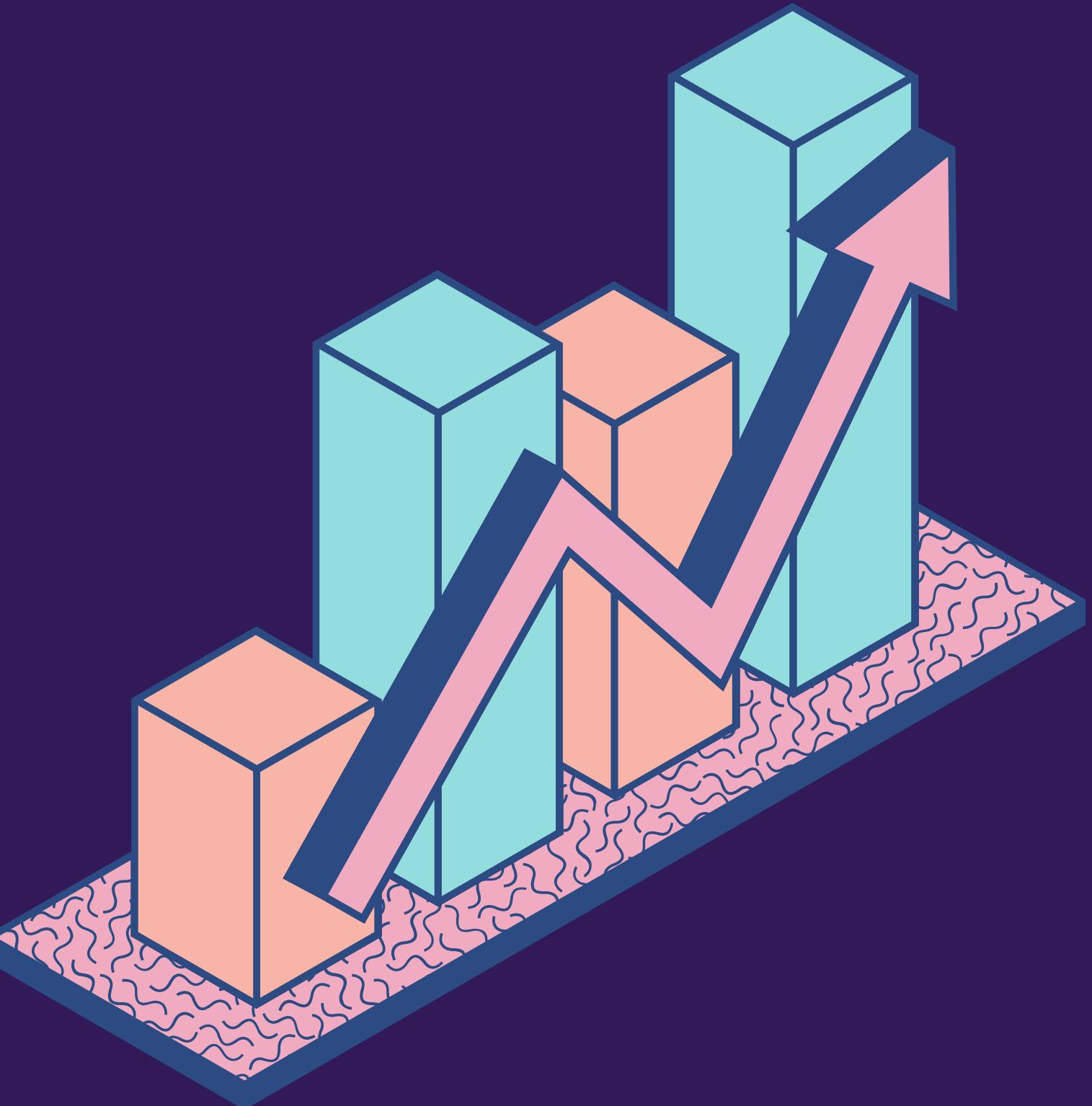
Multimodality

- synergistic predictive potential
- inputs with drastically varying sizes (i.e. images, text, sound)

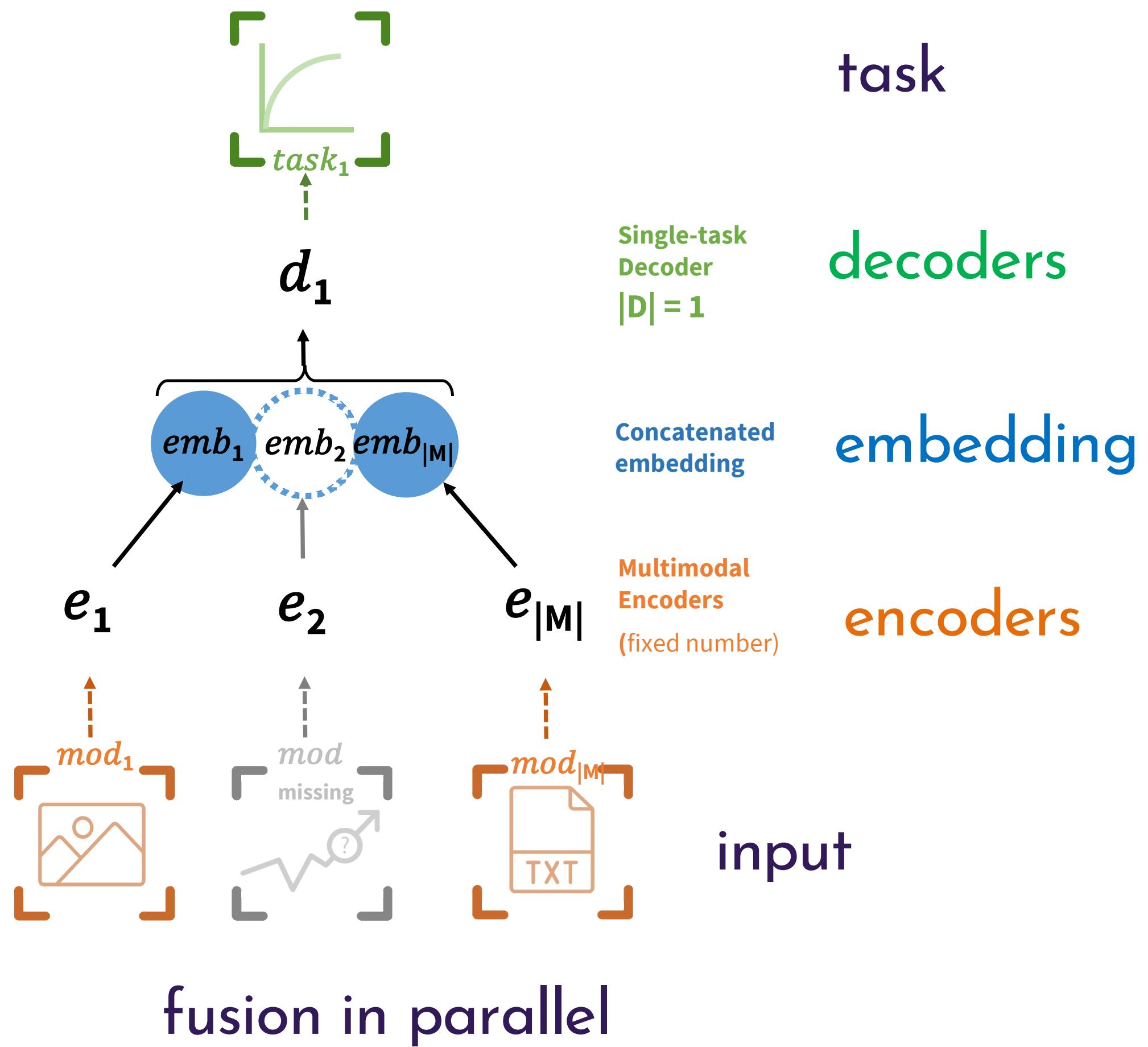
Current Approaches

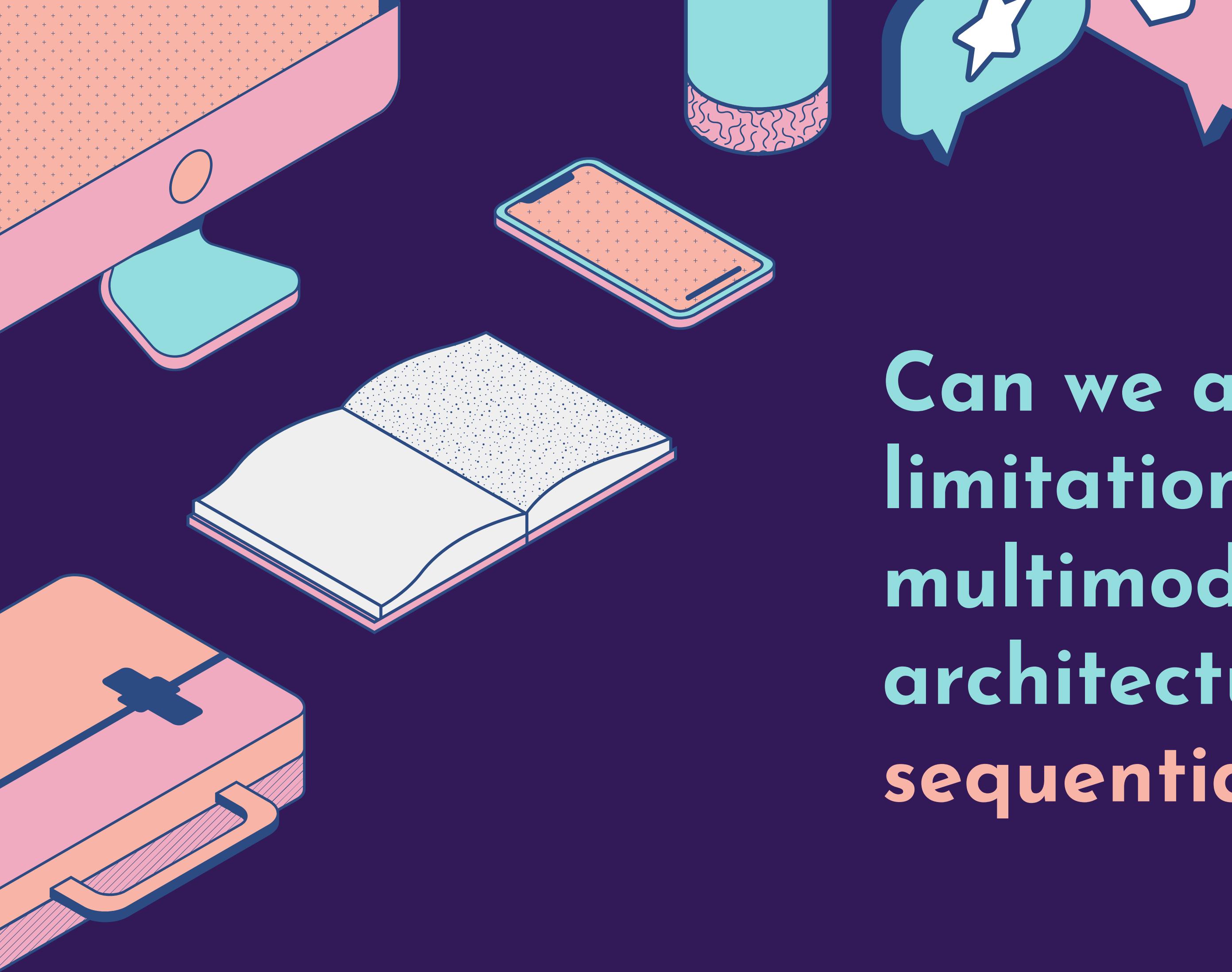
fuse representations in parallel

- limits interpretability
- dependency on modality availability



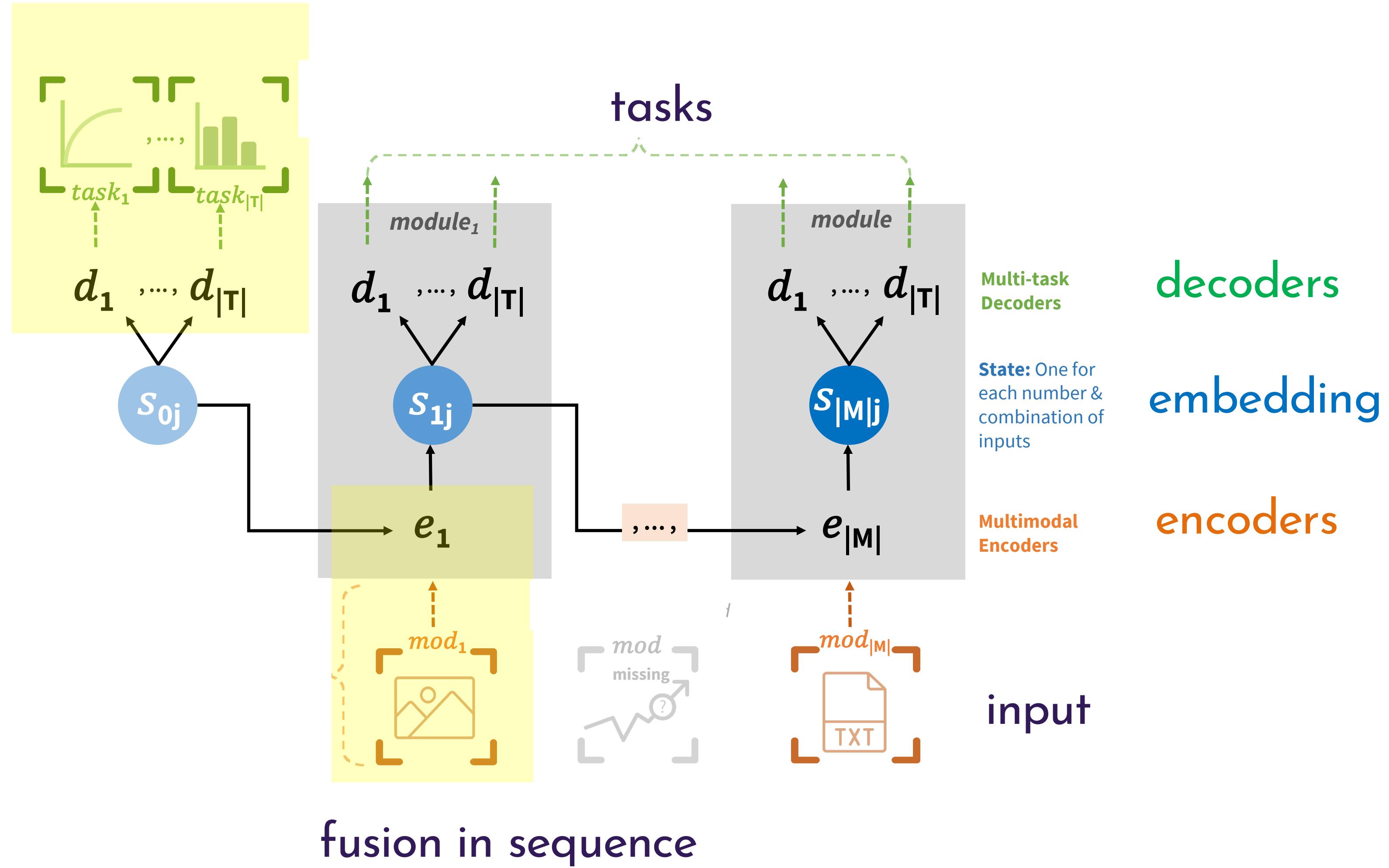
Traditional P-Fusion





Can we address these
limitations of current
multimodal
architectures with
sequential fusion?

MultiModN

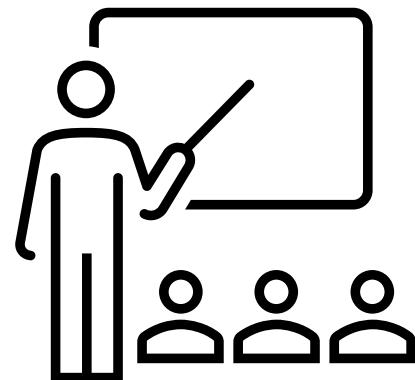
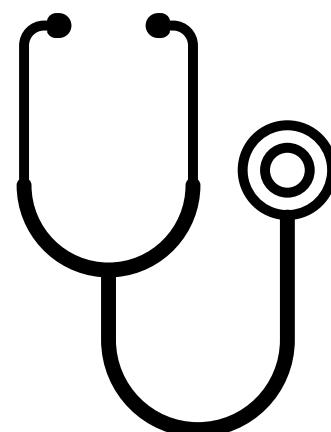


Real World Multimodal Settings

MULTIMODN

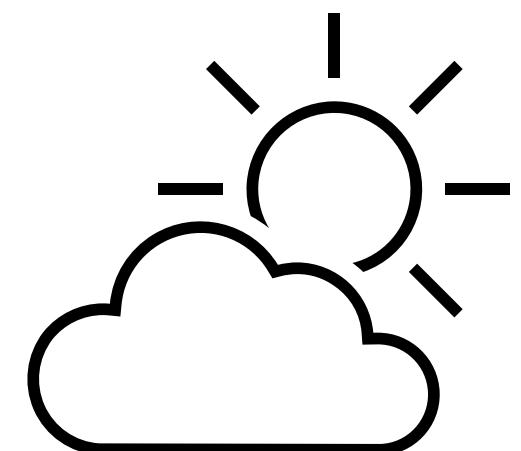
Evaluated on 10 real world tasks across multimodal healthcare (MIMIC IV + CXR-JPG), education (EPFL MOOC), and weather (Weather2k) benchmarks

921 patients in tertiary care in a Boston hospital



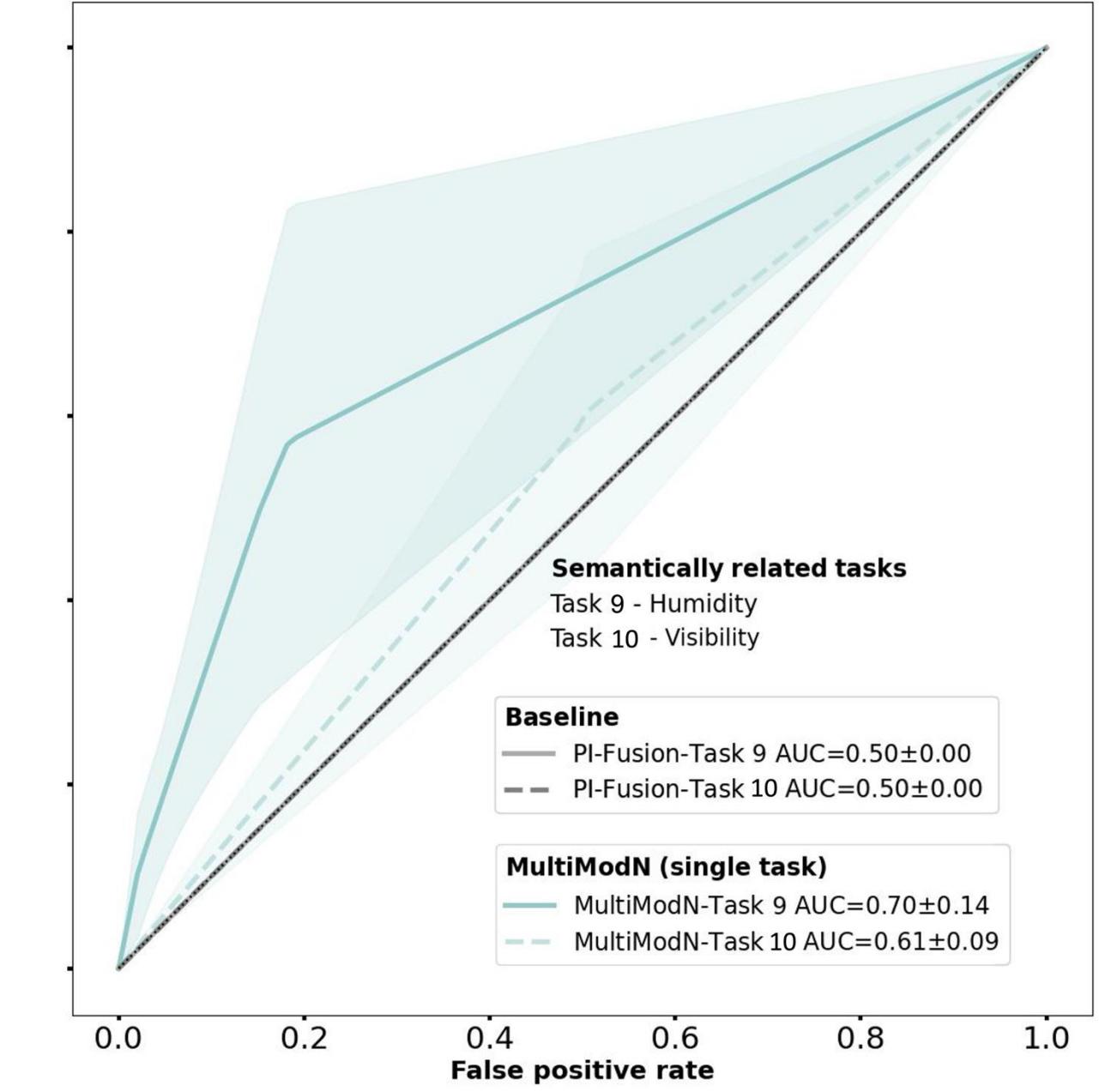
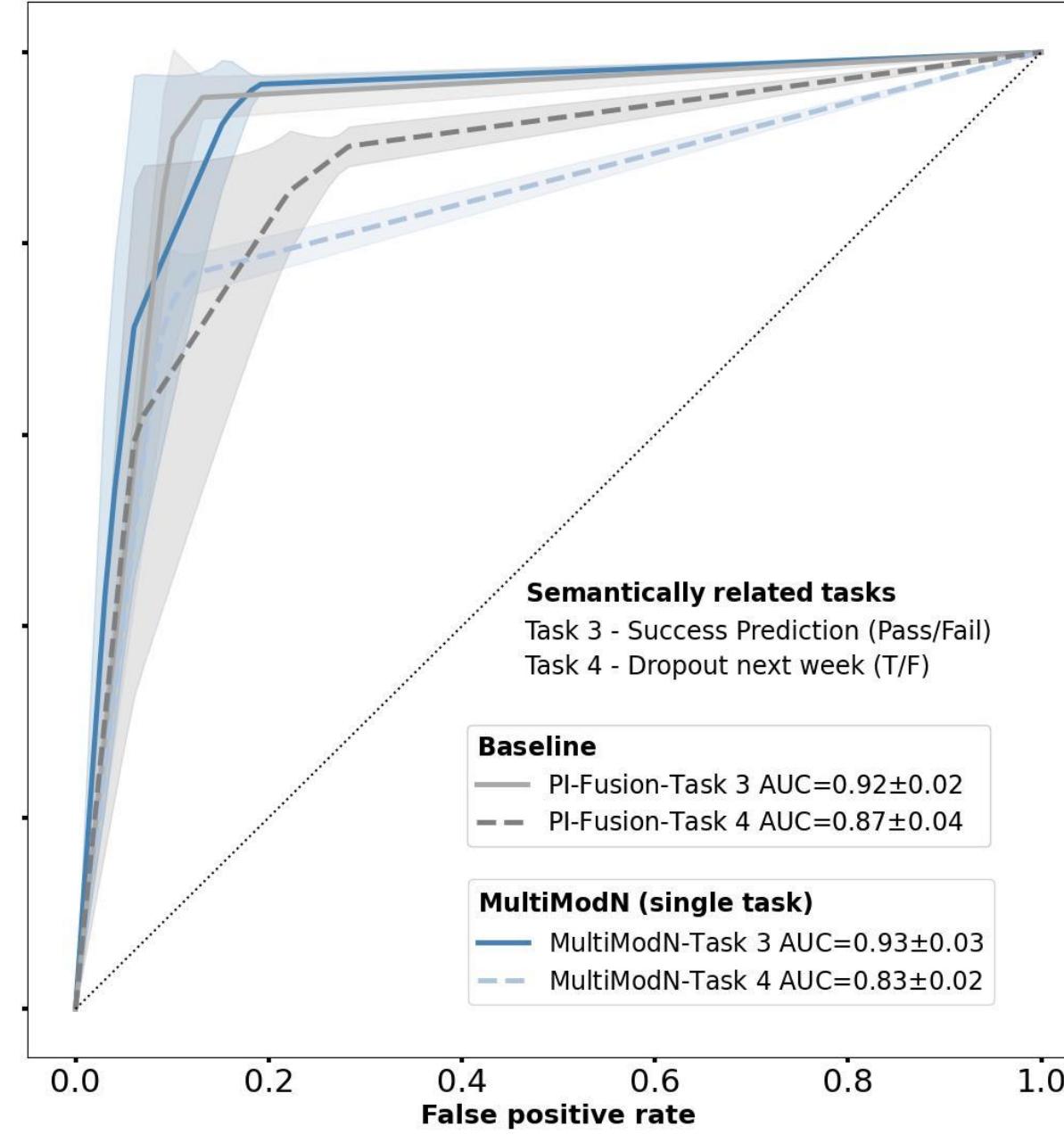
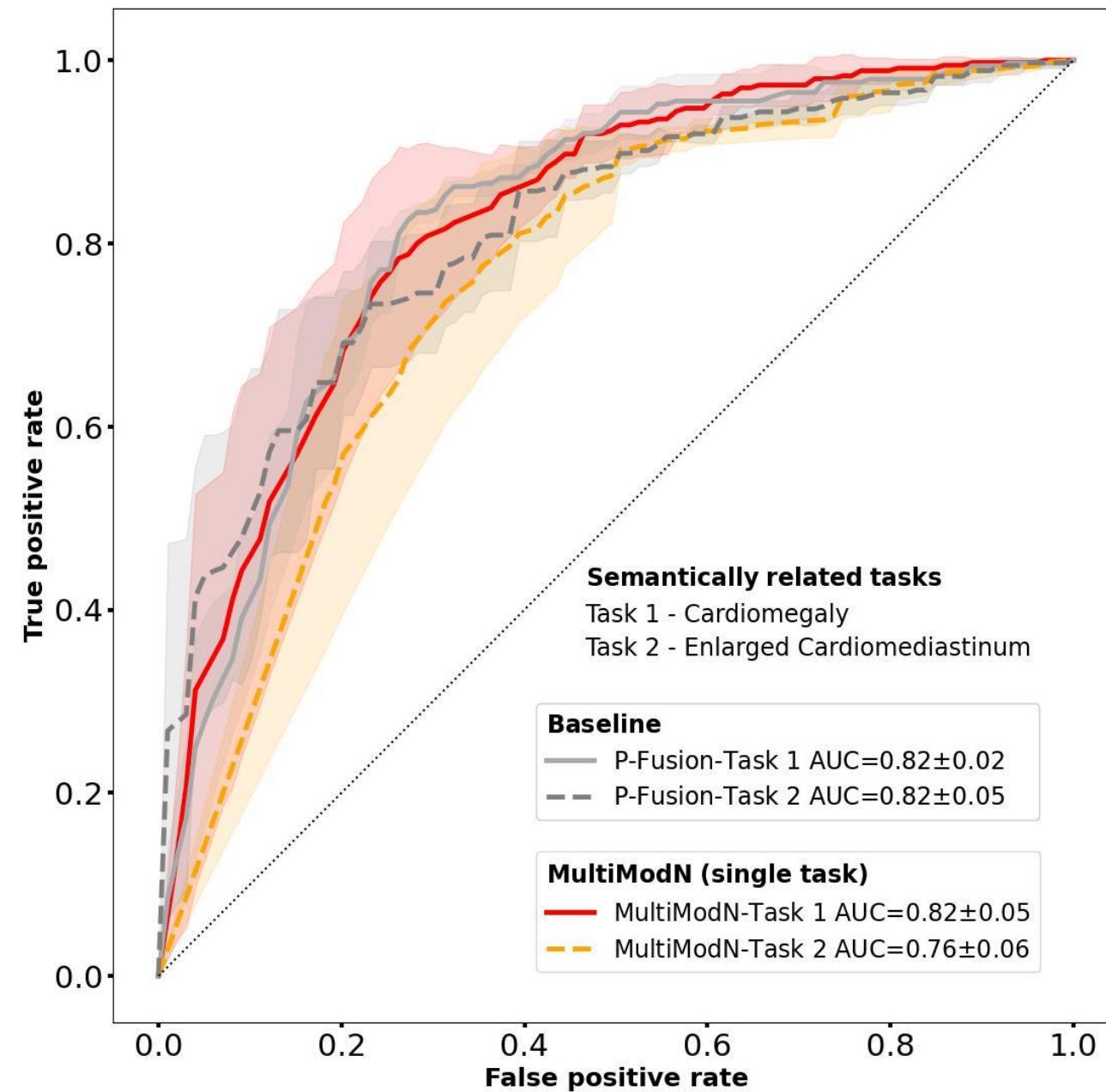
5,611 students in a MOOC with over 1 million interactions

1866 ground weather stations covering 6 million km²



Performance

MULTIMODN VS. P-FUSION



MIMIC

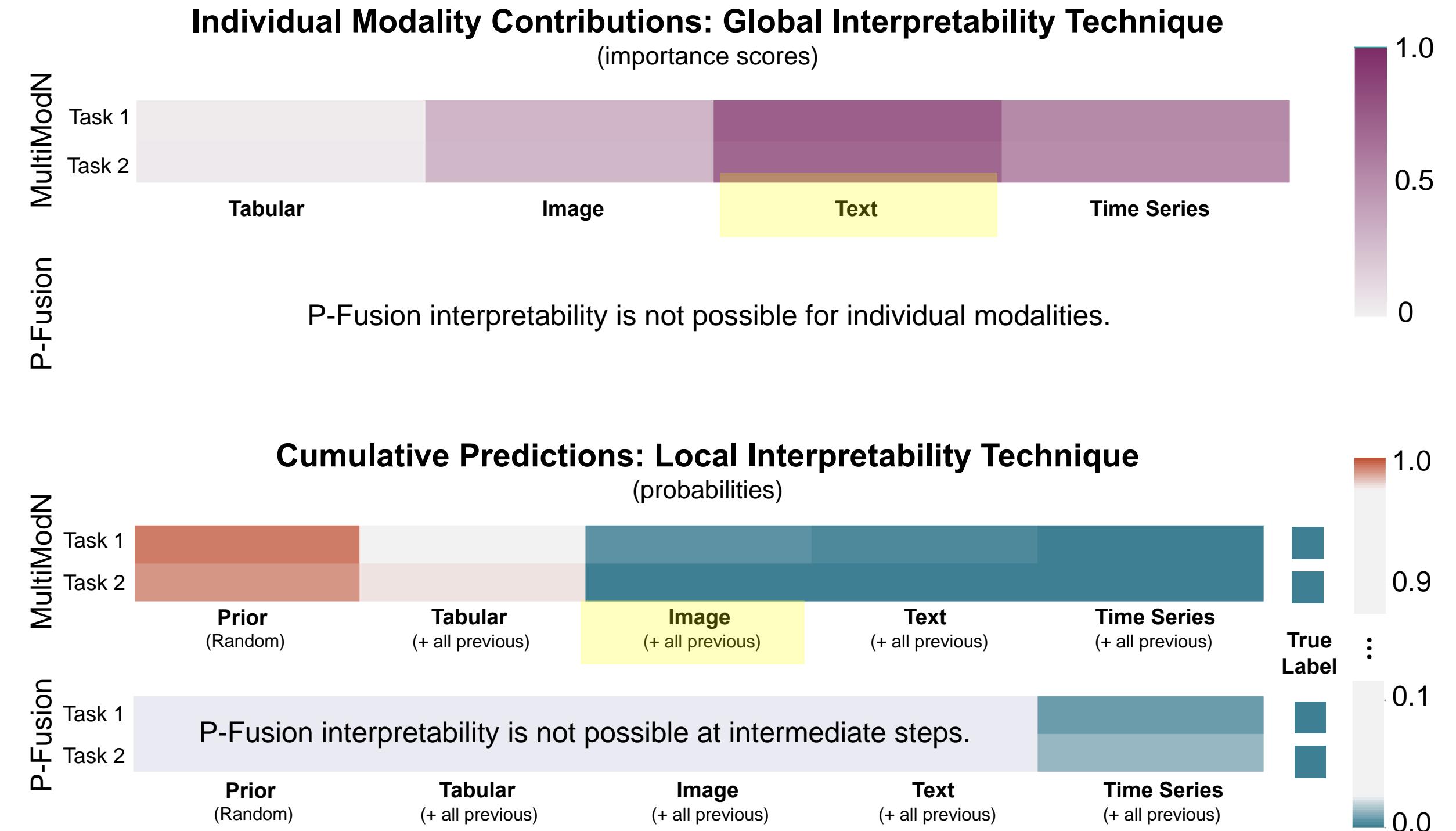
EDU

Weather2k

Interpretable-by-design: MIMIC

MULTIMODN VS. P-FUSION

Which modality
is the most
important for
this task?



Which modality
caused the
decision to flip?

MultiModN has modality-specific global (IMC) and local (CP) model explainability.

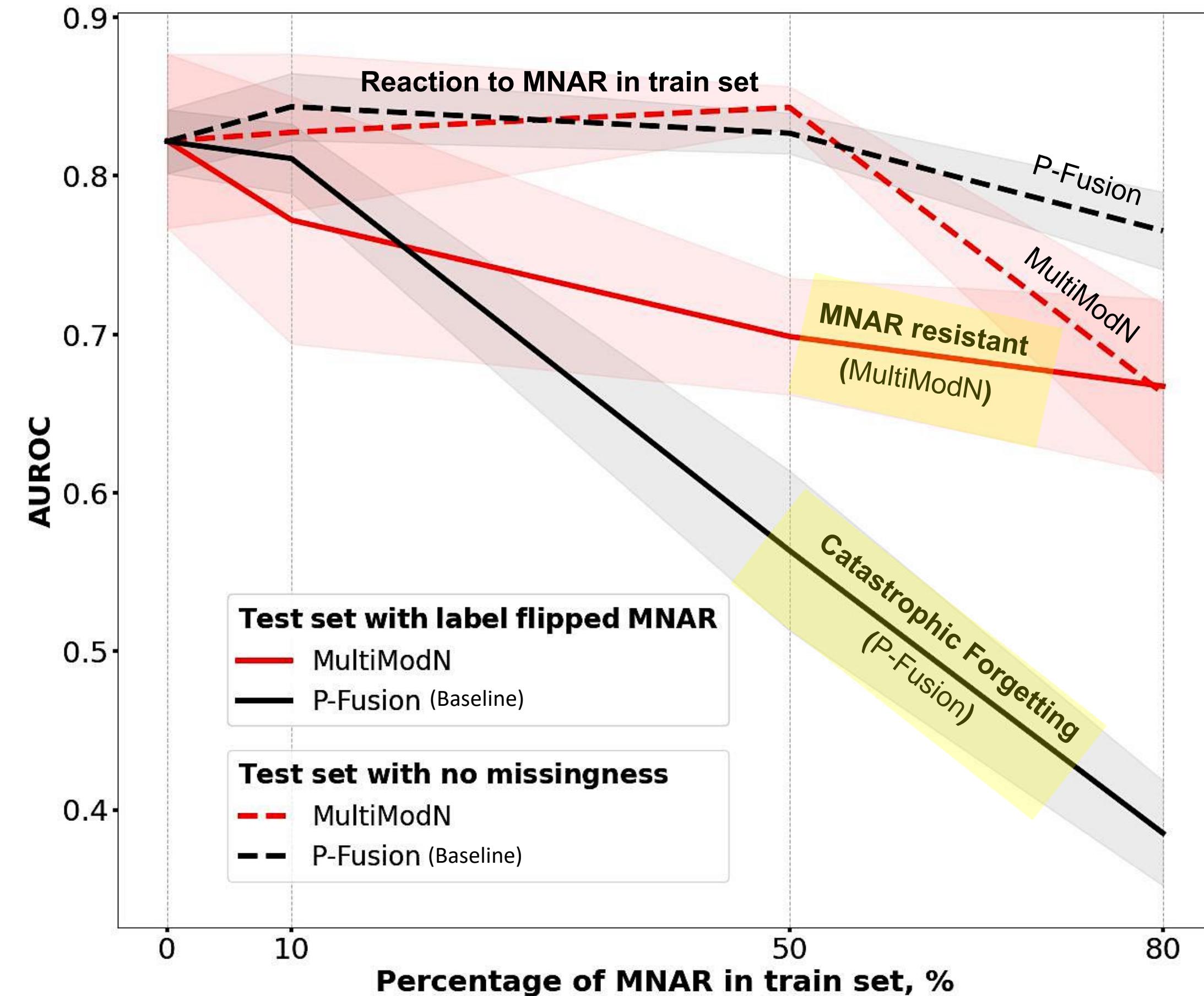
Robust to Missingness

MULTIMODN VS. P-FUSION

MultiModN is robust to bias from missing input modalities.

P-Fusion exhibits catastrophic MNAR failure when missingness patterns change.

Especially relevant in low-resource settings.



From the Lab to Real Life

MULTIMODN



Prof. Annie Hartley
(MD, PhD)



clinical studies using MultiModN for pneumonia and
tuberculosis diagnosis in **low resource settings**

- trained more than 100 doctors to collect multimodal data (images, ultrasound)
- recruiting 1000s of patients in South Africa, Tanzania, Namibia and Benin



Under
Review

InterpretCC

github.com/epfl-ml4ed/InterpretCC

Intrinsic, User-Centric Interpretability through Global MoE

EPFL



Vinitra Swamy



Syrielle Montariol



Julian Blackwell



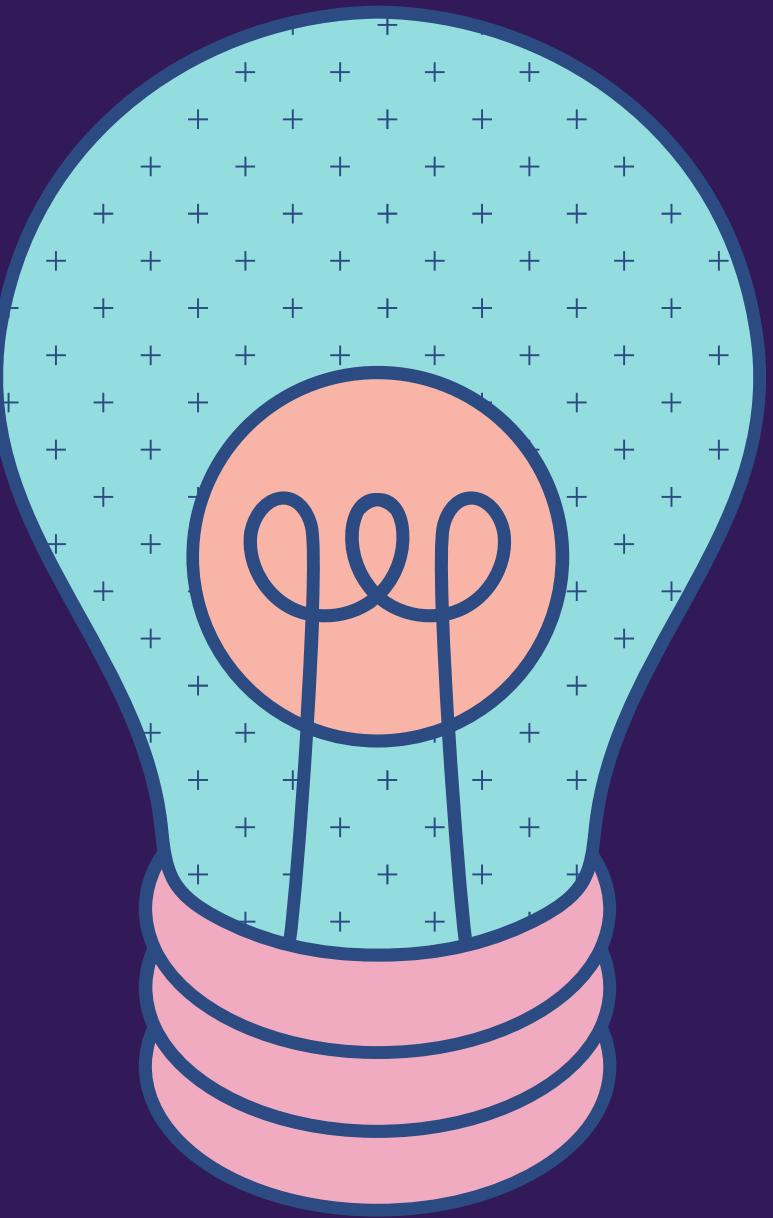
Jibril Frej



Martin Jaggi



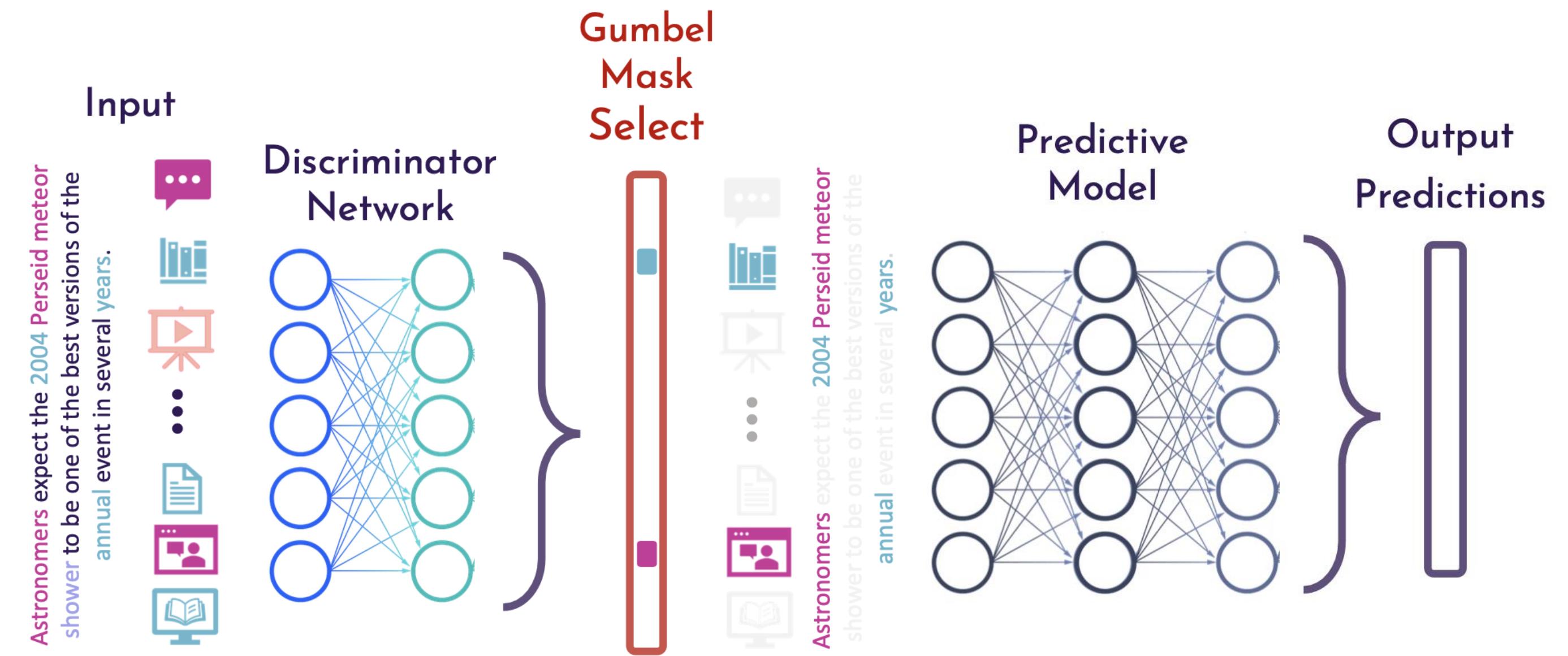
Tanja Käser



InterpretCC: Feature Gating for Interpretability

Swamy, Montariol, Blackwell, Frej,
Jaggi, Käser

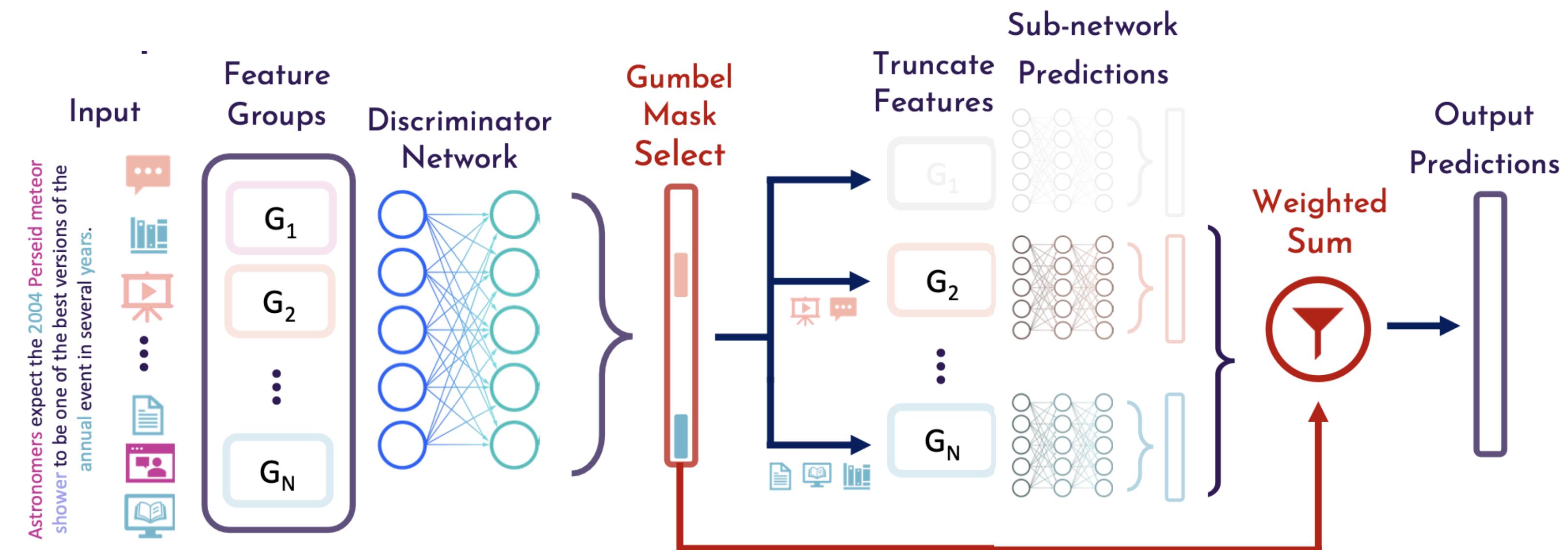
Adaptive Feature Gating - accuracy vs. interpretability tradeoff



InterpretCC: Mixture-of-Experts for Interpretability

Swamy, Montariol Blackwell, Frej,
Jaggi, Käser

Filter the feature space and send relevant parts to relevant experts



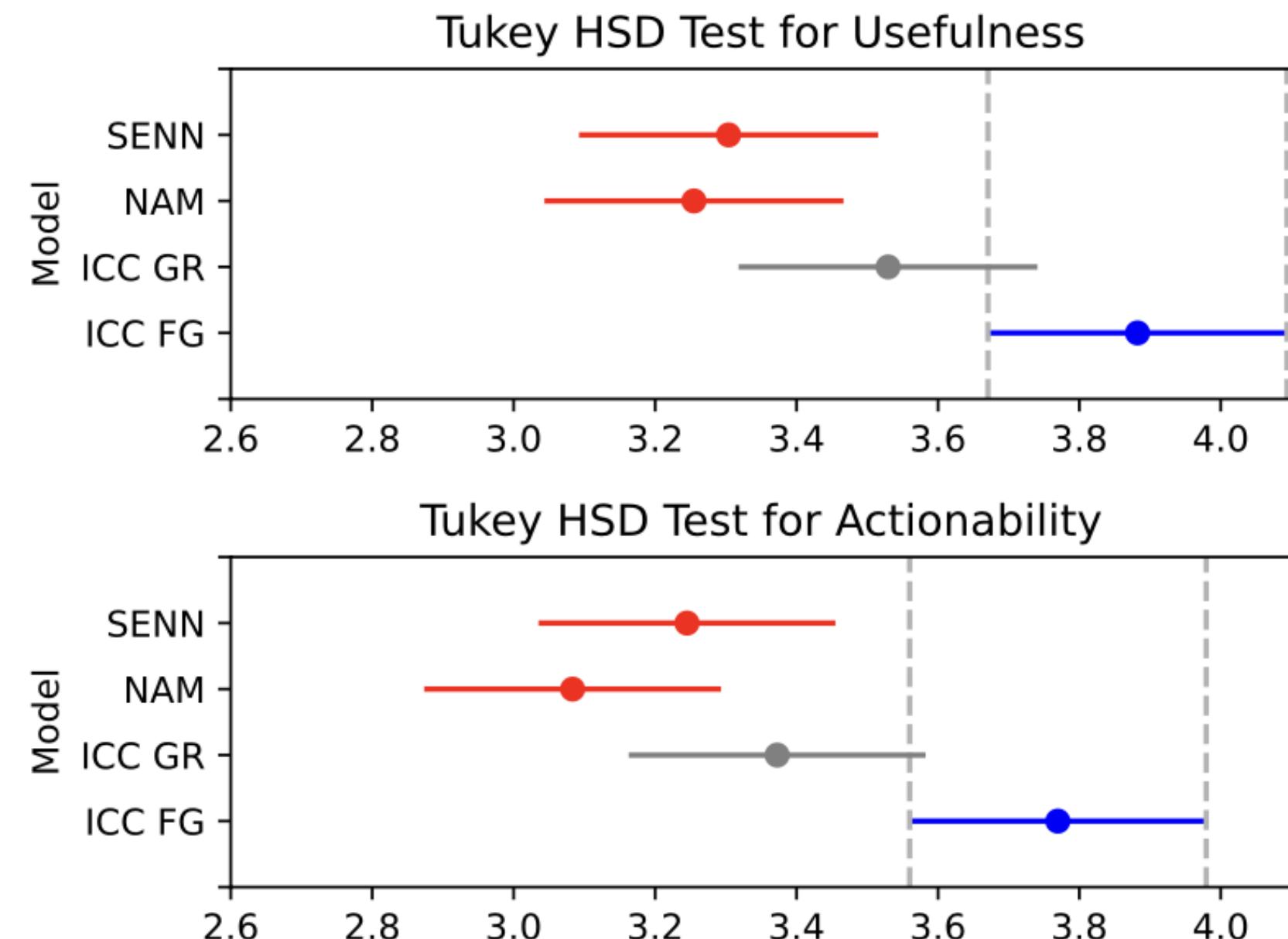
InterpretCC: Results



Swamy, Montariol Blackwell, Frej,
Jaggi, Käser

1) Maintains or improves model **performance** across 8 different benchmark datasets
(text, time-series, tabular)

2) Preferred by 56 teachers over other interpretable-by-design approaches



Under
Review

illuMinaTE

github.com/epfl-ml4ed/illuminate

From Explanations to Action: A Zero-Shot, Theory-Driven
LLM Framework for Student Performance Feedback



Vinitra Swamy*



Davide Romano*



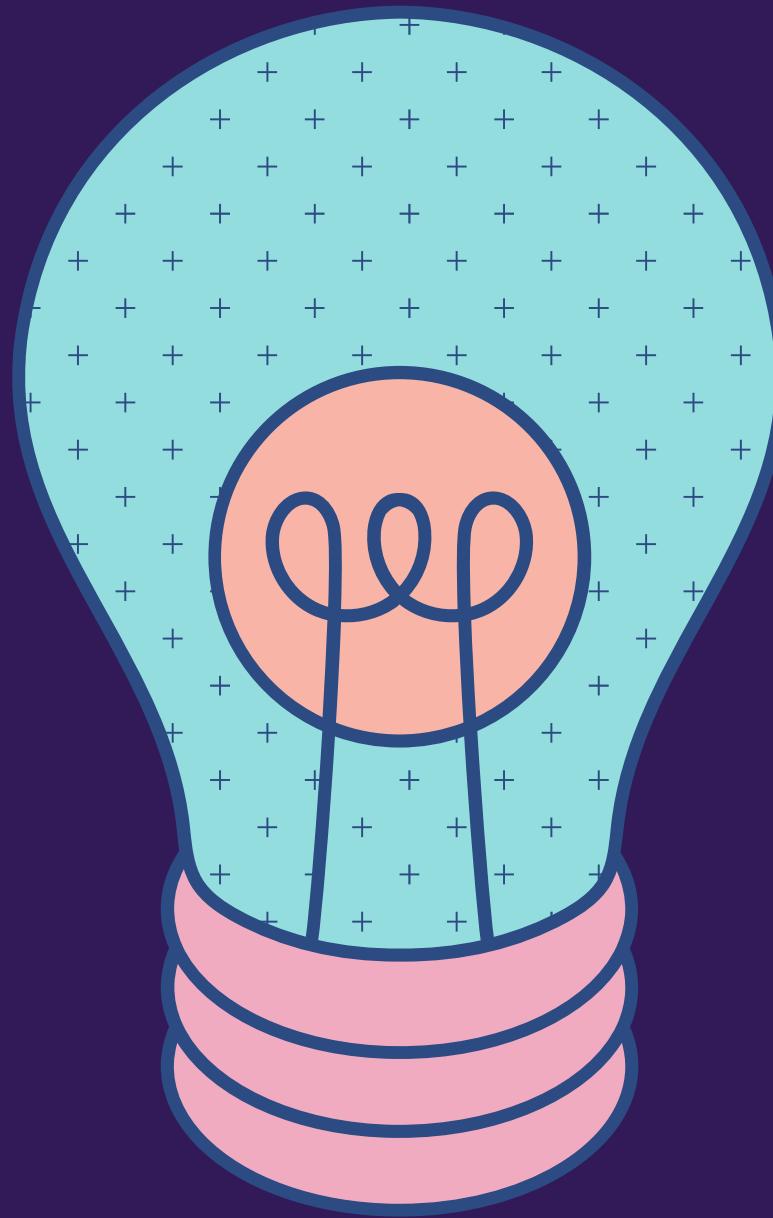
Bhargav Desikan



Oana-Maria Camburu



Tanja Käser

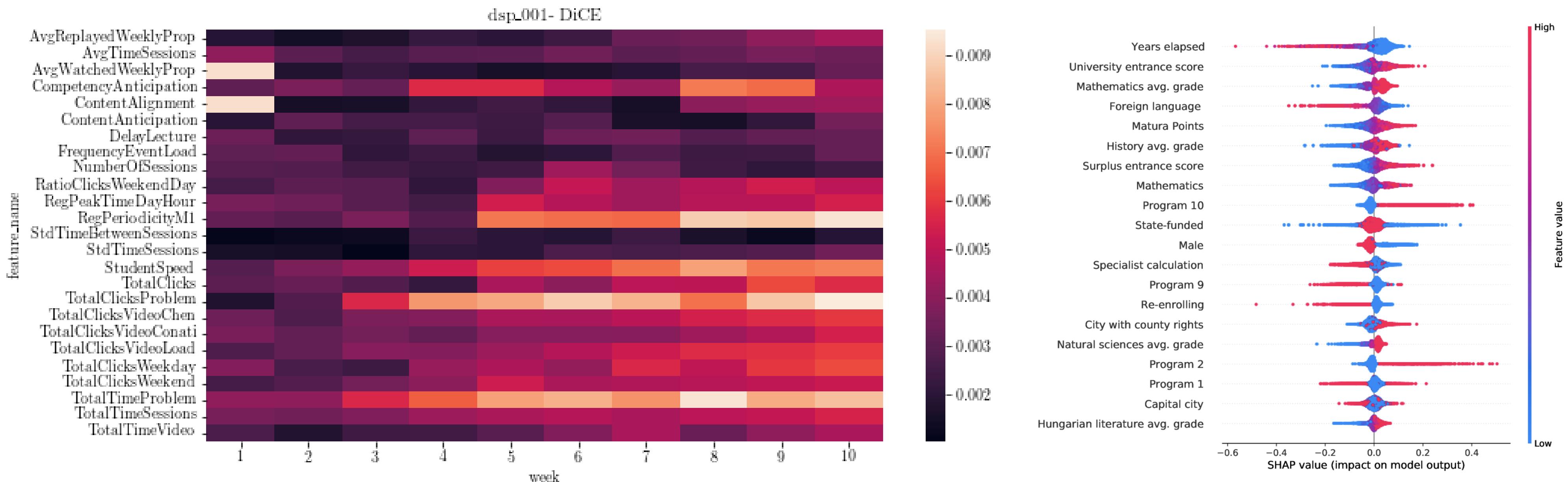
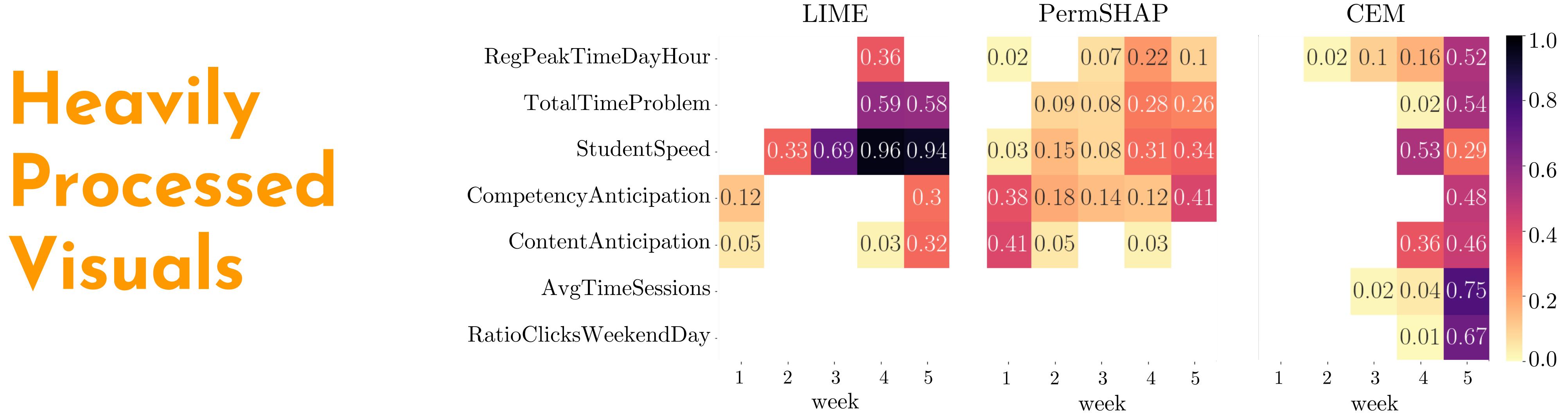


Processed Output of an Explainer

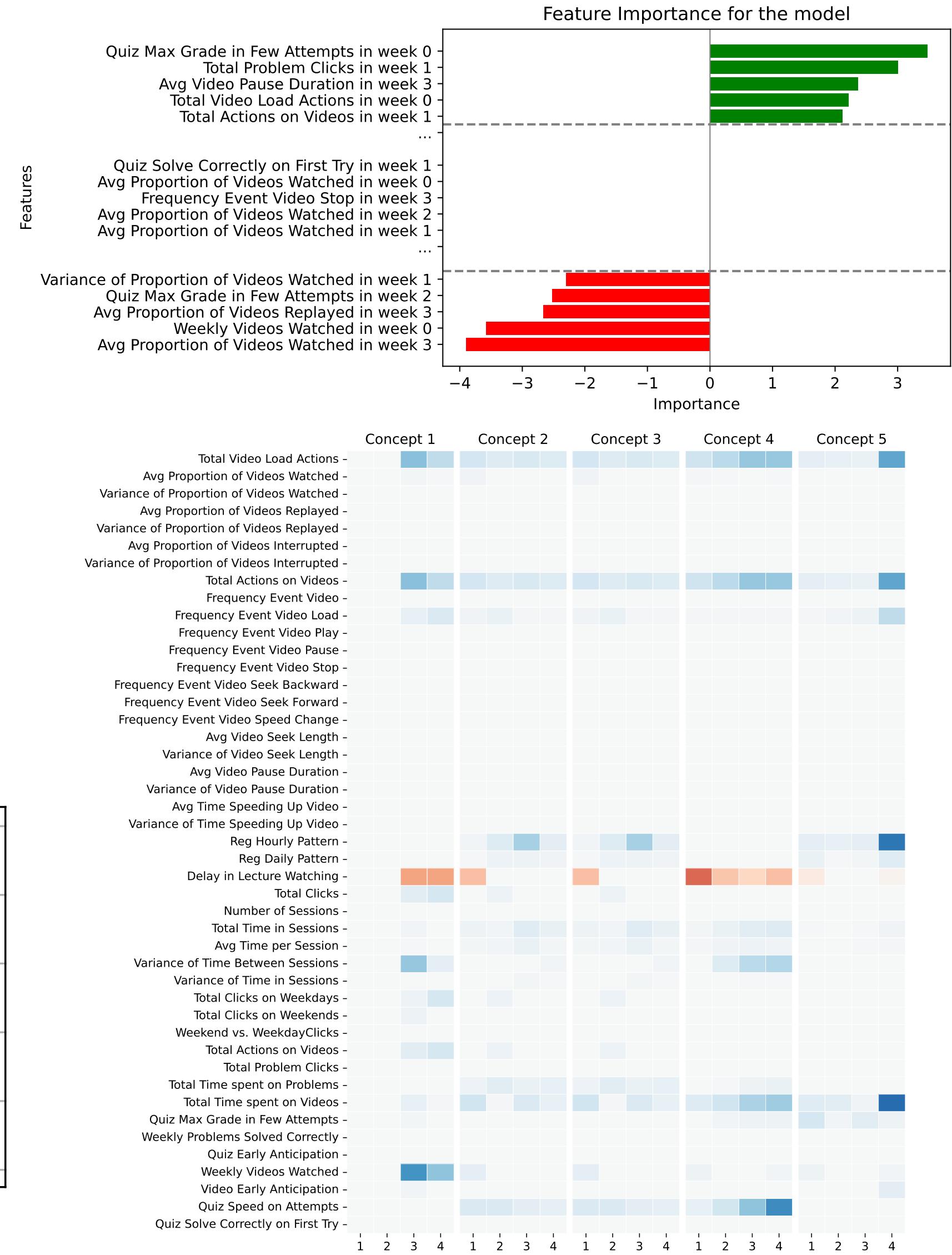
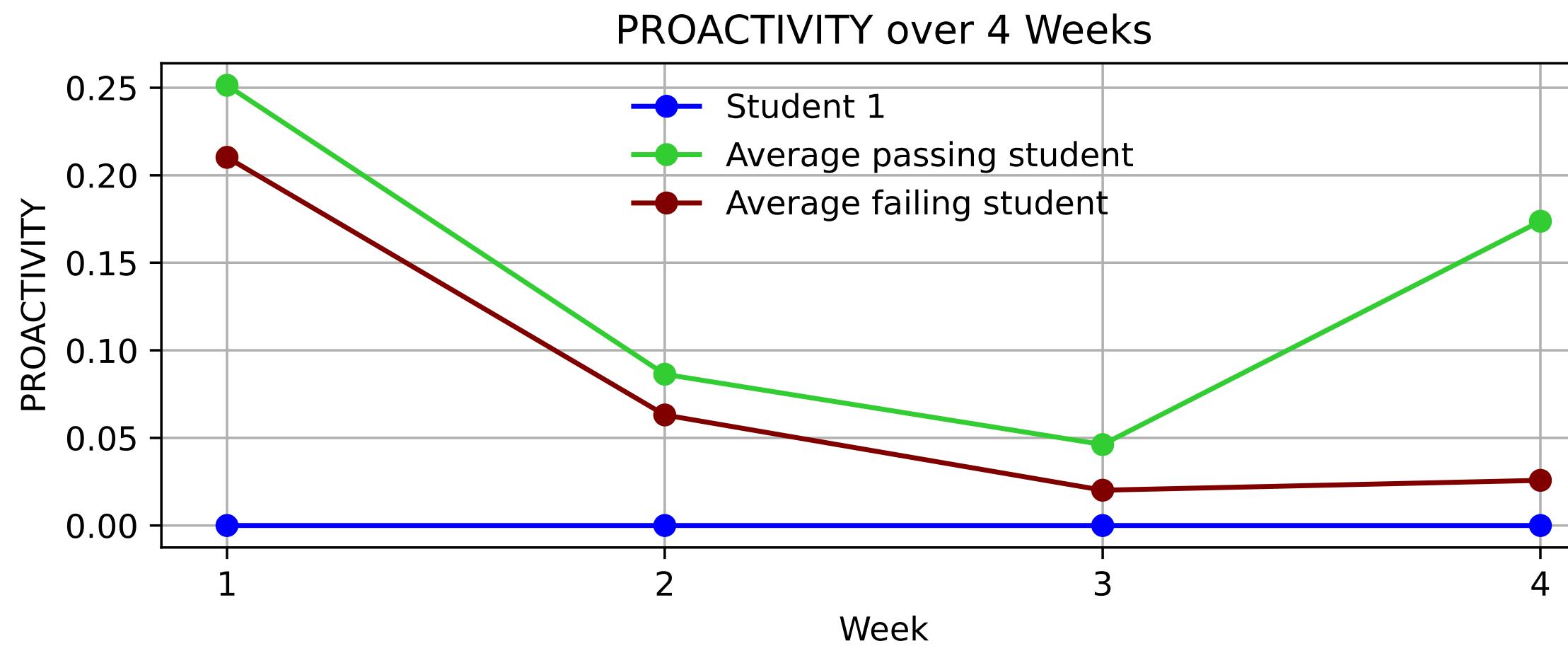
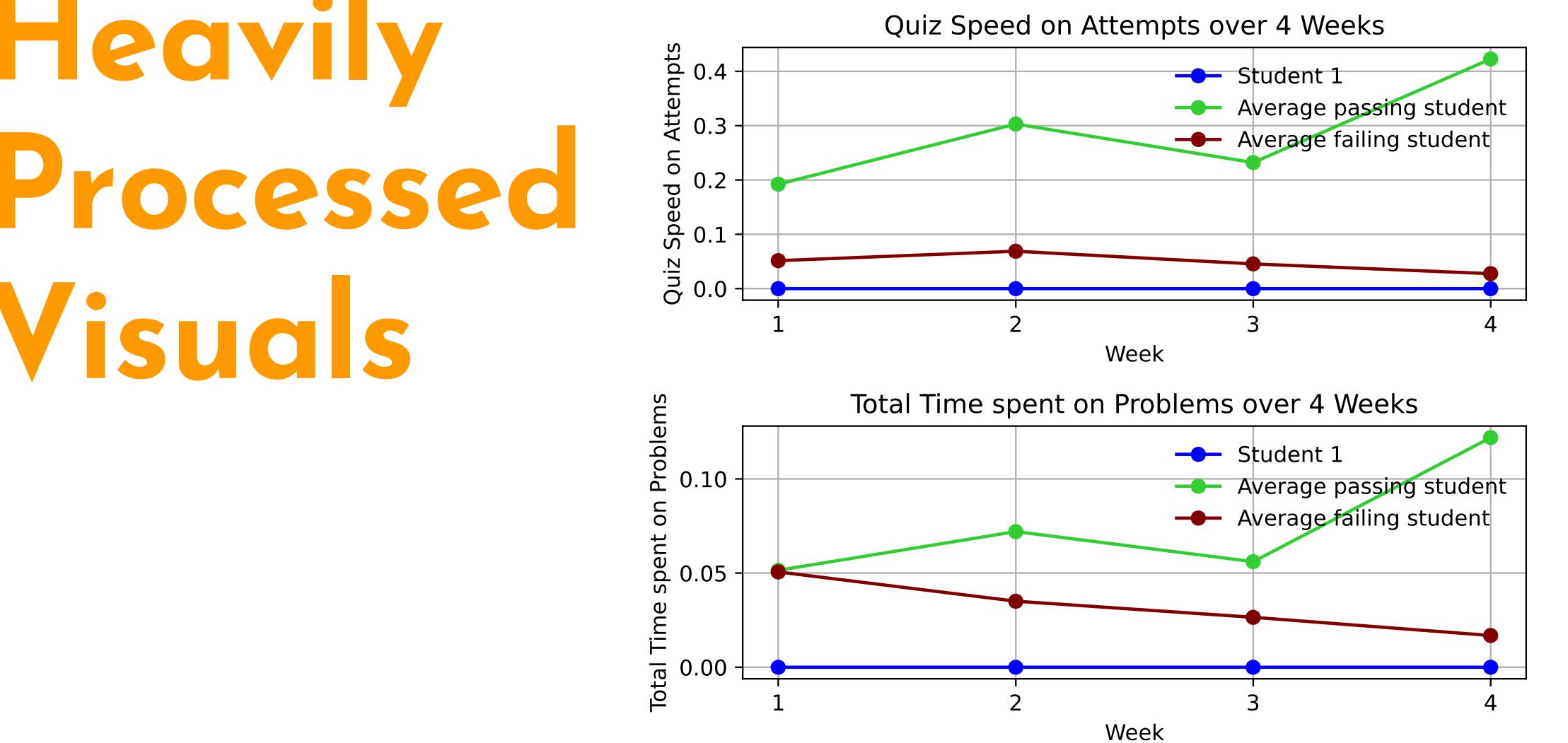
exp_num	total_clicks_InWeek1	number_sessions_InWeek1	time_in_video_sum_InWeek1	time_in_problem_sum_InWeek1
66	1575	-4.56E-09	-2.21E-10	0.074035813 1.3699654119458948e-09
67	5252	-5.60E-10	-2.21E-10	2.1067425364992842e-11 -2.06E-10
68	881	1.7555813192071668e-09	-8.83E-10	0.14653082042868498 -2.65E-10
69	2683	-3.42E-09	-2.21E-10	0.24476348871729897 1.264484322804904e-09
70	16963	-8.50E-10	0	0.045331784 -1.61E-09
71	5931	6.354141102171695e-10	1.3576282653637861e-08	0.17988949383993513 -7.48E-10
72	12145	0	0	0.9422763586044312 0
73	11999	0	0	0.03497038 0
74	1571	6.909329608451031e-09	-8.83E-10	0.13568544287717593 1.6364350777231529e-09
75	2220	-5.60E-10	0	0.057424471615632494 0
76	9184	0	0	0.025690399 0
77	9592	-2.28E-09	-2.21E-10	0 -1.84E-10
78	12309	5.789162378644352e-10	-2.21E-10	0 0 -4.92E-10
79	5394	-2.28E-09	0	0.026910097 0
80	730	-2.28E-09	-2.21E-10	0.029592504314488427 8.244599122332608e-13
81	3318	2.315664951457741e-09	6.7881413268189306e-09	0.014542981 -1.48E-09
82	290	-1.12E-09	0	0 -3.32E-09 0
83	1959	5.789162378644352e-10	0	1.8329978368480937e-09 4.3817753245245505e-10
84	1424	-1.70E-09	-4.42E-10	-4.82E-09 4.4265635254503444e-10
85	12925	0	0	0.021664523 0
86	11139	0	0	0 0 0
87	12603	0	0	0 0 0
88	16987	1.1578324757288705e-09	-2.21E-10	0 0 -1.46E-10
89	10284	0	0	0 0 0
90	1316	1.7555813192071668e-09	-8.83E-10	0.015449408282339427 -9.00E-10
91	2601	1.447290594661088e-10	0	0 0 0
92	2851	-2.84E-09	-4.42E-10	0.037342517841102124 0
93	16707	2.894581189322176e-10	0	0 9.635302451738159e-11
94	11108	0	0	0 0 0
95	5846	1.1578324757288705e-09	-2.21E-10	-3.86E-11 0
96	16538	1.447290594661088e-10	0	0 0 0
97	14230	0	0	0 0 0

10 weeks x 45 fts
= 450 columns

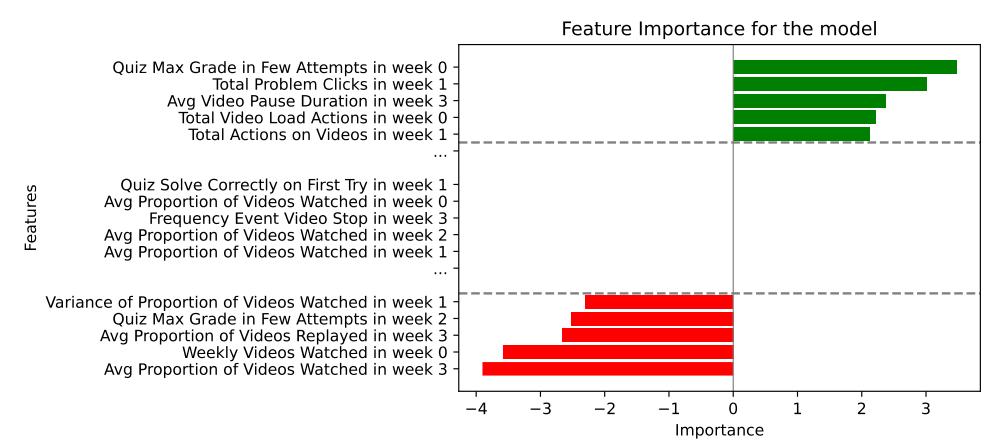
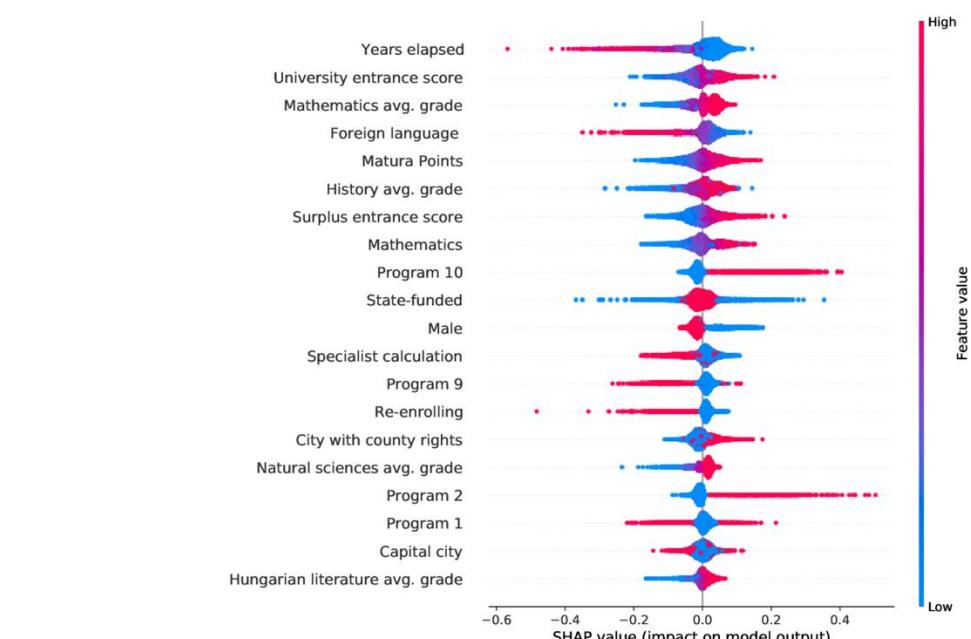
Heavily Processed Visuals



Heavily Processed Visuals

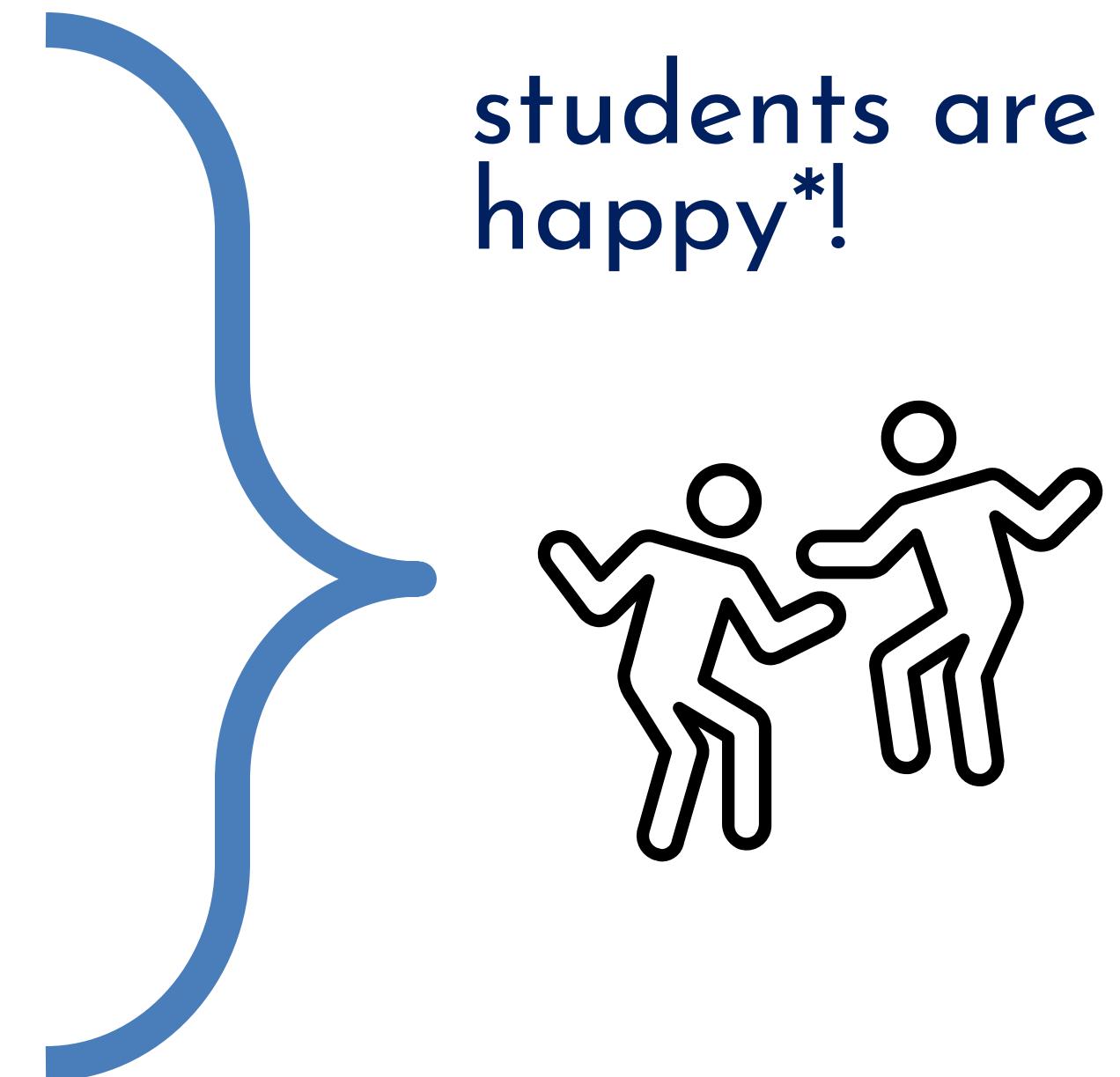
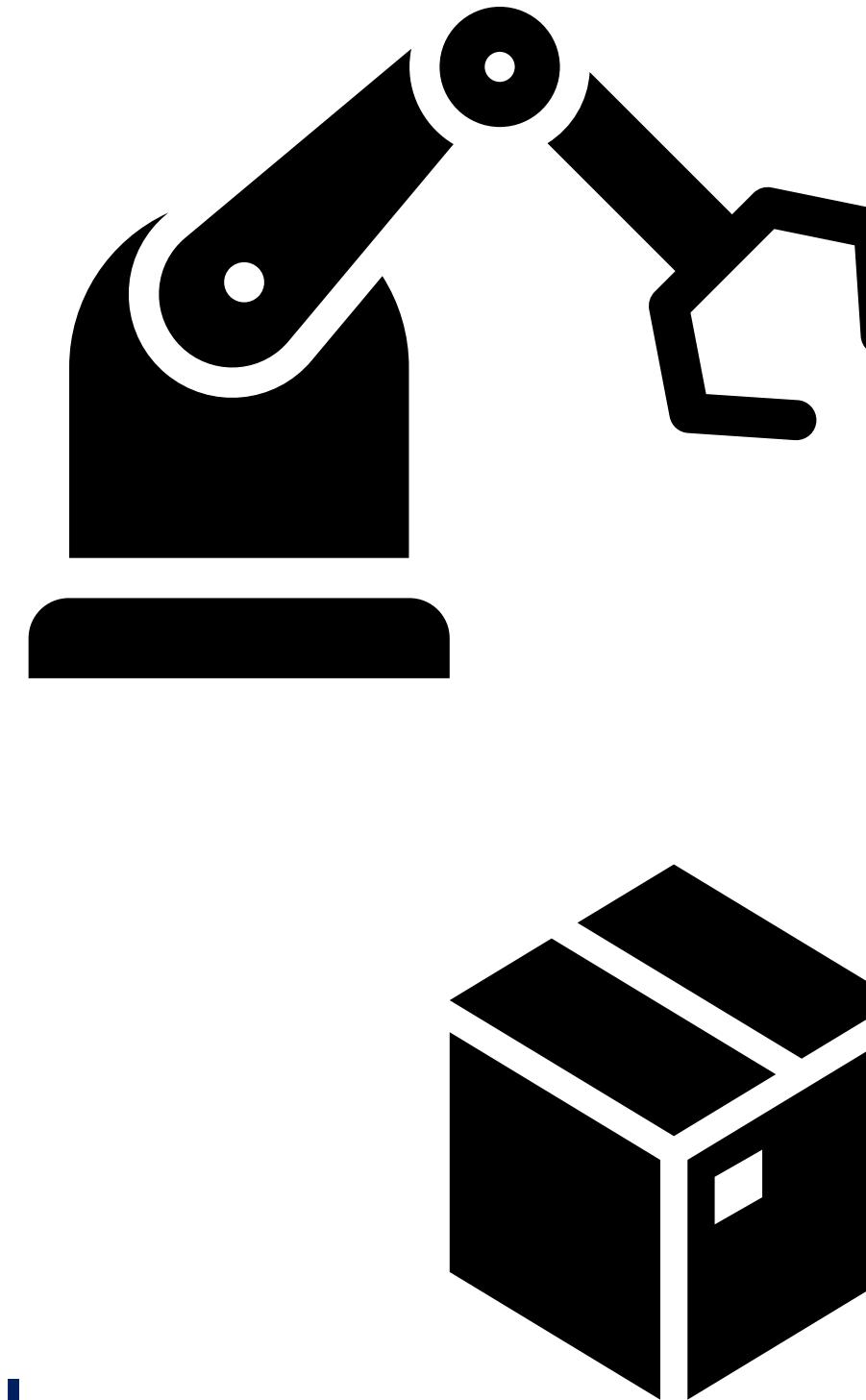


Select what's important

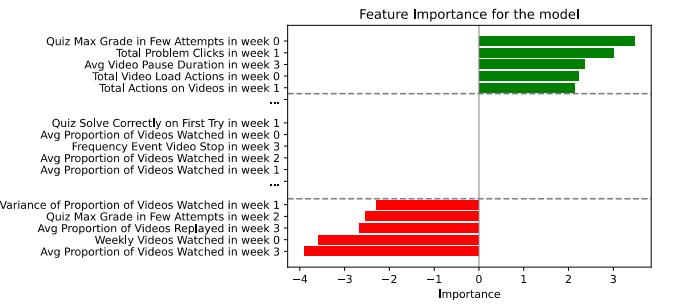
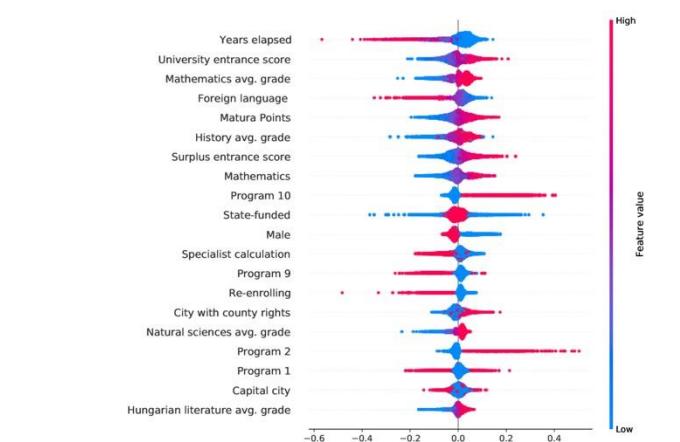


	exp_num	total_clicks_inWeek1	number_sessions_InWeek1	time_in_video_sum_inWeek1	time_in_problem_sum_inWeek1	...
66	1575	-4.56E-09	-2.21E-10	0.074035813	1.3699654119458948e-09	
67	5252	-5.60E-10	-2.21E-10	2.106742364992842e-11	-2.08E-10	
68	88	1.755581319207e-09	-8.83E-10	0.146530204286849	-2.65E-10	
69	2083	-3.42E-09	-2.21E-10	0.244763488712209	2.64484322804904e-09	
70	16963	-8.50E-10	-2.21E-10	0.45331784	1.61E-09	
71	5931	6.354141102171695e-10	1.3576282653637861e-08	0.1798884933993513	-7.48E-10	
72	12145	0	0	0.0422763586044312	0	
73	11999	0	0	0.03497038	0	
74	1571	6.909329608451031e-09	-8.83E-10	0.1356854428771795	1.6364350777231529e-09	
75	2220	-5.60E-10	-2.21E-10	0.0574244711562894	0	
76	9184	0	0	0.025690399	0	
77	9592	-2.28E-09	-2.21E-10	0	-1.84E-10	
78	12309	5.789162378644352e-10	-2.21E-10	0	-4.92E-10	
79	5394	-2.28E-09	0	0.026910097	0	
80	730	-2.28E-09	-2.21E-10	0.029592043144884e-09	6.24459912332608e-13	
81	3318	2.315664951457e-09	5.7881413268189306e-09	0.0174e-08	0.0135e-09	
82	229	-1.12E-09	-2.21E-10	0	0	
83	1959	5.789162378644352e-10	0.8329973868480937e-09	0.3817753245245505e-10	0	
84	1424	-1.70E-09	-4.42E-10	0	-4.82E-09	
85	12925	0	0	0.021664523	0	
86	11139	0	0	0	0	
87	12603	0	0	0	0	
88	16907	1.1578324757288705e-09	-2.21E-10	0	-1.46E-10	
89	10284	0	0	0	0	
90	1316	1.7555813192071688e-09	-8.83E-10	0.01544940828339427	-9.00E-10	
91	2601	1.447295054661088e-10	0	0	0	
92	2851	-2.84E-09	-4.42E-10	0.037342517841102124	0	
93	16707	2.89458118932176e-10	0	0.635302451738159e-11	0	
94	11139	0	0	0	0	
95	5046	1.1578324757288705e-09	-2.21E-10	-3.88E-11	0	
96	16538	1.447295054661088e-10	0	0	0	
97	14230	0	0	0	0	

Package it up nicely

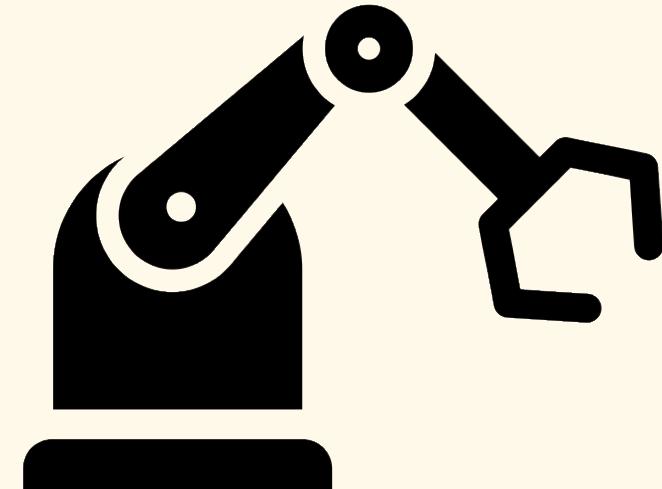


*with their improved learning outcomes

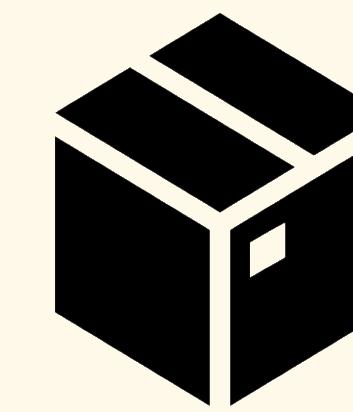


#	avg_num_total_clicks_inWeek1	number_sessions_inWeek1	one_h_video_avg_interest	time_n_problem_avg_interest
65	1129	-4.66e-10	-2.21e-10	309965119469484e-09
67	5202	-5.60e-10	-2.21e-10	30974356989242e-11
68	1000	-4.66e-10	-2.21e-10	30966119469484e-09
69	2802	-3.42e-09	-2.21e-10	26484222004004e-10
70	1665	5.765681179021e-10	5.34750887172942e-09	0.04301784
71	5.3041411017789e-10	5.37620520537679e-09	5.78853935049277e-09	-1.81e-09
72	1746	0	0	5.42273563604412e-09
73	1571	5.90329050481031e-09	5.15685428777198e-09	5.6364305777723152e-09
74	2205	-5.80e-09	-2.21e-10	0
75	9950	0	0	0.02987039e-09
76	1138	-2.28e-09	-2.21e-10	-4.98e-10
77	8384	-2.28e-09	0	0.02987039e-09
78	1024	-2.28e-09	-2.21e-10	-1.84e-10
79	1158	-2.28e-09	0	0.02987039e-09
80	1158	-1.12e-09	-3.32e-10	0
81	2318	5.2156685147741e-09	5.78814132693909e-09	5.15462091e-09
82	296	0	0	5.83209783604409e-10
83	1138	0	0	5.1775242535050e-10
84	1424	-1.70e-09	-4.42e-10	-4.82e-09
85	1138	0	0	0.021964622e-09
86	1138	0	0	0
87	12603	0	0	0
88	1024	0	0	-1.46e-09
89	1024	0	0	-6.00e-10
90	1158	0	0	0
91	2601	4.4729050461088e-10	4.63e-10	4.6154494682233947e-09
92	2601	5.84651893277124e-10	4.42e-10	5.837342517841190124e-09
93	1158	5.84651893277124e-10	0	5.63530245178150e-11
94	1158	0	0	0
95	1638	4.4729050461088e-10	-2.21e-10	-3.88e-11
96	14202	0	0	0

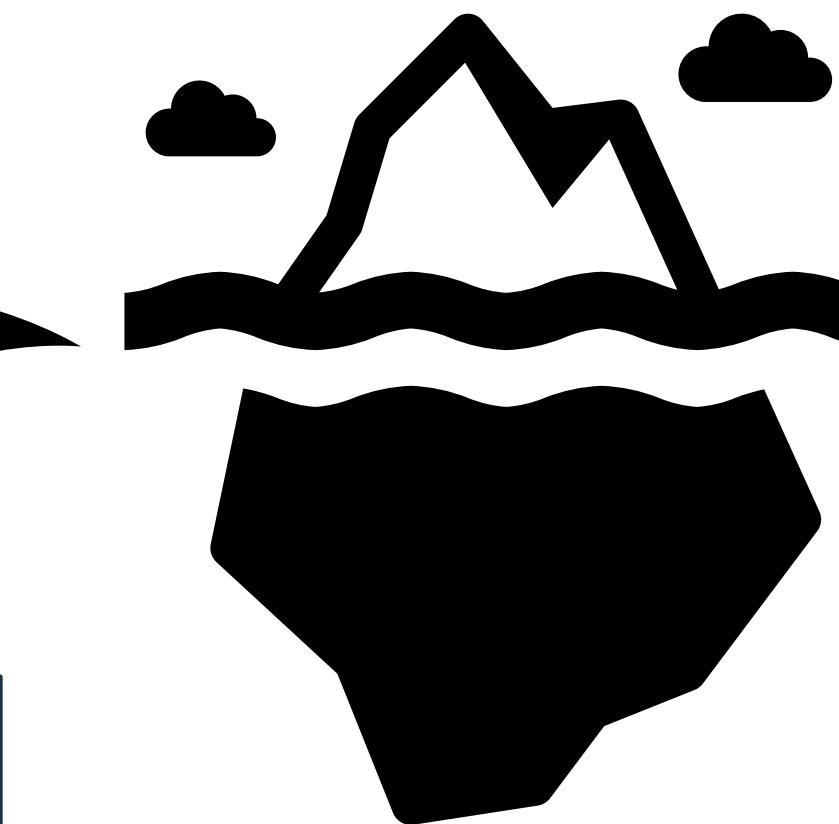
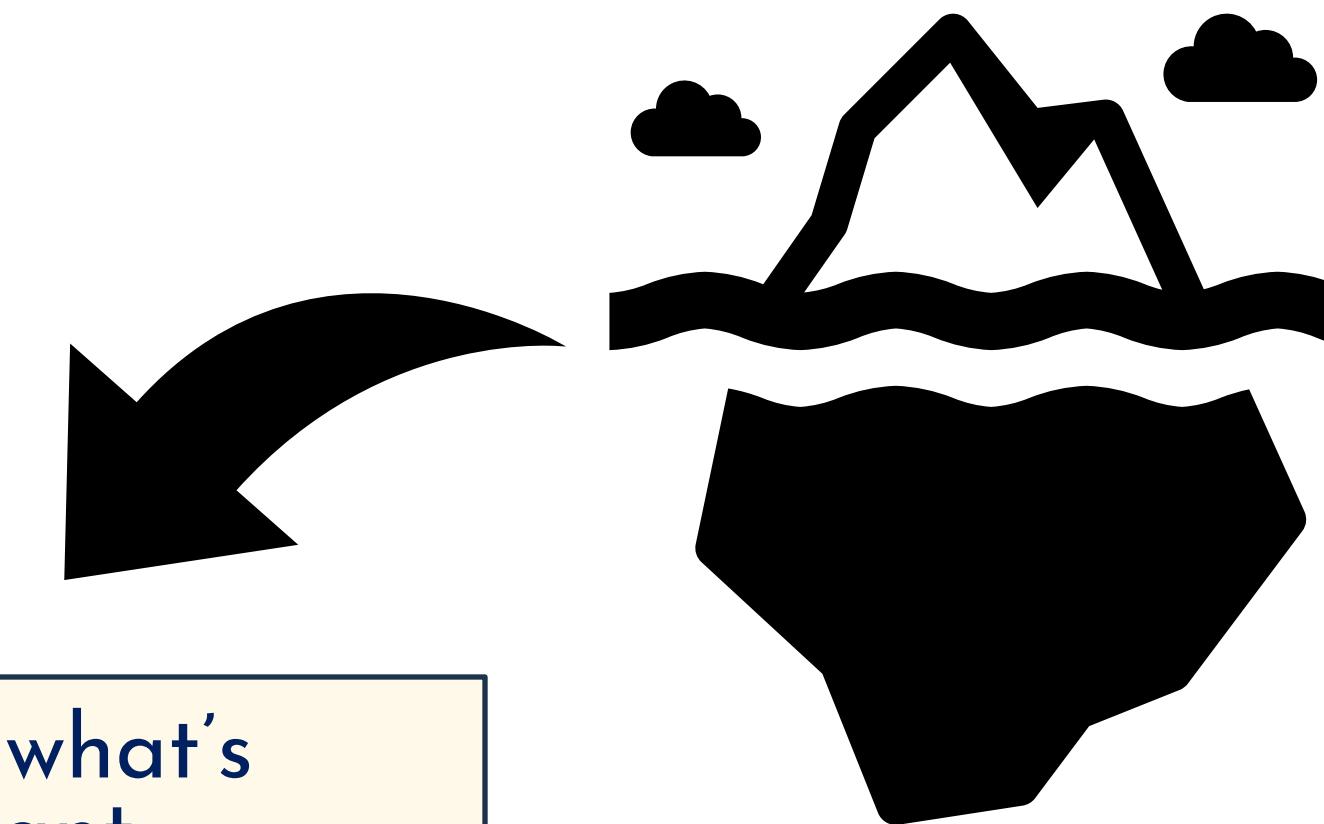
Select what's important



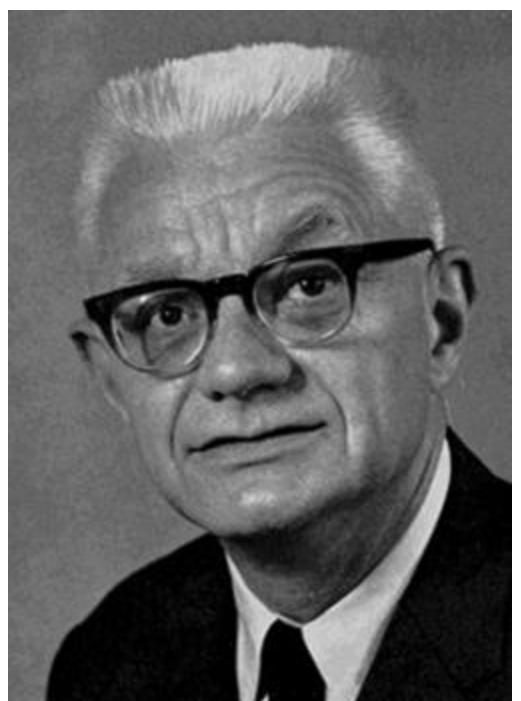
Package it up nicely



students are happy!



Leverage theories of explanation from philosophy and social sciences to define prompting strategies



Relevance
Selection

Hempel

DN model



C.S. Peirce
Abductive reasoning

Statistical
Relevance

Abnormal
Conditions



Malle
Social attribution
model

Necessity
Robustness



P. Lipton
Contrastive
explanation

RaR +
Contrastive

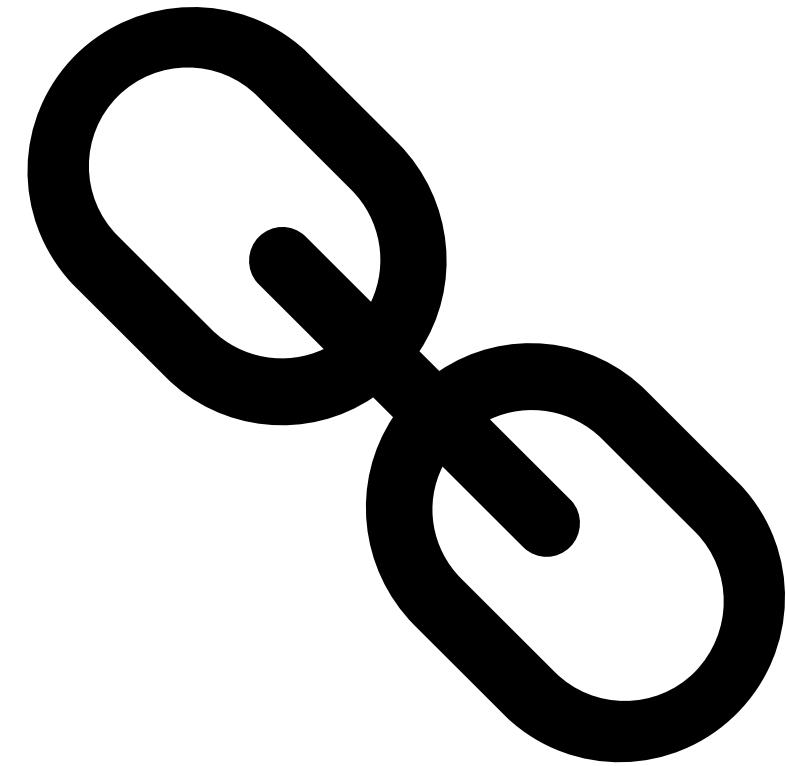


Lombrozo
Cognitive processes of
explanation

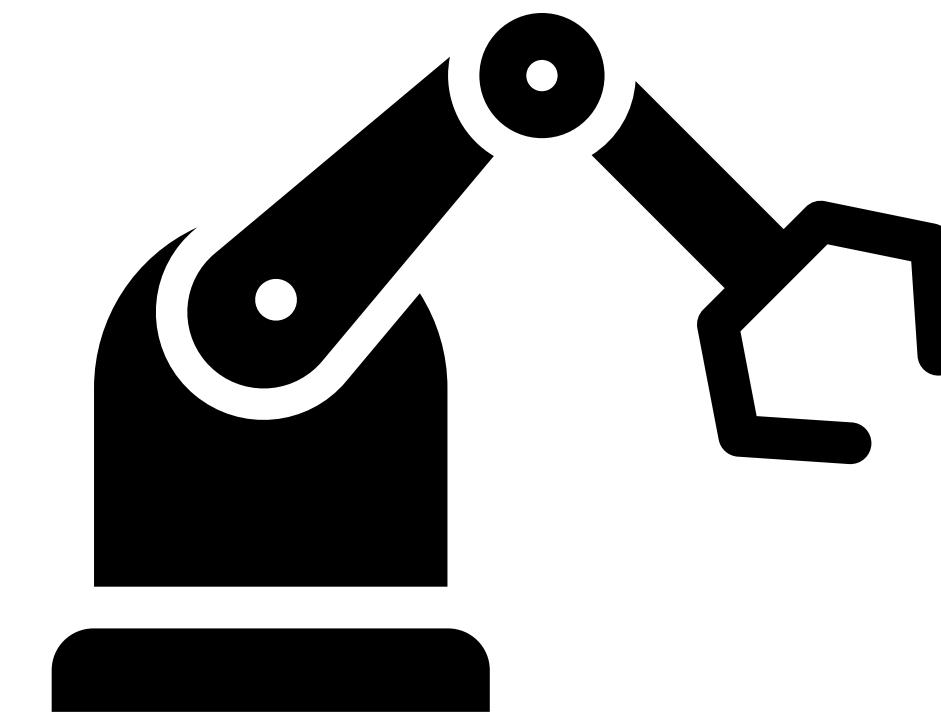
Pearl's
Explanation

Chain-of-
Thought

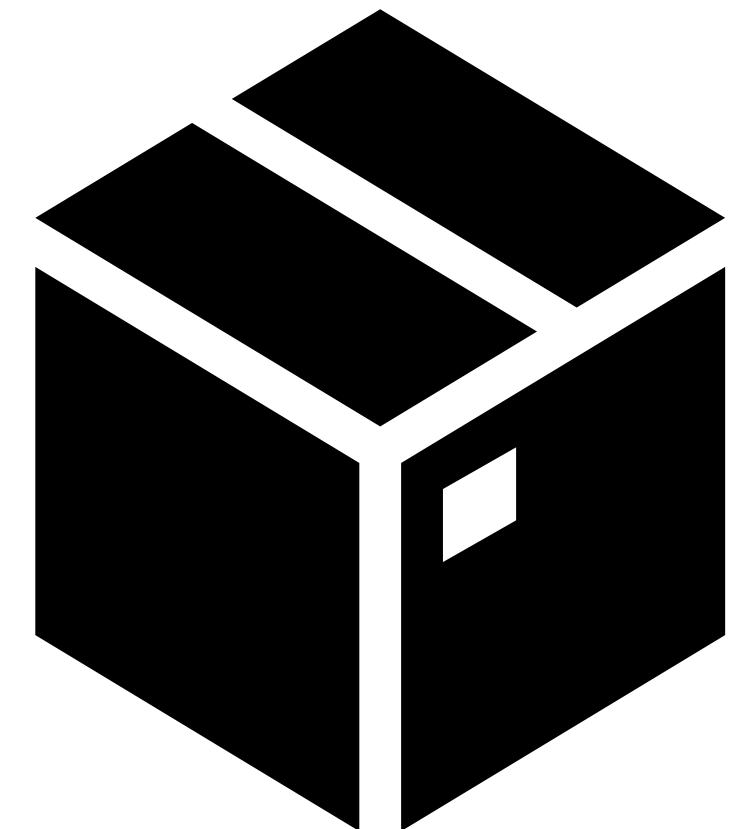
Cognitive Processes of Explanation



Causal
Connection



Explanation
Selection



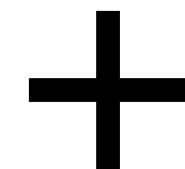
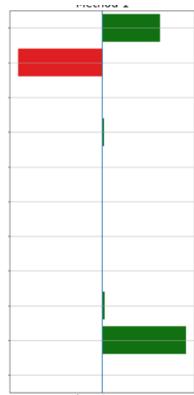
Explanation
Presentation

Explanation Selection Prompting Strategy

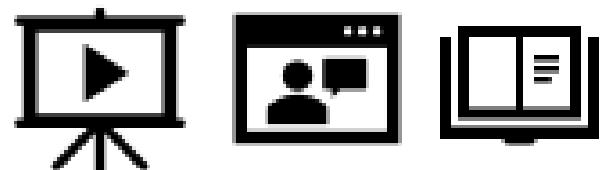
Model
Prediction



Explanation



Feature values



Course context

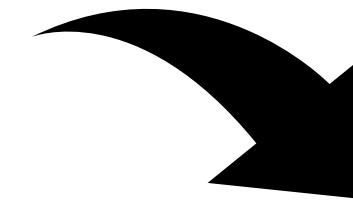
DSP 1

Theory Instructions

1. Select potential causes using these criteria:

- **Abnormality:** Tend to prefer abnormal causes.
- **Temporality:** Recent events are more relevant for the user and considered more mutable.
- **Controllability:** focus on the features that the student can control.

2. Select one explanation that **follows all of the criteria above** (Abnormality, Temporality, Controllability).



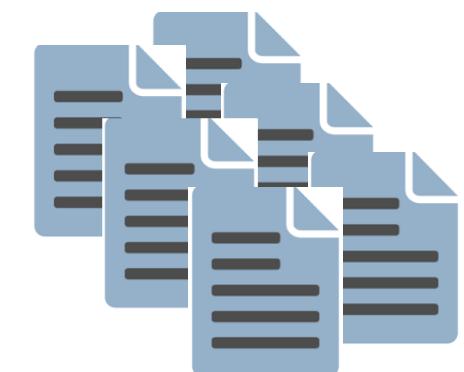
Language Model



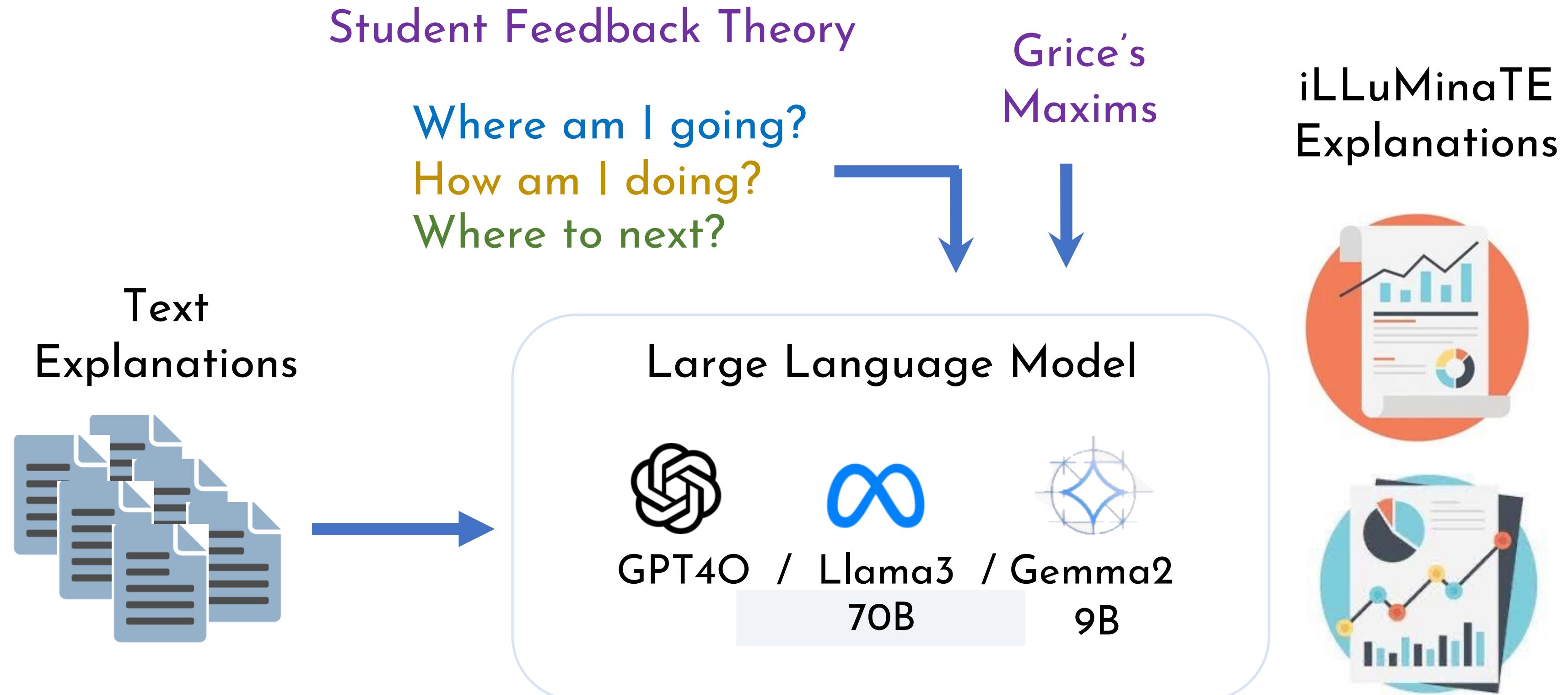
GPT4O / Llama3 / Gemma2
70B 9B



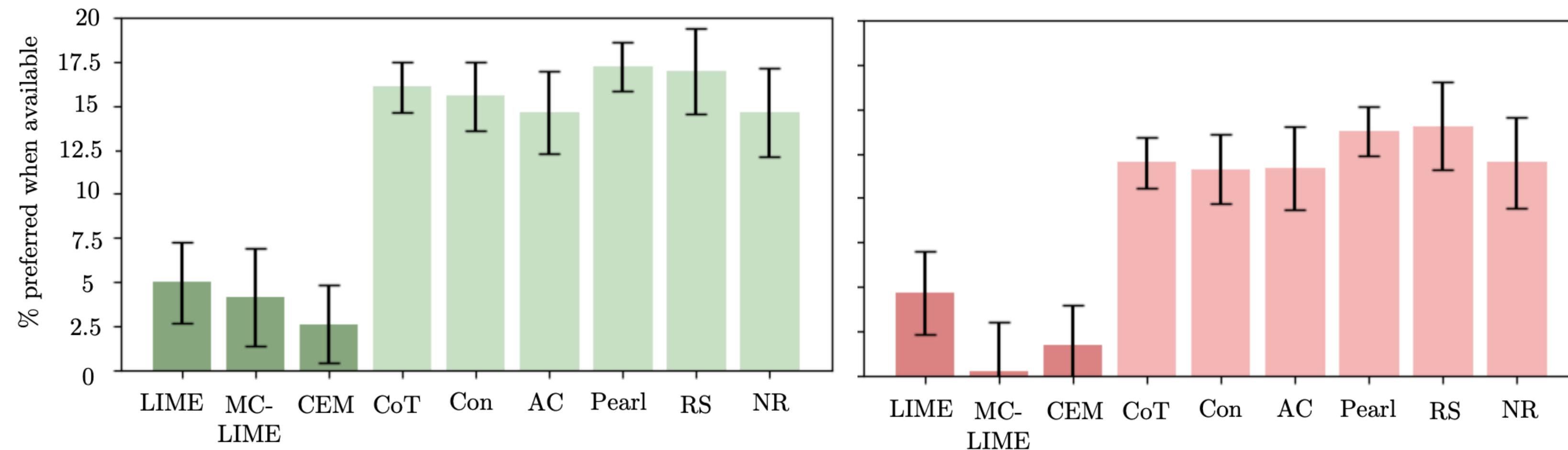
Text
Explanations



Explanation Presentation



Results

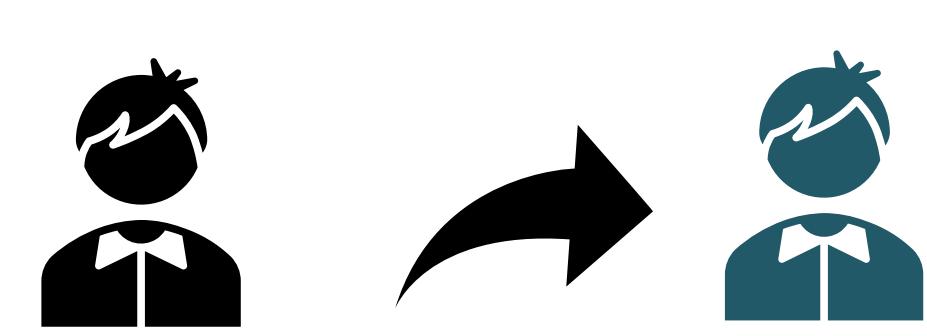


Students prefer iLLuMinaTE explanations significantly over expert processed post-hoc explanations!

text >> visual part of the explanation

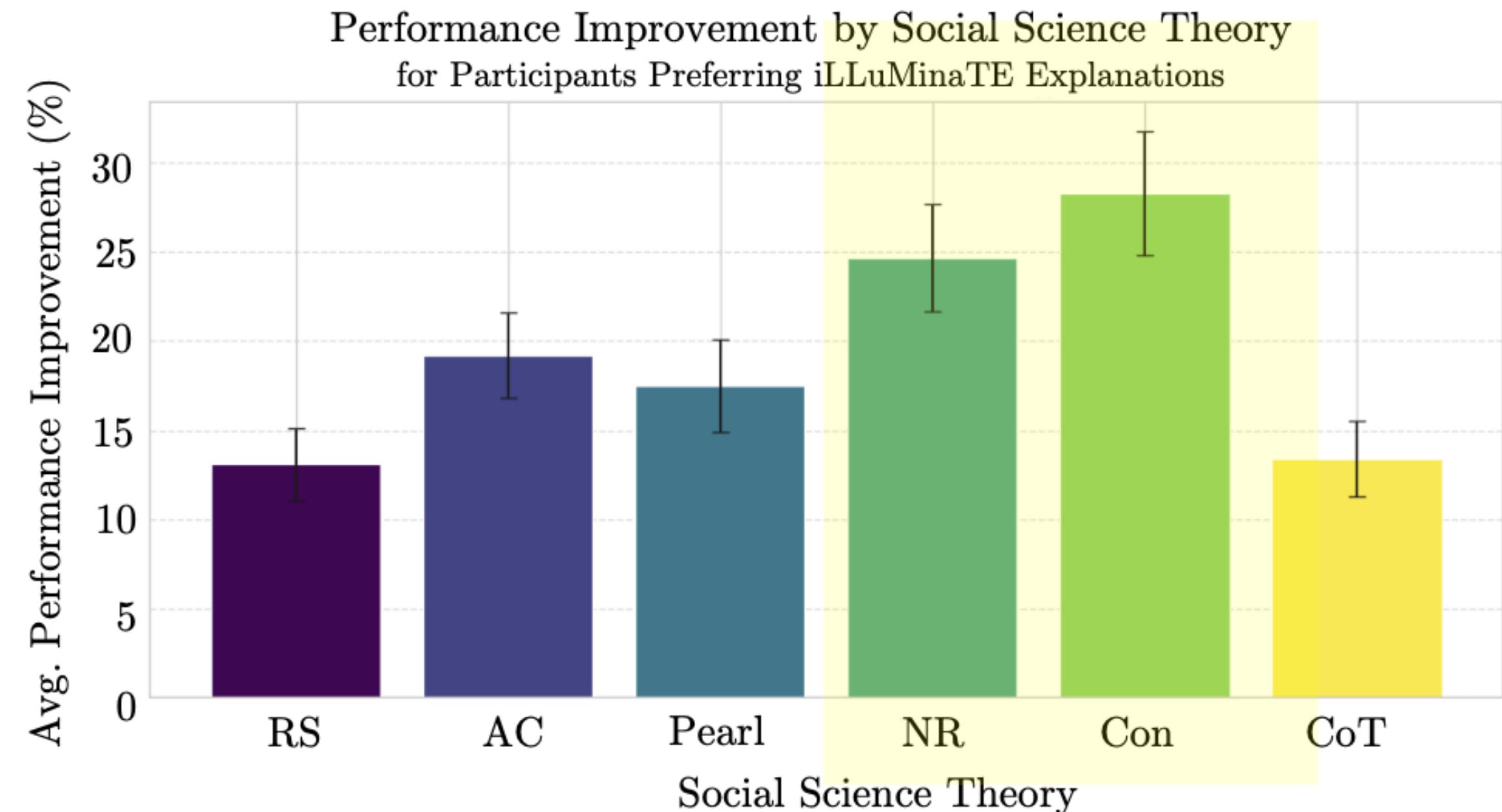
Actionability

students selected an action to take based on their preferred explanation

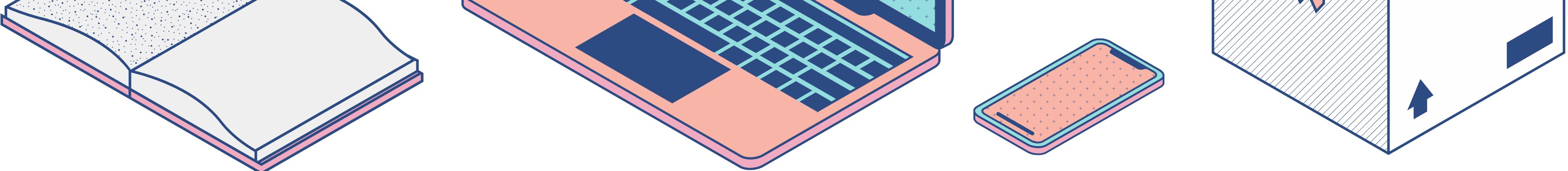


week 5

week 6
simulated intervention



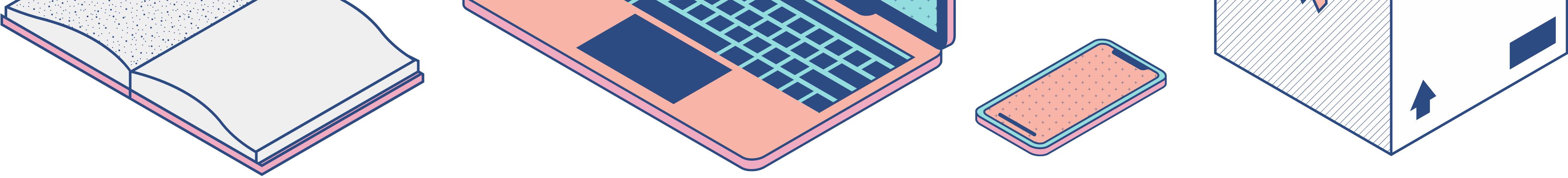
Contrastive, Necessity Robustness



Main Takeaways

EVALUATING THE EXPLAINERS
TRUSTING THE EXPLAINERS

Post-hoc explainers have serious problems,
and there is no effective way to validate them



Main Takeaways

MULTIMODN, INTERPRETCC, ILLUMINATE

With interpretable-by-design NNs,
guaranteed interpretability
does not have to come at the cost of performance
or human-understandability

Thank you!



Vinitra Swamy

vinitra.swamy@epfl.ch

vinitra.github.io

twitter: [@vinitra_s](https://twitter.com/vinitra_s)

