



Language
Technologies
Institute



Learning to Scaffold: Optimizing Model Explanations for Teaching

Patrick Fernandes*



Marcos Treviso*



Danish Pruthi



André Martins



Graham Neubig



How should we evaluate explanations?

- Explainability methods generally do not correlate with each other
- Most explanations do not help to predict the model's outputs and/or failures

How should we evaluate explanations?

- Explainability methods generally do not correlate with each other
- Most explanations do not help to predict the model's outputs and/or failures
- **Simulability:** "can we recover the model's output based on the explanation?"
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective

How should we evaluate explanations?

- Explainability methods generally do not correlate with each other
- Most explanations do not help to predict the model's outputs and/or failures
- **Simulability:** "can we recover the model's output based on the explanation?"
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective
- Pruthi et al. (2021) proposed a framework for measuring simulability that disregards **trivial protocols** 🥰

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

Simulability

(training time)

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(\underbrace{T(x)}_{\text{teacher}}, \underbrace{S_{\theta}(x)}_{\text{student}})]$$

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy

teacher *student*

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
↖

teacher *student*

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy

teacher *student*

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
teacher student

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Introducing explanations: Teacher and Student explainers $E_T(x), E_S(x)$

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
teacher student

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Introducing explanations: Teacher and Student explainers $E_T(x), E_S(x)$

$$\theta_E^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_{\theta}}(x))]$$

simulability loss

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
teacher student

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Introducing explanations: Teacher and Student explainers $E_T(x), E_S(x)$

$$\theta_E^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_{\theta}}(x))]$$

simulability loss *explainer regularizer (e.g.. KL)*

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
teacher student

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Introducing explanations: Teacher and Student explainers $E_T(x), E_S(x)$

$$\theta_E^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_{\theta}}(x))]$$

simulability loss *explainer regularizer (e.g.. KL)*

(standard simulability) $\text{SIM}(T, S_{\theta^*})$

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
teacher student

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Introducing explanations: Teacher and Student explainers $E_T(x), E_S(x)$

$$\theta_E^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_{\theta}}(x))]$$

simulability loss *explainer regularizer (e.g.. KL)*

(standard simulability) $\text{SIM}(T, S_{\theta^*}) < \text{SIM}(T, S_{\theta_E^*})$ (scaffolded simulability)

Simulability

(training time) $\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x))]$

cross entropy
teacher student

(test time) $\text{SIM}(T, S_{\theta}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [1\{T(x) = S_{\theta}(x)\}]$

agreement

Introducing explanations: Teacher and Student explainers $E_T(x), E_S(x)$

$$\theta_E^* = \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta}(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_{\theta}}(x))]$$

simulability loss *explainer regularizer (e.g.. KL)*

(standard simulability) $\text{SIM}(T, S_{\theta^*}) < \text{SIM}(T, S_{\theta_E^*})$ (scaffolded simulability)

Can we **learn explainers** $\phi(E)$ that optimize **simulability**?

(scaffolded simulability) $\text{SIM}(T, S_{\theta_E^*}) < \text{SIM}(T, S_{\theta_{\phi(E)}})$ (optim. scaffolded simulability)

Optimizing Explainers for Teaching

- Scaffold-**Maximizing Training (SMaT)** framework

$$\mathcal{L}_{\text{student}}(x; T, E_T, S_\theta, E_S) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_\theta}(x))$$

Optimizing Explainers for Teaching

- Scaffold-**Maximizing Training (SMaT)** framework

$$\mathcal{L}_{\text{student}}(x; T, E_T, S_\theta, E_S) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_\theta}(x))$$

Optimizing Explainers for Teaching

- Scaffold-**Maximizing Training (SMaT)** framework

$$\mathcal{L}_{\text{student}}(x; T, E_T, S_\theta, E_S) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_\theta}(x))$$

$$\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_\theta, E_{\phi_S}) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_{\phi_T}(x), E_{\phi_S}(x))$$

simulability loss

parameterized explainers

Optimizing Explainers for Teaching

- Scaffold-**Maximizing Training (SMaT)** framework

$$\mathcal{L}_{\text{student}}(x; T, E_T, S_\theta, E_S) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_\theta}(x))$$

$$\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_\theta, E_{\phi_S}) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_{\phi_T}(x), E_{\phi_S}(x))$$

simulability loss *parameterized explainers*

- Bi-level optimization:

(inner opt.) $\theta^*(\phi_T), \phi_S^*(\phi_T) = \arg \max_{\theta, \phi_S} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_\theta, E_{\phi_S})]$

student parameters and student explainer parameters

Optimizing Explainers for Teaching

- Scaffold-**Maximizing Training (SMaT)** framework

$$\mathcal{L}_{\text{student}}(x; T, E_T, S_\theta, E_S) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_T(x), E_{S_\theta}(x))$$

$$\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_\theta, E_{\phi_S}) = \mathcal{L}_{\text{sim}}(T(x), S_\theta(x)) + \beta \mathcal{L}_{\text{expl}}(E_{\phi_T}(x), E_{\phi_S}(x))$$

simulability loss *parameterized explainers*

- Bi-level optimization:

(inner opt.) $\theta^*(\phi_T), \phi_S^*(\phi_T) = \arg \max_{\theta, \phi_S} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_\theta, E_{\phi_S})]$

student parameters and student explainer parameters

(outer opt.) $\phi_T^* = \arg \max_{\phi_T} \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta^*(\phi_T)})]$

teacher explainer parameters

Optimizing Explainers for Teaching

How can we optimize this?

- Bi-level optimization:

(inner opt.) $\theta^*(\phi_T), \phi_S^*(\phi_T) = \arg \max_{\theta, \phi_S} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_{\theta}, E_{\phi_S})]$

student parameters and student explainer parameters

(outer opt.) $\phi_T^* = \arg \max_{\phi_T} \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta^*(\phi_T)})]$

teacher explainer parameters

Optimizing Explainers for Teaching

How can we optimize this?

- Assume the explainers are differentiable

- Bi-level optimization:

(inner opt.) $\theta^*(\phi_T), \phi_S^*(\phi_T) = \arg \max_{\theta, \phi_S} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_{\theta}, E_{\phi_S})]$

student parameters and student explainer parameters

(outer opt.) $\phi_T^* = \arg \max_{\phi_T} \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta^*(\phi_T)})]$

teacher explainer parameters

Optimizing Explainers for Teaching

How can we optimize this?

- Assume the explainers are differentiable
- Explicit differentiation with a truncated gradient update

- Bi-level optimization:

(inner opt.) $\theta^*(\phi_T), \phi_S^*(\phi_T) = \arg \max_{\theta, \phi_S} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_{\theta}, E_{\phi_S})]$

student parameters and student explainer parameters

(outer opt.) $\phi_T^* = \arg \max_{\phi_T} \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta^*(\phi_T)})]$

teacher explainer parameters

Optimizing Explainers for Teaching

How can we optimize this?

- Assume the explainers are differentiable
- Explicit differentiation with a truncated gradient update
- Diff. through a gradient operation \Leftrightarrow JAX for Hessian-vector products

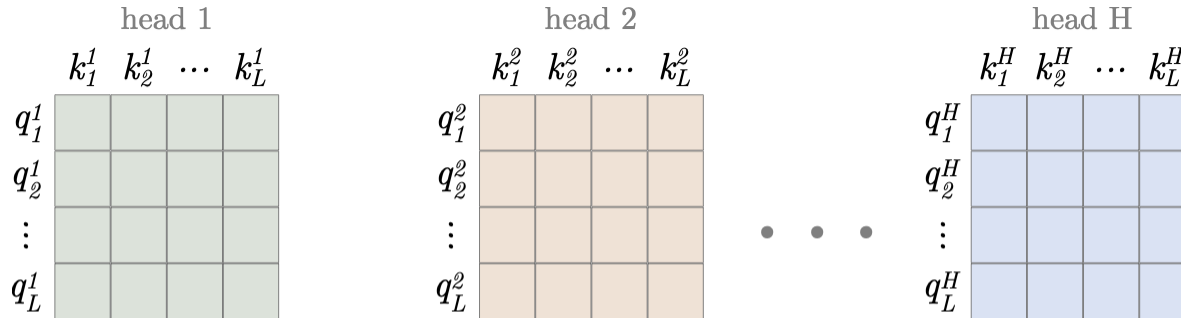
- Bi-level optimization:

(inner opt.) $\theta^*(\phi_T), \phi_S^*(\phi_T) = \arg \max_{\theta, \phi_S} \mathbb{E}_{x \sim \mathcal{D}_{\text{train}}} [\mathcal{L}_{\text{student}}(x; T, E_{\phi_T}, S_{\theta}, E_{\phi_S})]$
student parameters and student explainer parameters

(outer opt.) $\phi_T^* = \arg \max_{\phi_T} \mathbb{E}_{x \sim \mathcal{D}_{\text{test}}} [\mathcal{L}_{\text{sim}}(T(x), S_{\theta^*(\phi_T)})]$
teacher explainer parameters

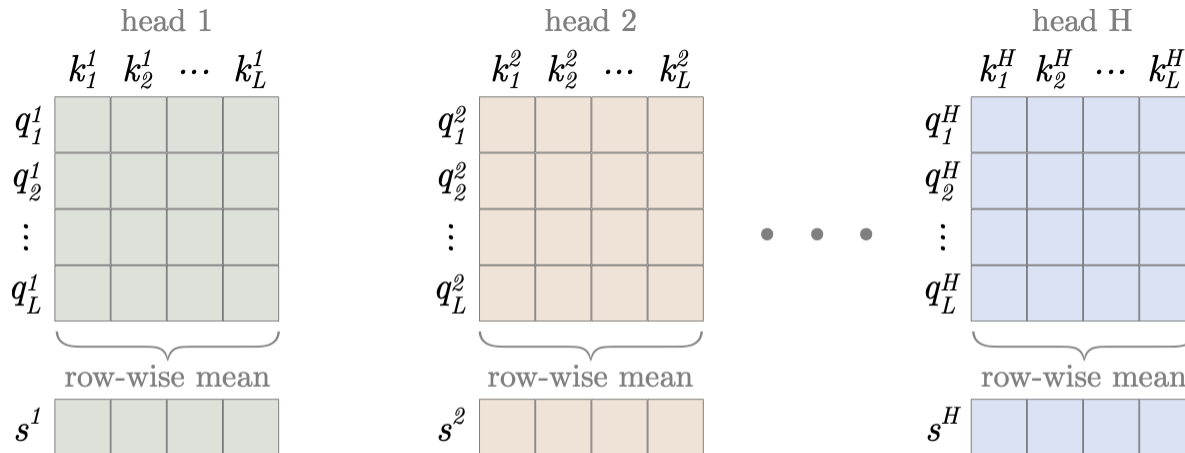
Differentiable, Parameterized Explainer

- Head-level parameterization:



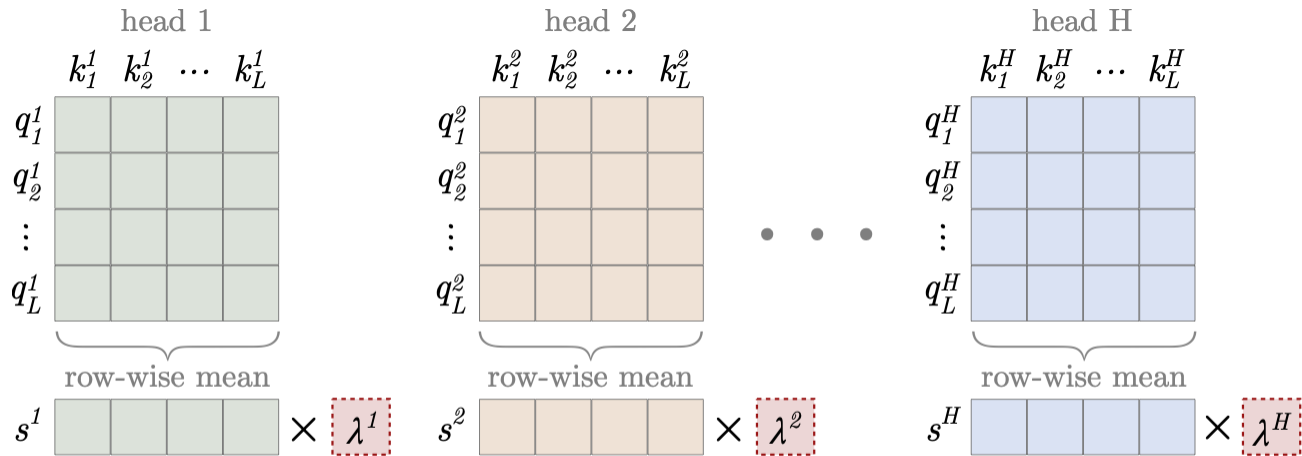
Differentiable, Parameterized Explainer

- Head-level parameterization:



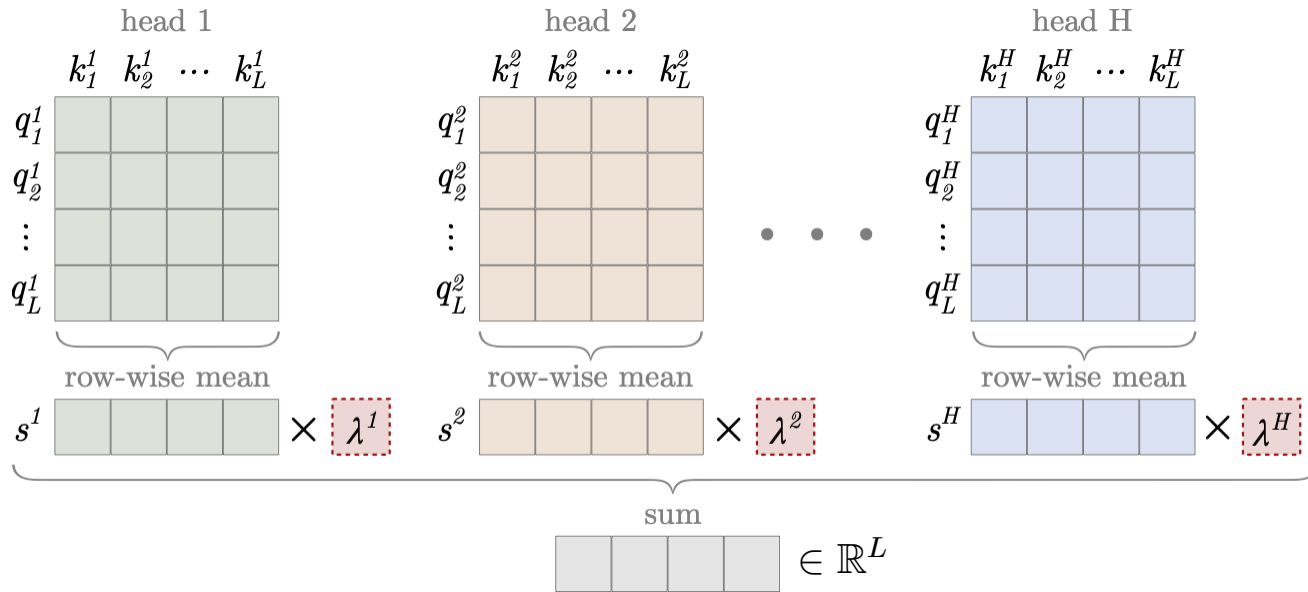
Differentiable, Parameterized Explainer

- Head-level parameterization:



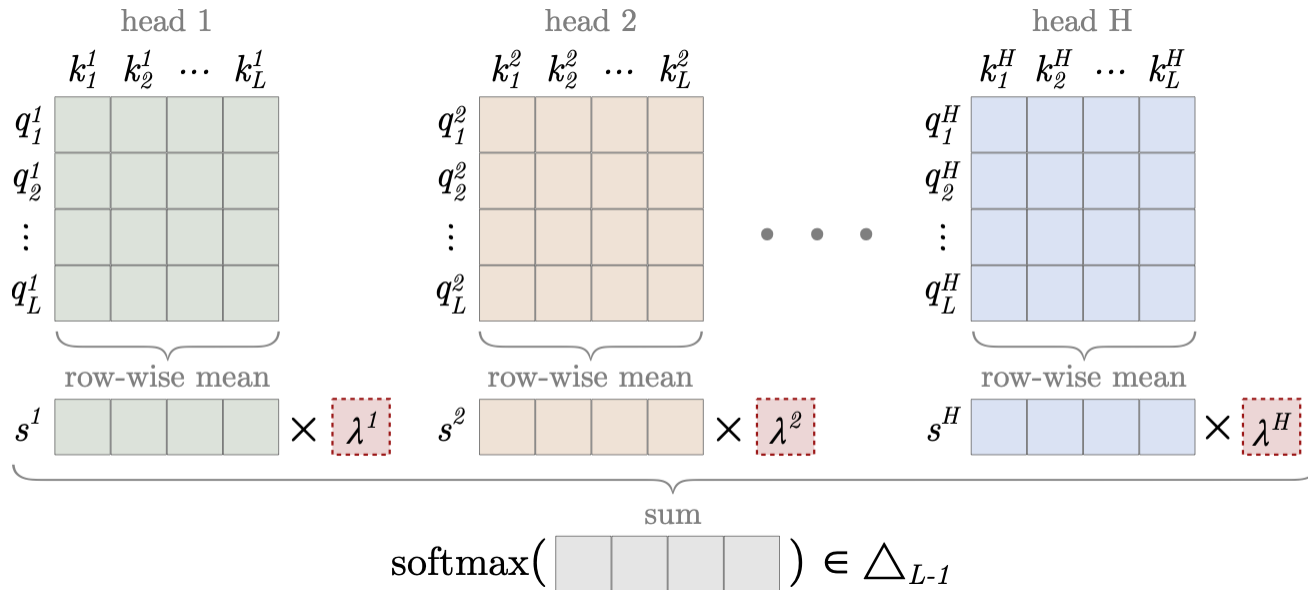
Differentiable, Parameterized Explainer

- Head-level parameterization:



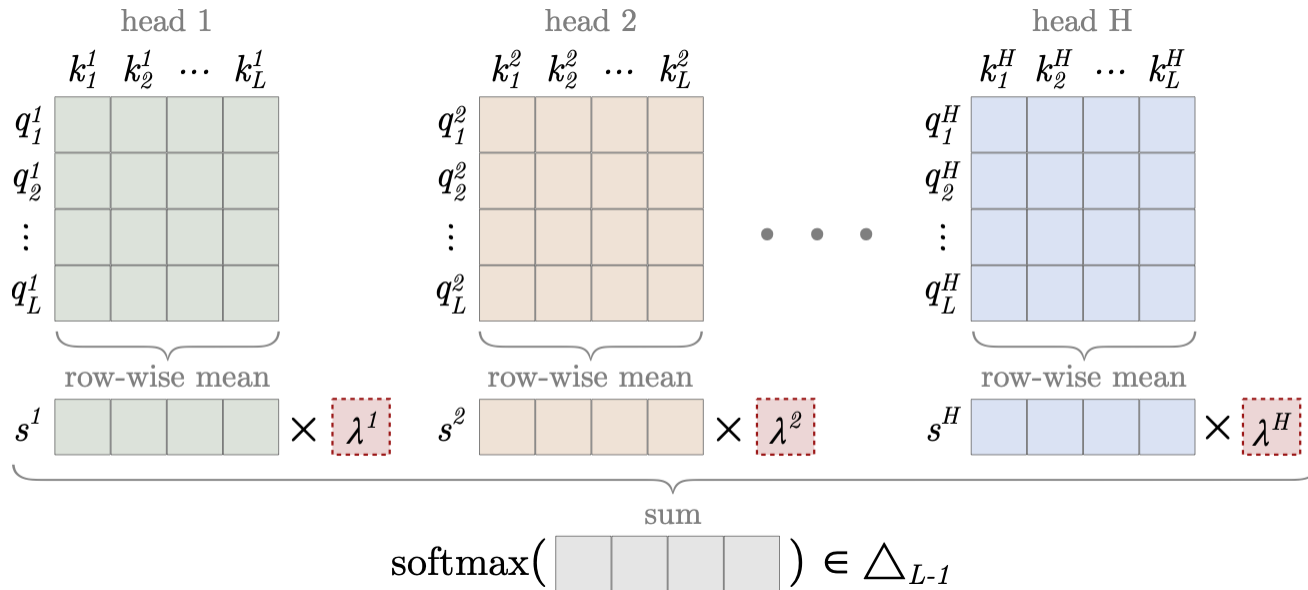
Differentiable, Parameterized Explainer

- Head-level parameterization:



Differentiable, Parameterized Explainer

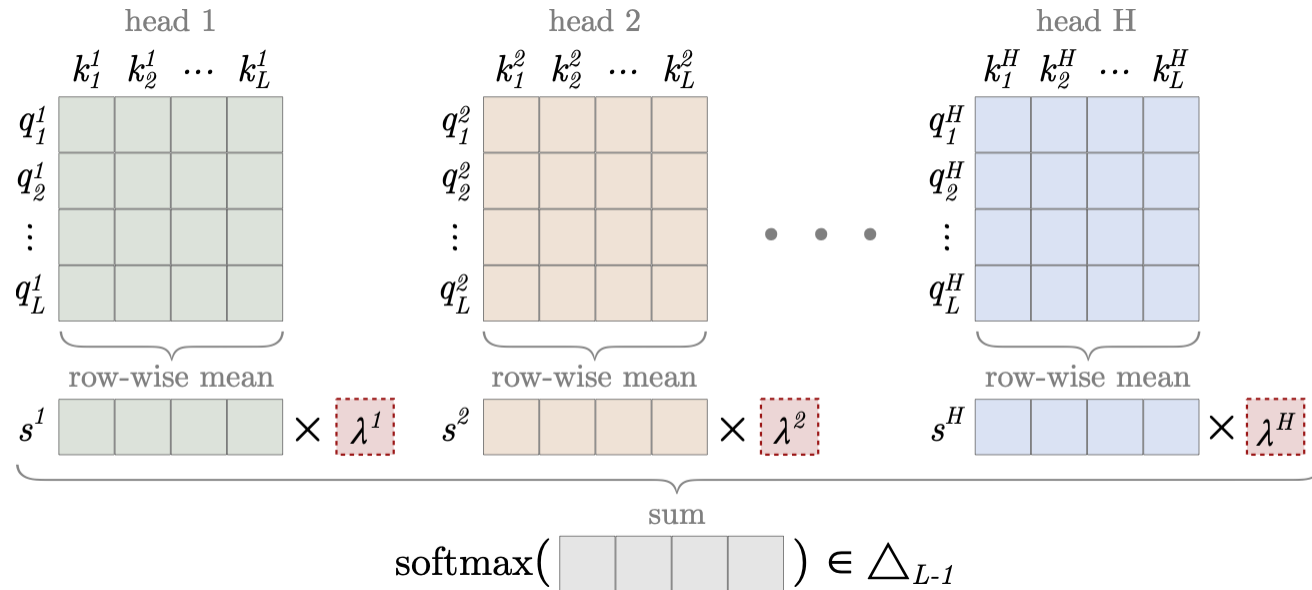
- Head-level parameterization:



$$\lambda_T = \text{normalize}(\phi_T) \in \Delta_{H-1}$$

Differentiable, Parameterized Explainer

- Head-level parameterization:

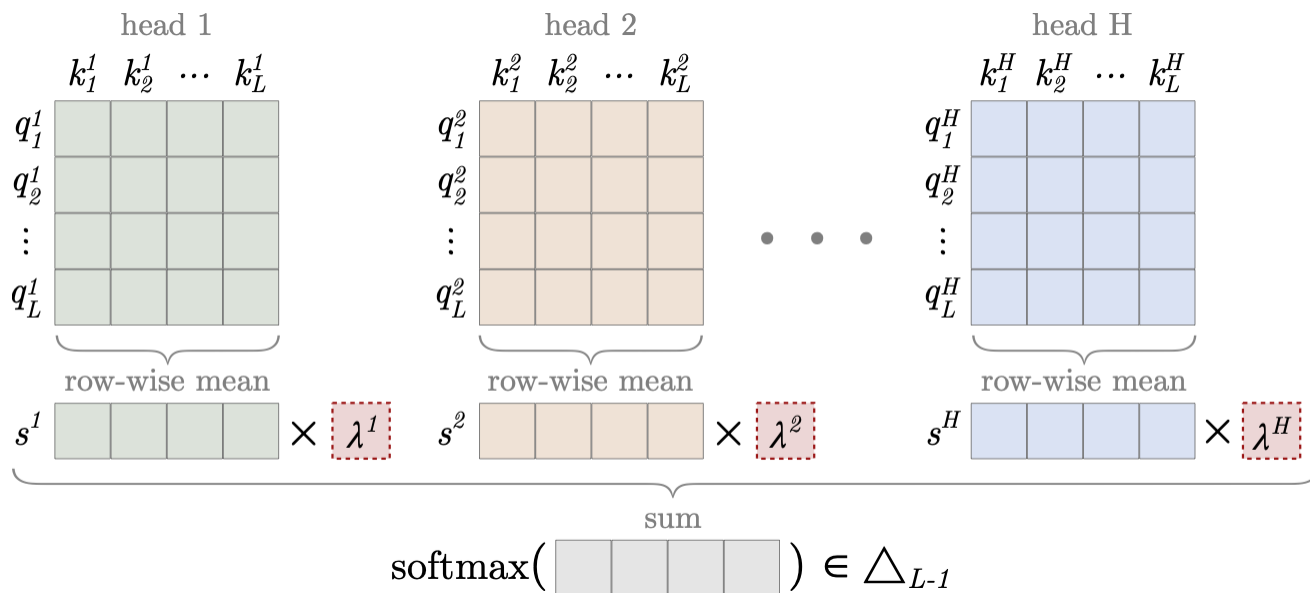


$$\lambda_T = \text{normalize}(\phi_T) \in \Delta_{H-1}$$

$$\text{sparsemax}(z) = \arg \min_{p \in \Delta_{H-1}} \|p - z\|_2$$

Differentiable, Parameterized Explainer

- Head-level parameterization:



$$\lambda_T = \text{normalize}(\phi_T) \in \Delta_{H-1}$$

$$\text{sparsemax}(z) = \arg \min_{p \in \Delta_{H-1}} \|p - z\|_2$$



Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



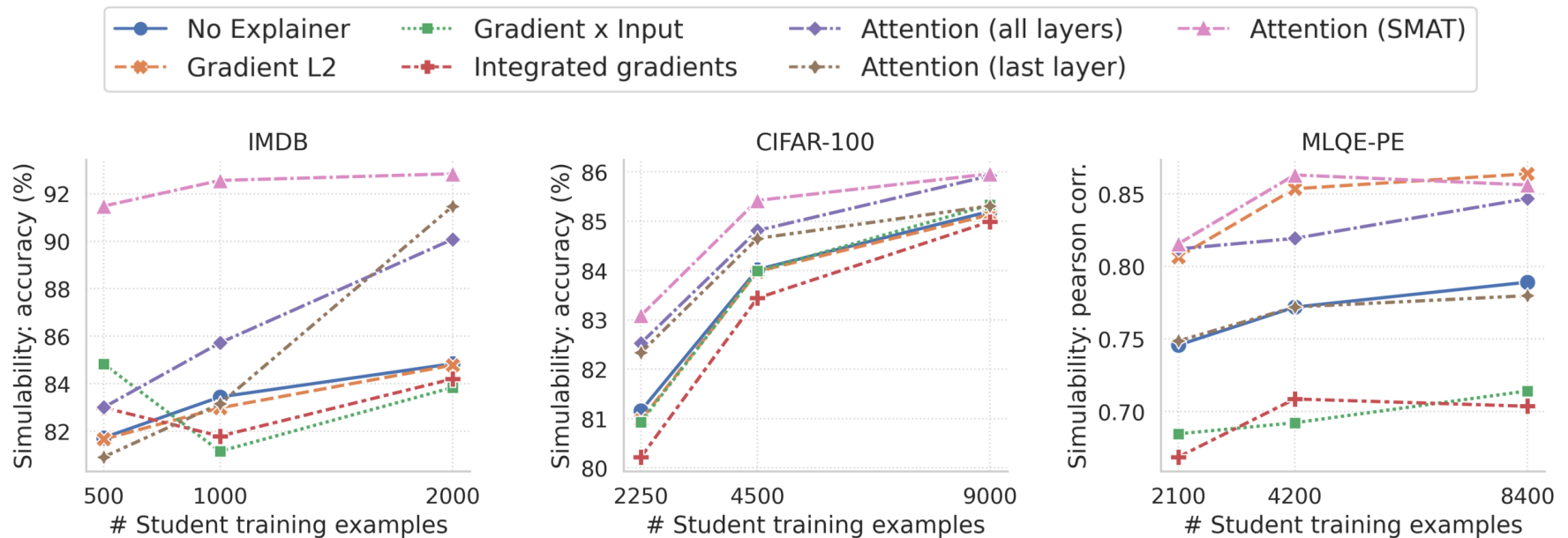
Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



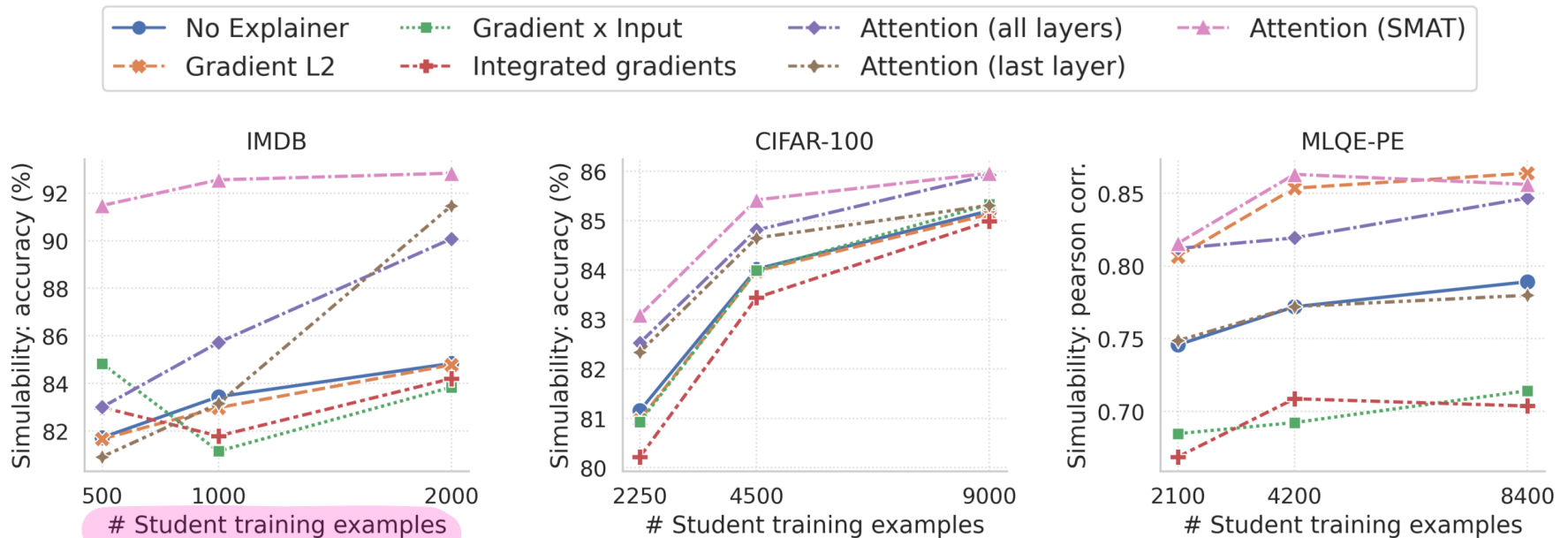
Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



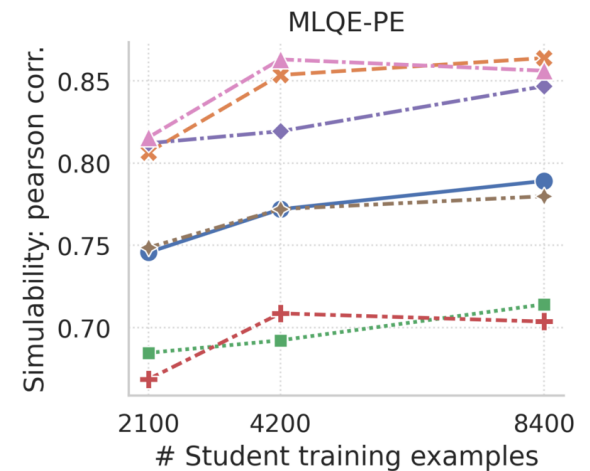
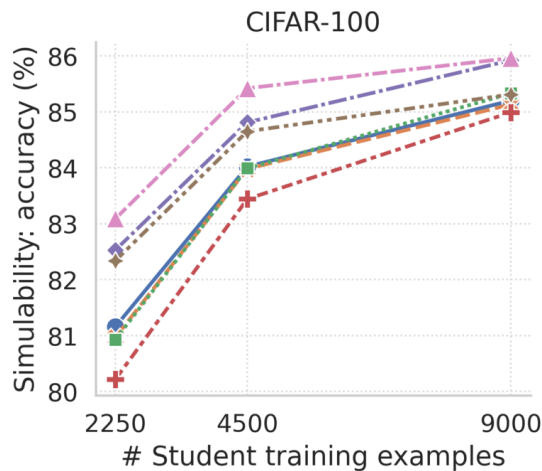
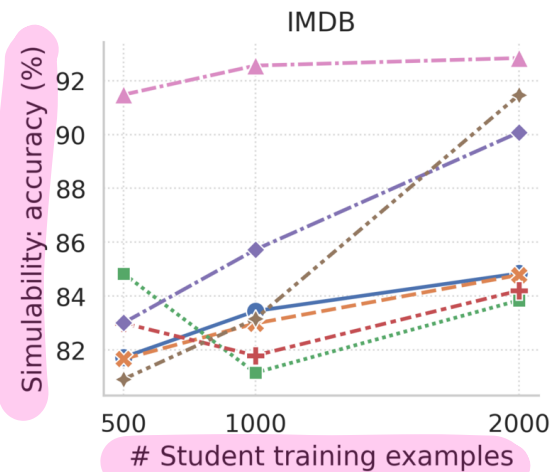
Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



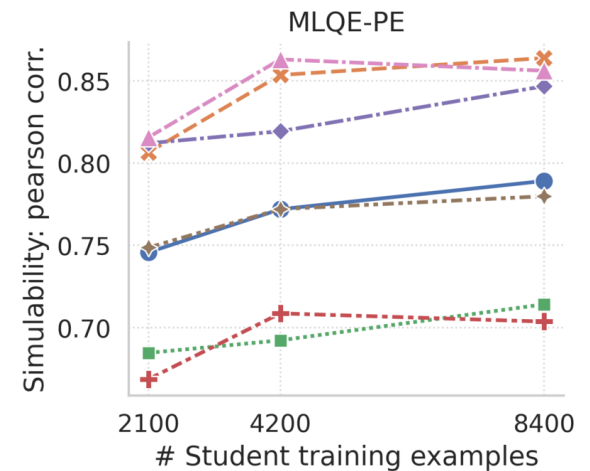
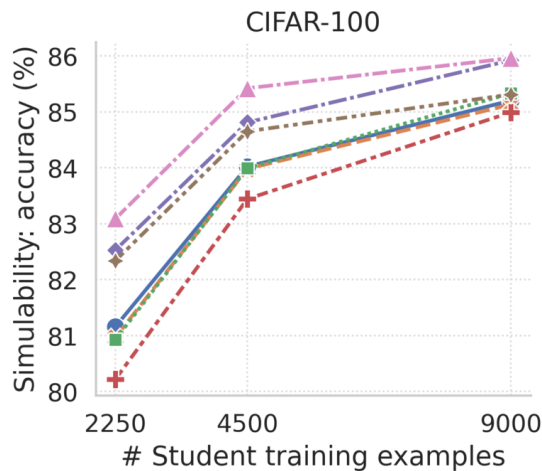
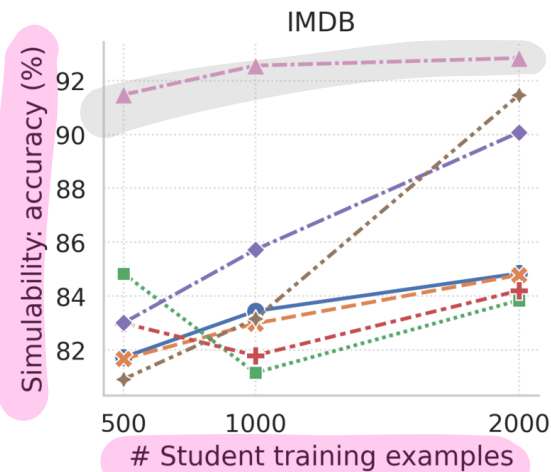
Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



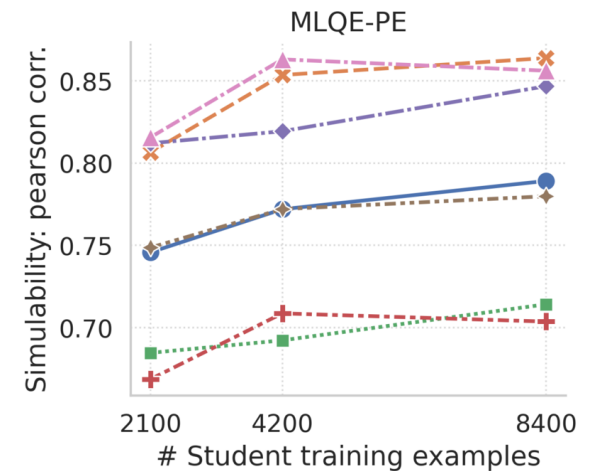
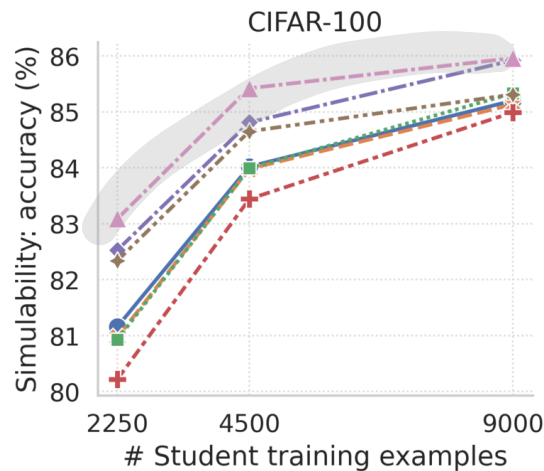
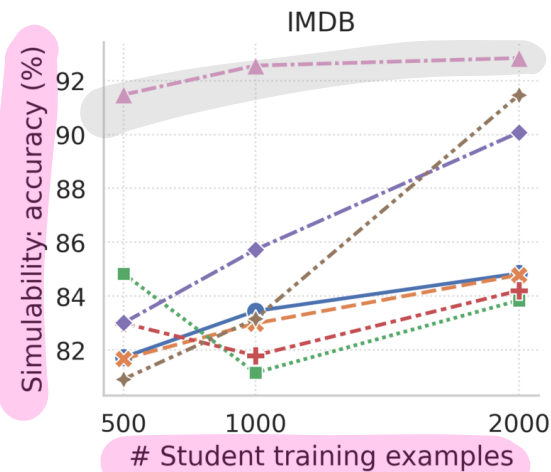
Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



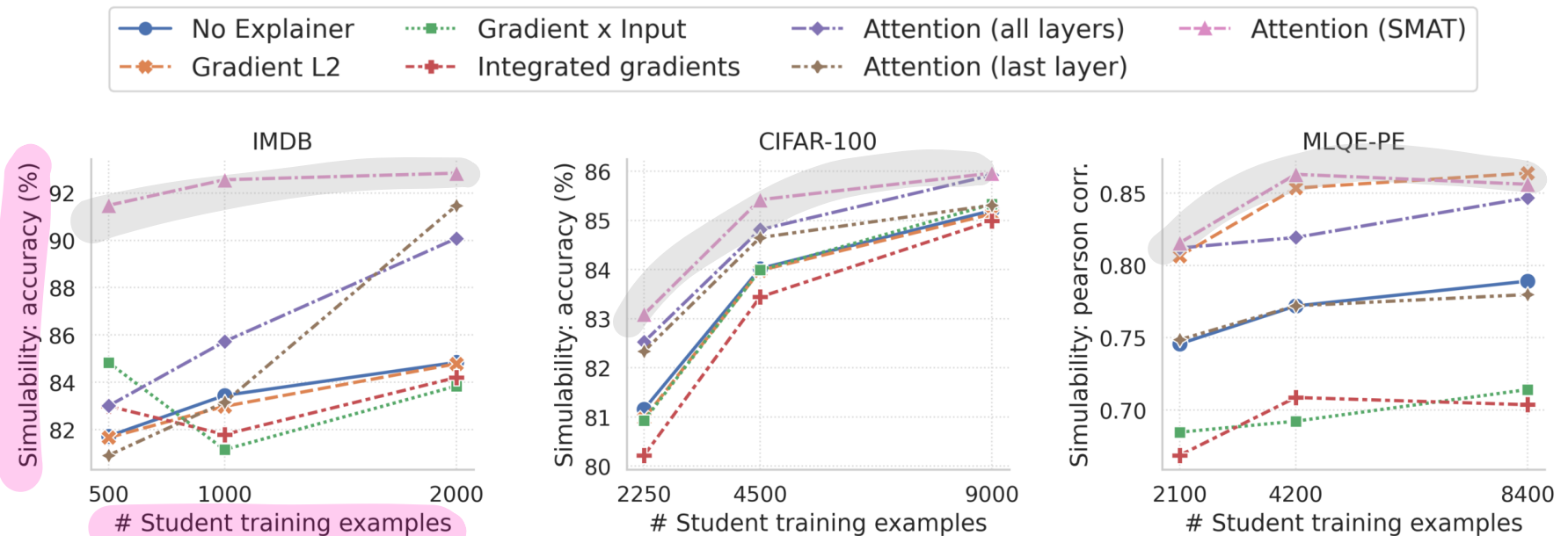
Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



Experiments: simulability

- Text classification (IMDB)
- Image classification (CIFAR-100)
- Machine Translation Quality Estimation (MLQE-PE)



Experiments: plausibility

- *Plausibility (human-likeness)* of the explainers

Text Classification	
	AUC
Grad. L2	0.65
Grad. \times Input	0.51
Integrated Grad.	0.53
Attn. (<i>all layers</i>)	0.68
Attn. (<i>last layer</i>)	0.61
Attn. (SMaT)	0.73
<hr/>	
Attn. (<i>best layer</i>)*	0.75
Attn. (<i>best head</i>)*	0.75

Image Classification		
	Rank	TrueSkill
Grad. \times Input	3-4	-2.7 \pm .67
Integ. Grad.	3-4	-2.1 \pm .67
Attn. (<i>all lx.</i>)	2	0.7 \pm .67
Attn. (SMaT)	1	4.3\pm.70

Quality Estimation		
	OVERALL	
	src.	tgt.
Gradient L2	0.67	0.59
Gradient \times Input	0.61	0.54
Integrated Gradients	0.62	0.53
Attention (<i>all layers</i>)	0.62	0.59
Attention (<i>last layer</i>)	0.54	0.50
Attention (SMaT)	0.66	0.60
<hr/>		
Attention (<i>best layer</i>)*	0.65	0.65
Attention (<i>best head</i>)*	0.67	0.66

Experiments: plausibility

- *Plausibility (human-likeness)* of the explainers

Text Classification	
	AUC
Grad. L2	0.65
Grad. \times Input	0.51
Integrated Grad.	0.53
Attn. (<i>all layers</i>)	0.68
Attn. (<i>last layer</i>)	0.61
Attn. (SMaT)	0.73
<hr/>	
Attn. (<i>best layer</i>)*	0.75
Attn. (<i>best head</i>)*	0.75

Image Classification		
	Rank	TrueSkill
Grad. \times Input	3-4	-2.7 \pm .67
Integ. Grad.	3-4	-2.1 \pm .67
Attn. (<i>all lx.</i>)	2	0.7 \pm .67
Attn. (SMaT)	1	4.3\pm.70

Quality Estimation		
	OVERALL	
	src.	tgt.
Gradient L2	0.67	0.59
Gradient \times Input	0.61	0.54
Integrated Gradients	0.62	0.53
Attention (<i>all layers</i>)	0.62	0.59
Attention (<i>last layer</i>)	0.54	0.50
Attention (SMaT)	0.66	0.60
<hr/>		
Attention (<i>best layer</i>)*	0.65	0.65
Attention (<i>best head</i>)*	0.67	0.66

Experiments: plausibility

- *Plausibility (human-likeness)* of the explainers

Text Classification	
	AUC
Grad. L2	0.65
Grad. \times Input	0.51
Integrated Grad.	0.53
Attn. (<i>all layers</i>)	0.68
Attn. (<i>last layer</i>)	0.61
Attn. (SMaT)	0.73
Attn. (<i>best layer</i>)*	0.75
Attn. (<i>best head</i>)*	0.75

Image Classification		
	Rank	TrueSkill
Grad. \times Input	3-4	-2.7 \pm .67
Integ. Grad.	3-4	-2.1 \pm .67
Attn. (<i>all lx.</i>)	2	0.7 \pm .67
Attn. (SMaT)	1	4.3\pm.70

Quality Estimation		
	OVERALL	
	src.	tgt.
Gradient L2	0.67	0.59
Gradient \times Input	0.61	0.54
Integrated Gradients	0.62	0.53
Attention (<i>all layers</i>)	0.62	0.59
Attention (<i>last layer</i>)	0.54	0.50
Attention (SMaT)	0.66	0.60
Attention (<i>best layer</i>)*	0.65	0.65
Attention (<i>best head</i>)*	0.67	0.66

attention (all layers): i ' ve seen river ##dance in person and nothing compares to the video , but the show is awesome , the dancers are amazing . the music is impact ##ing . and the overall performance is outstanding . i ' ve never seen anything like it ! i suggest that you see this show if you can ! ! !

attention (SMaT): i ' ve seen river ##dance in person and nothing compares to the video , but the show is awesome . the dancers are amazing . the music is impact ##ing . and the overall performance is outstanding . i ' ve never seen anything like it ! i suggest that you see this show if you can ! ! !

Experiments: plausibility

- *Plausibility (human-likeness)* of the explainers

Text Classification	
	AUC
Grad. L2	0.65
Grad. \times Input	0.51
Integrated Grad.	0.53
Attn. (<i>all layers</i>)	0.68
Attn. (<i>last layer</i>)	0.61
Attn. (SMaT)	0.73
Attn. (<i>best layer</i>)*	0.75
Attn. (<i>best head</i>)*	0.75

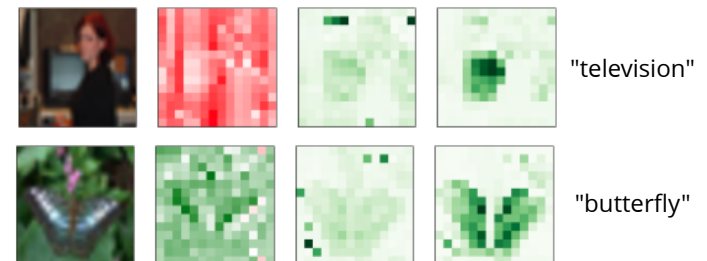
Image Classification		
	Rank	TrueSkill
Grad. \times Input	3-4	-2.7 \pm .67
Integ. Grad.	3-4	-2.1 \pm .67
Attn. (<i>all lx.</i>)	2	0.7 \pm .67
Attn. (SMaT)	1	4.3\pm.70

Quality Estimation		
	OVERALL	
	src.	tgt.
Gradient L2	0.67	0.59
Gradient \times Input	0.61	0.54
Integrated Gradients	0.62	0.53
Attention (<i>all layers</i>)	0.62	0.59
Attention (<i>last layer</i>)	0.54	0.50
Attention (SMaT)	0.66	0.60
Attention (<i>best layer</i>)*	0.65	0.65
Attention (<i>best head</i>)*	0.67	0.66

attention (all layers): i ' ve seen river ##dance in person and nothing compares to the video , but the show is awesome , the dancers are amazing . the music is impact ##ing . and the overall performance is outstanding . i ' ve never seen anything like it ! i suggest that you see this show if you can ! ! !

attention (SMaT): i ' ve seen river ##dance in person and nothing compares to the video , but the show is awesome . the dancers are amazing . the music is impact ##ing . and the overall performance is outstanding . i ' ve never seen anything like it ! i suggest that you see this show if you can ! ! !

Input image Integ. Grad. Attn. (all lx.) Attn. (SMaT)



Experiments: head projection

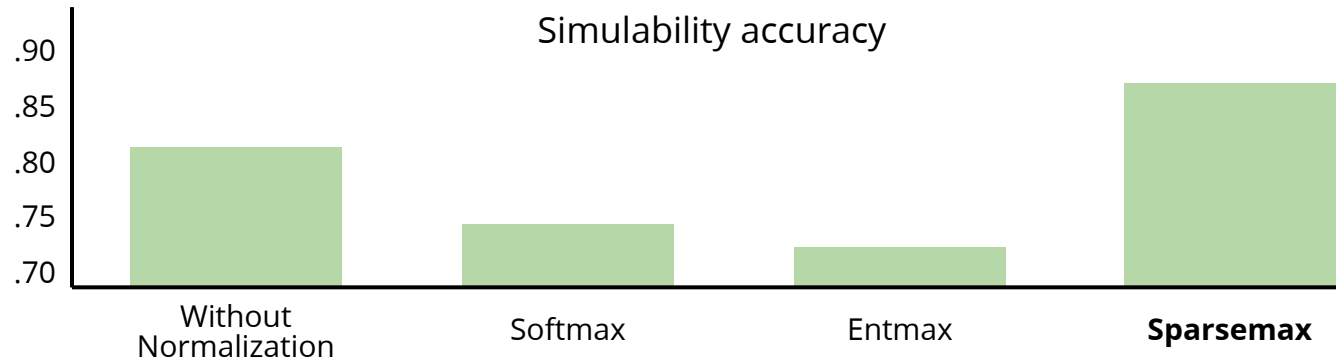
- Normalization functions

$$\lambda_T = \text{normalize}(\phi_T) \in \Delta_{H-1}$$

Experiments: head projection

- Normalization functions

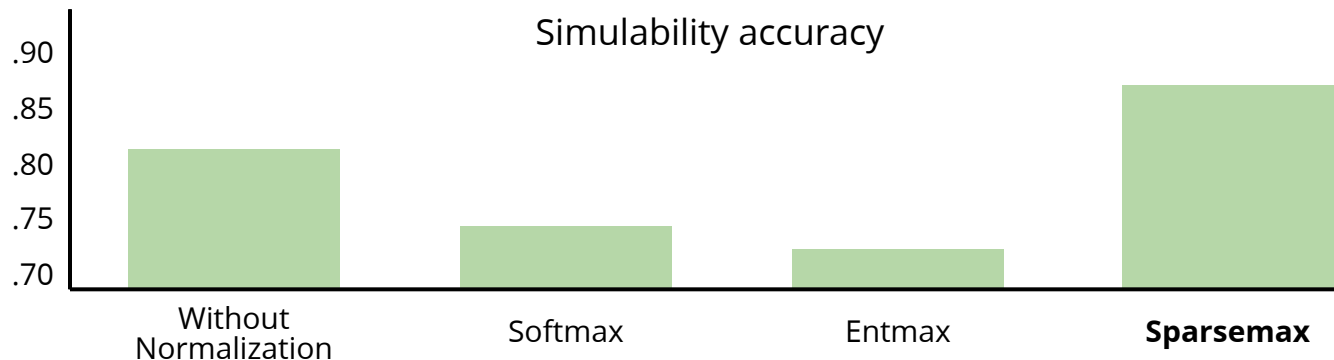
$$\lambda_T = \text{normalize}(\phi_T) \in \Delta_{H-1}$$



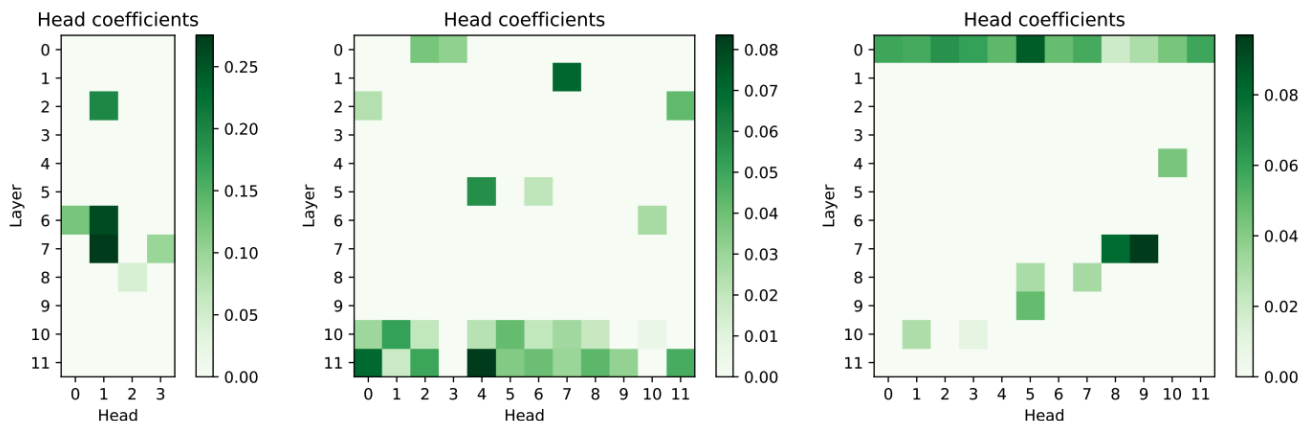
Experiments: head projection

- Normalization functions

$$\lambda_T = \text{normalize}(\phi_T) \in \Delta_{H-1}$$

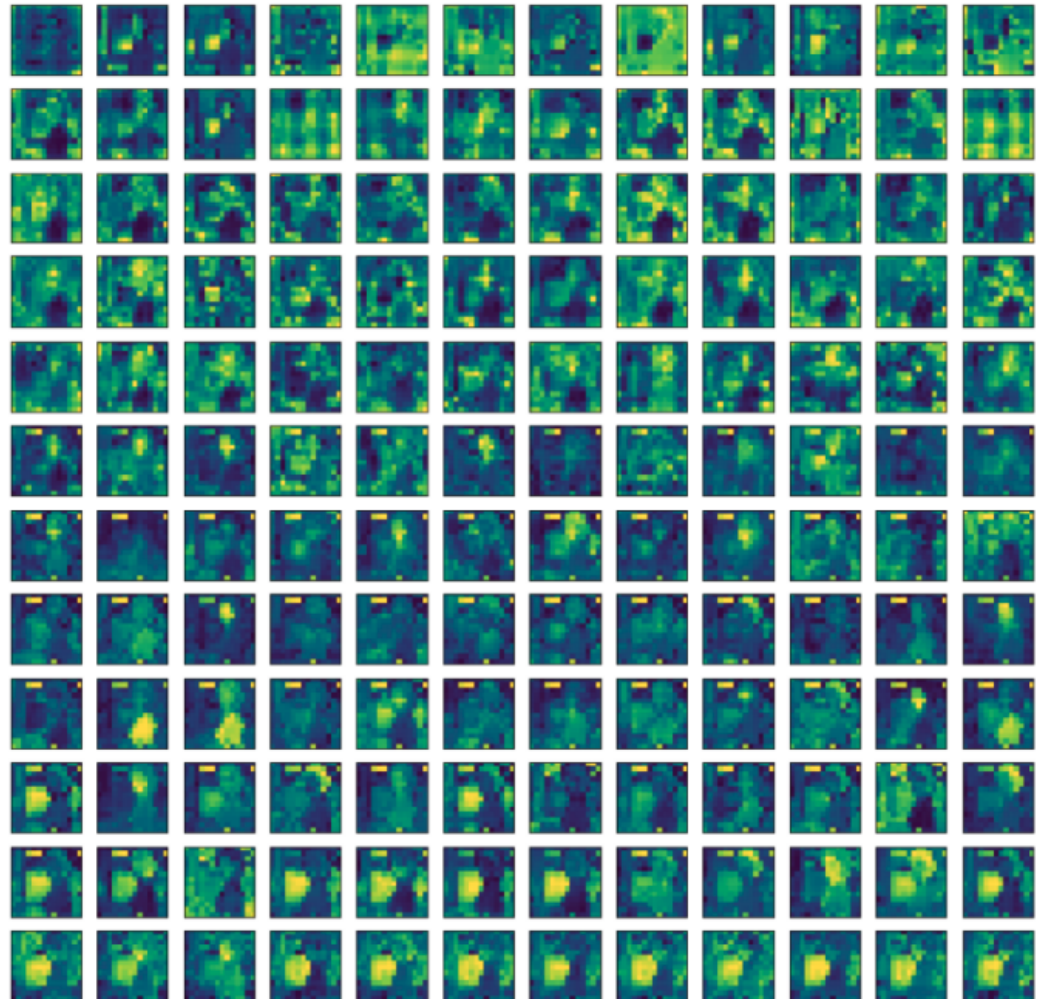
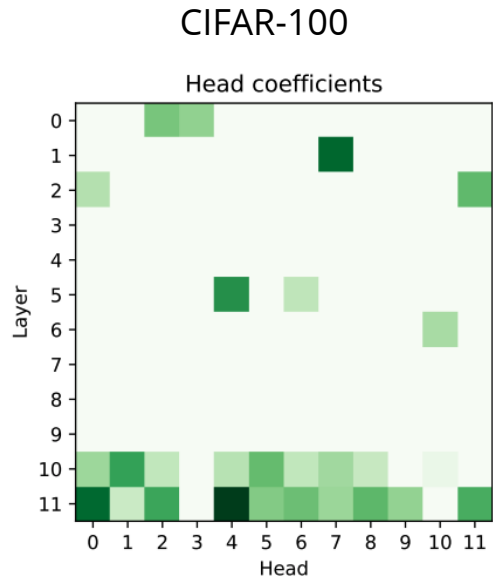


- Only a small subset of attention heads are deemed relevant by SMA_T



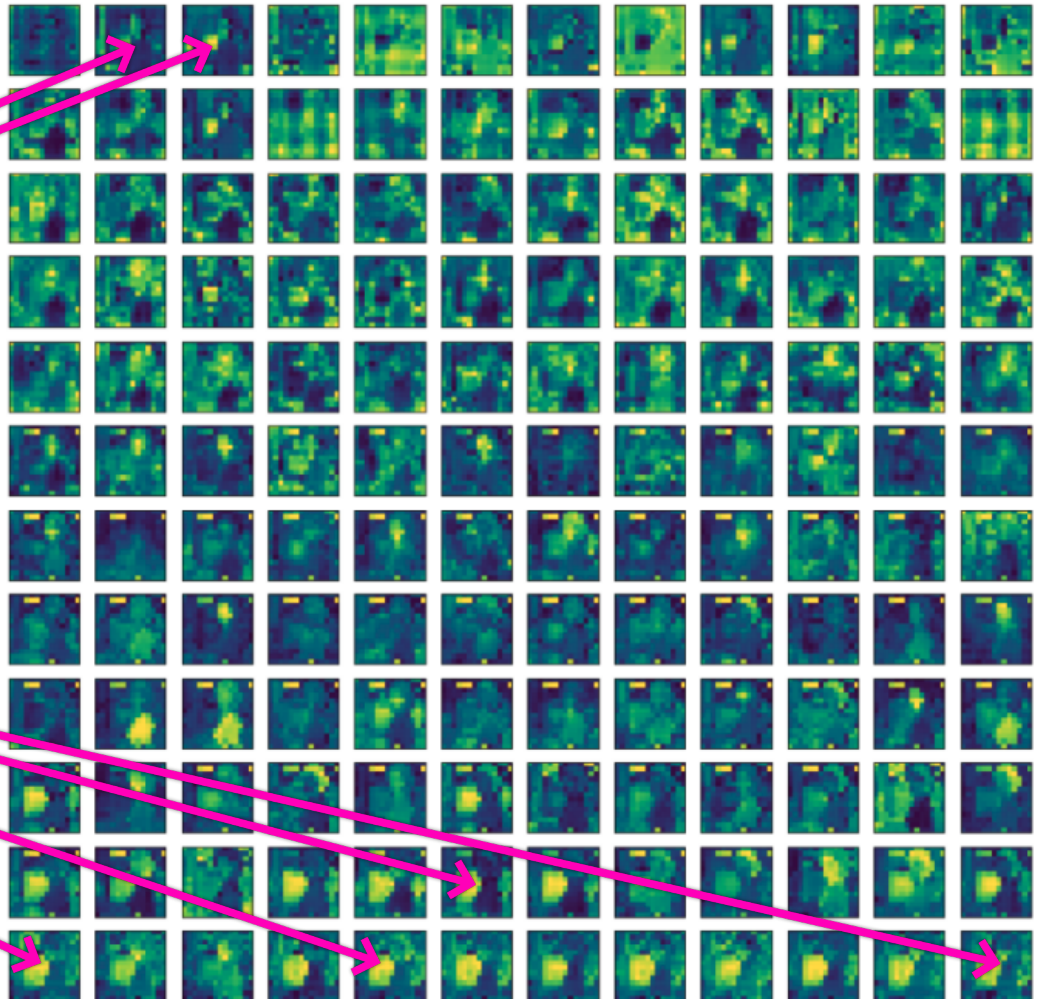
Experiments: head projection

Original image (“television”)



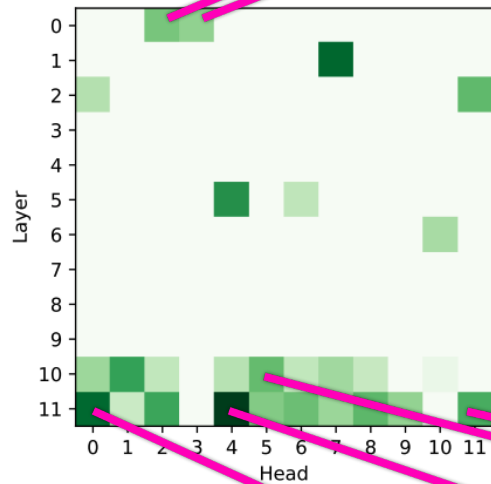
Experiments: head projection

Original image ("television")



CIFAR-100

Head coefficients



Conclusions

- SMaT is a framework that optimizes explanations for teaching students
 - SMaT leads to **high simulability**
 - SMaT learns **plausible explanations**
- We hope this work motivates the interpretability community to consider **scaffolding** as valuable criterion **for evaluating and designing new methods**



(paper) arxiv.org/abs/2204.10810

(code) github.com/CoderPat/learning-scaffold

Introduction

- **Simulability** is particularly appealing for evaluating explanations
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective

Introduction

- **Simulability** is particularly appealing for evaluating explanations
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective
- Pruthi et al. (2021) proposed a framework for measuring simulability that

Introduction

- **Simulability** is particularly appealing for evaluating explanations
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective
- Pruthi et al. (2021) proposed a framework for measuring simulability that
 - ★ disregards **trivial protocols**

Introduction

- **Simulability** is particularly appealing for evaluating explanations
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective
- Pruthi et al. (2021) proposed a framework for measuring simulability that

★ disregards **trivial protocols**



punctuation symbols ⇒ positive
stop words ⇒ negative

Introduction

- **Simulability** is particularly appealing for evaluating explanations
 - ✓ aligns with the goal of communicating the underlying model behavior
 - ✓ is easily measurable (both manually and automatically)
 - ✓ puts all explainability methods under a single perspective
- Pruthi et al. (2021) proposed a framework for measuring simulability that
 - ★ disregards **trivial protocols** →
 - 🏀 requires **an optimization procedure**

punctuation symbols ⇒ positive
stop words ⇒ negative