# Robustness and Compositionality in ML with Insights from Cognitive Bottlenecks

PhD thesis defense of Ankit Vani







#### Collaborators



Aaron Courville



Aishwarya Agrawal



Alessandro Sordoni



Arian Hosseini



Bac Nguyen



Tsirigotis



Dzmitry Bahdanau



Eeshan Dhekane



Frederick Tung



Gabriel L. Oliveira



Hattie Zhou



Hossein Sharifi-Noghabi



Hugo

Jose Gallego



Max Schwarzer



Michael Noukhovitch



Ranjay Krishna







Sarvjeet Singh Ghotra



Larochelle

Simon Lacoste-Julien



Yuchen Lu





#### **Publications**

GAIT: A Geometric Approach to Information Theory. Jose Gallego, Ankit Vani, Max Schwarzer, Simon Lacoste-Julien. AISTATS 2020.

Iterated Learning for Emergent Systematicity in VQA. Ankit Vani, Max Schwarzer, Yuchen Lu, Eeshan Dhekane, Aaron Courville. ICLR 2021.

Fortuitous Forgetting in Connectionist Networks. Hattie Zhou, Ankit Vani, Hugo Larochelle, Aaron Courville. ICLR 2022.

On the Compositional Generalization Gap of In-Context Learning. Arian Hosseini, Ankit Vani, Dzmitry Bahdanau, Alessandro Sordoni, Aaron Courville. BlackboxNLP Workshop 2022.

Simplicial Embeddings in Self-Supervised Learning and Downstream Classification. Samuel Lavoie, Christos Tsirigotis, Max Schwarzer, Ankit Vani, Michael Noukhovitch, Kenji Kawaguchi, Aaron Courville. ICLR 2023.

Forget Sharpness: Perturbed Forgetting of Model Biases Within SAM Dynamics. Ankit Vani, Frederick Tung, Gabriel L. Oliveira, Hossein Sharifi-Noghabi. ICML 2024.

SPARO: Selective Attention for Robust and Compositional Transformer Encodings for Vision. Ankit Vani, Bac Nguyen, Samuel Lavoie, Ranjay Krishna, Aaron Courville. ECCV 2024.



Article 1

#### **Publications**

GAIT: A Geometric Approach to Information Theory. Jose Gallego, Ankit Vani, Max Schwarzer, Simon Lacoste-Julien. AISTATS 2020.

Iterated Learning for Emergent Systematicity in VQA. Ankit Vani, Max Schwarzer, Yuchen Lu, Eeshan Dhekane, Aaron Courville. ICLR 2021.

Fortuitous Forgetting in Connectionist Networks. Hattie Zhou, Ankit Vani, Hugo Larochelle, Aaron Courville. ICLR 2022.

On the Compositional Generalization Gap of In-Context Learning. Arian Hosseini, Ankit Vani, Dzmitry Bahdanau, Alessandro Sordoni, Aaron Courville. BlackboxNLP Workshop 2022.

Simplicial Embeddings in Self-Supervised Learning and Downstream Classification. Samuel Lavoie, Christos Tsirigotis, Max Schwarzer, Ankit Vani, Michael Noukhovitch, Kenji Kawaguchi, Aaron Courville. ICLR 2023.

Forget Sharpness: Perturbed Forgetting of Model Biases Within SAM Dynamics. Ankit Vani, Frederick Tung, Gabriel L. Oliveira, Hossein Sharifi-Noghabi. ICML 2024.	Article 2
SPARO: Selective Attention for Robust and Compositional Transformer Encodings for Vision. Ankit Vani, Bac Nguyen, Samuel Lavoie, Ranjay Krishna, Aaron Courville. ECCV 2024.	Article 3





### Introduction

Humans are pretty good at Systematic Generalization

Humans are pretty good at Systematic Generalization

- If you know what "a cat chases a mouse" means, no trouble understanding what "a Roomba chases an elephant" means
  - $\circ$   $\quad$  Even if you have never encountered the phrase before
- Systematic generalization<sup>[1]</sup> is the consistent application of learned rules to unseen situations

Humans are pretty good at Systematic Generalization

- If you know what "a cat chases a mouse" means, no trouble understanding what "a Roomba chases an elephant" means
  - Even if you have never encountered the phrase before
- Systematic generalization<sup>[1]</sup> is the consistent application of learned rules to unseen situations

Humans are pretty good at *Compositional Generalization* 

If you understand "red sphere" and "blue cube," you understand "red cube"



But machine learning models can be pretty bad at behaving systematically and understanding compositionality



But machine learning models can be pretty bad at behaving systematically and understanding compositionality

- Standard i.i.d. assumption in machine learning: Test data comes from the same distribution as the training data
  - Possible catastrophic drops in performance when violated
- For a model capable of systematic generalization
  - Instead of performance depending on divergence between training and test distributions,
  - Performance depends on the discrepancy in the mechanistic processes generating the training and testing observations



## How a model behaves on samples outside the training distribution depends on its *inductive biases*



How a model behaves on samples outside the training distribution depends on its *inductive biases* 

Inductive biases: assumptions by a model to generalize to unseen data
 Shaped by network architecture, training objective, regularization, optimization, etc.



How a model behaves on samples outside the training distribution depends on its *inductive biases* 

Inductive biases: assumptions by a model to generalize to unseen data
 Shaped by network architecture, training objective, regularization, optimization, etc.

What are the inductive biases behind human learning that allows us to systematically generalize?



#### Cognitive bottlenecks

• Overwhelming abundance of stimuli around us, and human mental resources are limited

- Our cognitive processes impose bottlenecks for what information is stored, retained, and consciously attended
  - The inductive biases imparted by these bottlenecks are crucial for systematicity



This thesis explores how inductive biases in machine learning can be informed and developed through insights from cognitive bottlenecks behind human learning and knowledge representation

Cognitive bottlenecks we will discuss:

- 1. Iterated learning
- 2. Forgetting and relearning
- 3. Selective attention









(Based on [1,2])

[1] Kirby et al. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. National Academy of Sciences (2008).
 [2] Kirby et al. Iterated learning and the evolution of language. Current opinion in neurobiology 28 (2014).





(Based on [1,2])

[1] Kirby et al. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. National Academy of Sciences (2008).
 [2] Kirby et al. Iterated learning and the evolution of language. Current opinion in neurobiology 28 (2014).





[1] Kirby et al. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. National Academy of Sciences (2008).
[2] Kirby et al. Iterated learning and the evolution of language. Current opinion in neurobiology 28 (2014).





[1] Kirby et al. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. National Academy of Sciences (2008).[2] Kirby et al. Iterated learning and the evolution of language. Current opinion in neurobiology 28 (2014).



#### Learning bottleneck

- Need to learn an highly expressive language
  - Through limited supervision
- Language properties likely to pass through become universal
  - Compositional rules are more likely to survive the transmission <sup>[1]</sup>
    - Easier to learn
    - Faster to learn



## Article 1

#### Iterated Learning for Emergent Systematicity in VQA



International Conference on Learning Representations (ICLR), 2021



### Iterated learning in machine learning

- In 2020, iterated learning (IL) in ML stayed close to its Cog Sci roots
  - Most of the work involved agents playing very simplistic referential games <sup>[1,2,3,4,5]</sup>

- Li and Bowling. Ease-of-teaching and language structure from emergent communication. NeurIPS 2019.
- [2] Ren et al. Compositional languages emerge in a neural iterated learning model. ICLR 2020.
- [3] Dagan et al. Co-evolution of language and agents in referential games. EACL 2021.
- [4] Cogswell et al. Emergence of compositional language with deep generational transmission. arXiv:1904.09067 (2019).
- [5] Guo et al. The emergence of compositional languages for numeric concepts through iterated learning in neural agents. arXiv:1910.05291 (2019).
- [6] Zheng et al. Iterated learning improves compositionality in large vision-language models. CVPR 2024.



24

## Iterated learning in machine learning

- In 2020, iterated learning (IL) in ML stayed close to its Cog Sci roots
  - Most of the work involved agents playing very simplistic referential games <sup>[1,2,3,4,5]</sup>
- Claim: Learning bottleneck is a fundamental way to recover structure
- We demonstrated broader applicability through the more complex task of visual question-answering (VQA)
  - Since then, IL has been applied to improve compositional generalization in larger and real-world settings, *e.g.*, training vision-language models like CLIP <sup>[6]</sup>

- [1] Li and Bowling. *Ease-of-teaching and language structure from emergent communication*. NeurIPS 2019.
- [2] Ren et al. Compositional languages emerge in a neural iterated learning model. ICLR 2020.
- [3] Dagan et al. Co-evolution of language and agents in referential games. EACL 2021.
- [4] Cogswell et al. Emergence of compositional language with deep generational transmission. arXiv:1904.09067 (2019).
- [5] Guo et al. The emergence of compositional languages for numeric concepts through iterated learning in neural agents. arXiv:1910.05291 (2019).
- [6] Zheng et al. Iterated learning improves compositionality in large vision-language models. CVPR 2024.



### Visual question-answering (VQA)



#### Q: Is a green shape left of a square?



#### Visual question-answering (VQA)



#### Q: Is a green shape left of a square? A: Yes



#### Visual question-answering (VQA)







## Neural module networks (NMNs)<sup>[1,2]</sup>



[1] Andreas et al. Neural module networks. CVPR 2016.

[2] Johnson et al. Inferring and executing programs for visual reasoning. CVPR 2017.



#### Core idea

- With the right layout, NMNs exhibit compositionality <sup>[1]</sup>; other methods fail
- But: The right layout does not emerge naturally
  - Bahdanau et al., 2019<sup>[1]</sup> needed to provide correct tree-structured layouts
  - Learned layouts only converged to be robust under a strong prior for the correct structure



#### Core idea

- With the right layout, NMNs exhibit compositionality <sup>[1]</sup>; other methods fail
- But: The right layout does not emerge naturally
  - Bahdanau et al., 2019<sup>[1]</sup> needed to provide correct tree-structured layouts
  - Learned layouts only converged to be robust under a strong prior for the correct structure

Use IL to encourage structured layouts towards systematic generalization



Article 1 Method



#### Standard NMN training





#### Standard NMN training

Viewed as two agents with an emergent language to solve the task





#### Standard NMN training

Viewed as two agents with an emergent language to solve the task





#### Iterated learning for NMNs





#### Iterated learning for NMNs




#### Iterated learning for NMNs



![](_page_37_Picture_0.jpeg)

#### Iterated learning for NMNs

![](_page_37_Figure_2.jpeg)

![](_page_38_Picture_0.jpeg)

### Learning bottleneck

Limit the length of the learning phases: Early stopping

• Generally a sweet spot for the number of steps

Program generator learning phase steps

- *Too few:* Low confidence in utterances, high variance
- Too many: Overfitting to transmitted data

Execution engine learning phase steps

- *Too few:* High variance gradients at the start of interacting phase
- *Too many:* Overfitting to an imperfect program generator

![](_page_38_Figure_10.jpeg)

![](_page_38_Figure_11.jpeg)

![](_page_39_Picture_0.jpeg)

#### Article 1 Results with CLEVR

![](_page_40_Picture_0.jpeg)

#### CLEVR/CLOSURE example

![](_page_40_Picture_2.jpeg)

<sup>[1]</sup>**Q1 (CLEVR):** There is another cube that is the same size as the brown cube; what is its color?

<sup>11</sup>Q2 (CLEVR): There is a thing that is in front of the yellow thing; does it have the same color as cylinder? <sup>21</sup>Q3 (CLOSURE): There is another rubber object that is the

<sup>[2]</sup>Q3 (CLOSURE): There is another rubber object that is the same size as the gray cylinder; does it have the same color as the tiny shiny block?

(Figure taken from [2])

![](_page_41_Picture_0.jpeg)

42

### CLEVR with 100 GT programs

![](_page_41_Figure_2.jpeg)

**Tensor-NMN:** modules are residual blocks<sup>[1]</sup> with separate parameters, from [2] **Vector-NMN:** modules share parameters with separate FiLM<sup>[3]</sup> adapters, from [4]

- All models reach similar training accuracies
- IL leads to significantly higher program accuracy
- The generalization difference is more apparent in the OOD CLOSURE dataset

![](_page_42_Picture_0.jpeg)

## CLOSURE categories: 100 GT programs

Evaluation set	Tensor	r-NMN	Vector-NMN		
	Without IL	With IL	Without IL	With IL	
CLEVR-Val	$0.912 \pm 0.07$	$0.964 \pm 0.01$	$0.960\pm0.01$	$\boldsymbol{0.964 \pm 0.00}$	
and_mat_spa	$\boldsymbol{0.278 \pm 0.17}$	$0.264 \pm 0.16$	$0.400 \pm 0.13$	$0.335 \pm 0.18$	
or_mat	$0.327 \pm 0.11$	$0.481 \pm 0.24$	$0.367 \pm 0.11$	$0.563 \pm 0.23$	
or_mat_spa	$0.286 \pm 0.13$	$0.405 \pm 0.22$	$0.330 \pm 0.11$	$0.444 \pm 0.24$	
compare_mat	$0.793 \pm 0.11$	$0.851 \pm 0.17$	$0.660\pm0.16$	$0.873 \pm 0.12$	
compare_mat_spa	$0.746 \pm 0.13$	$0.853 \pm 0.15$	$0.677 \pm 0.14$	$0.871 \pm 0.12$	
embed_spa_mat	$0.824 \pm 0.07$	$0.947 \pm 0.03$	$0.863 \pm 0.07$	$\boldsymbol{0.900 \pm 0.08}$	
embed_mat_spa	$0.739 \pm 0.14$	$0.941 \pm 0.02$	$0.894 \pm 0.03$	$0.936 \pm 0.03$	

- IL improves performance on all but one CLOSURE category
- Tensor-NMN with IL leads to CLEVR accuracy similar to previous works with far fewer programs
  - 18,000 for Johnson et al., 2017<sup>[1]</sup>, 1,000 for Vedantam et al., 2019<sup>[2]</sup>

<sup>[1]</sup> Johnson et al. Inferring and executing programs for visual reasoning. CVPR 2017.

<sup>[2]</sup> Vedantam et al. Probabilistic neural symbolic models for interpretable visual question answering. ICML 2019.

![](_page_43_Picture_0.jpeg)

Article 1 Conclusion

![](_page_44_Picture_0.jpeg)

### Key takeaways

Iterated learning is more broadly applicable where hard-to-learn compositional latents are desired, beyond simple referential games

Iterated learning amplifies compositionality when a preference for it exists

- In addition to model architecture, this bias can be imparted through few-shot supervision
- Emergent compositional structure can improve downstream systematic generalization
  - We find similar gains in simpler (SHAPES) and larger real-world image (GQA) settings

![](_page_45_Picture_0.jpeg)

![](_page_45_Figure_1.jpeg)

# Forgetting and Relearning

## Forgetting and relearning in humans

- Forgetting can be frustrating, but it is beneficial
- Spacing effect<sup>[1]</sup>
  - More effective learning with spaced-out study sessions
- Forgetting allows integration of new perspectives and knowledge unburdened by biases of the past

![](_page_46_Figure_5.jpeg)

i ∰Mila

Universit

## Forgetting and relearning in humans

- Forgetting can be frustrating, but it is beneficial
- Spacing effect<sup>[1]</sup>
  - More effective learning with spaced-out study sessions
- Forgetting allows integration of new perspectives and knowledge unburdened by biases of the past

![](_page_47_Figure_5.jpeg)

Univers

Under continuous forgetting and relearning:

Less relevant and specialized information fades away, while more pertinent and structurally coherent knowledge is retained

<sup>(</sup>Figure based on [1])

![](_page_48_Picture_0.jpeg)

## Forgetting and relearning in ML

- The forget-and-relearn<sup>[1]</sup> paradigm typically partially resets weights to "forget" before continuing training for "relearning"
- Iterated learning can be seen as an instance of forget-and-relearn
- Iterative resetting counters the primacy bias in reinforcement learning <sup>[2]</sup>

![](_page_48_Figure_5.jpeg)

(Figure taken from [2])

![](_page_49_Picture_0.jpeg)

### Article 2

#### Forget Sharpness: Perturbed Forgetting of Model Biases Within SAM Dynamics

![](_page_49_Picture_3.jpeg)

International Conference on Machine Learning (ICML), 2024

![](_page_50_Picture_0.jpeg)

## Sharpness-Aware Minimization (SAM)

Assumption: Flatter loss surface regions generalize better

SAM<sup>[1]</sup> is motivated to guide optimization to flatter regions of the loss surface

- Reduce sharpness by approximately minimizing a PAC-Bayes upper bound
- Algorithm:
  - Perturb by maximizing loss
  - Compute minimization gradient
  - Step from unperturbed weights

![](_page_50_Figure_9.jpeg)

![](_page_51_Picture_0.jpeg)

### Unresolved concerns with SAM

Improved generalization with SAM deviates from its motivating theory

- The PAC-Bayes bound does not sufficiently explain SAM's benefits <sup>[1]</sup>
  - Bound derived for random perturbations, which perform worse in practice
  - Steepest ascent perturbations loosen the bound but perform better
  - *m*-Sharpness: Why do smaller perturbing batches perform better?
- Empirically, flatness and generalization do not correlate in general <sup>[2]</sup>

[1] Andriushchenko et al. Towards understanding sharpness-aware minimization. ICML 2022.

[2] Andriushchenko et al. A modern look at the relationship between sharpness and generalization. ICML 2023.

![](_page_52_Picture_0.jpeg)

Article 2 Method

![](_page_53_Picture_0.jpeg)

## "Perturbed forgetting" perspective

- SAM dynamics perform forgetting and relearning without erasing learned state
- Perturbations can discard undesirable model biases

Steepest ascent gradients for a small set of *m* samples can reveal undesirable shortcuts the model learned for them

- Perturbing discards them
- Predictions then utilize the global learned structure
- Gradients strengthen this global structure

![](_page_54_Picture_0.jpeg)

### Information-theoretic argument

Generalization gap bound from [1]:  $\Delta(\boldsymbol{\theta}) \in \tilde{O}\left(\sqrt{\frac{I(\boldsymbol{X}; \hat{\boldsymbol{Y}} \mid Y) + I(\boldsymbol{\theta}; D)}{n}}\right)$ 

#### Avoid increasing $I(\boldsymbol{\theta}; D)$

- Minimizing the loss makes it easier to identify D as the training dataset
- Perturbing does not minimize loss

Decrease  $I(\boldsymbol{X}; \hat{\boldsymbol{Y}} \mid Y)$ 

- Model biases can be exposed via outputs  $\hat{Y}$
- Perturbing can target and discard them

#### Wila Université de Montré

# Output bias forgetting (OBF)

Gradient of cross-entropy (CE) with logits *z*:

$$\nabla_{\boldsymbol{\theta}} L^{\text{CE}}(\boldsymbol{y}, \boldsymbol{\hat{y}}) = \mathbb{E}_{i \sim \boldsymbol{\hat{y}}} \left[ \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_i \right] - \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_{\boldsymbol{y}}$$

Minimizing (for learning) pushes down more on worse predictions 🗸

Maximizing (for perturbing) pulls up more on worse predictions

Forgetting some biases but amplifying others 😣

# Output bias forgetting (OBF)

Gradient of cross-entropy (CE) with logits z:

 $\alpha = 0$ 

$$\nabla_{\boldsymbol{\theta}} L^{\text{CE}}(\boldsymbol{y}, \boldsymbol{\hat{y}}) = \mathbb{E}_{i \sim \boldsymbol{\hat{y}}} \left[ \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_i \right] - \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_{\boldsymbol{y}}$$

Minimizing (for learning) pushes down more on worse predictions 🗸

Maximizing (for perturbing) pulls up more on worse predictions

Forgetting some biases but amplifying others 🔀

The OBF perturbation minimizes target likelihood without bias amplification

$$\nabla_{\boldsymbol{\theta}} L^{\mathrm{BF}}(\boldsymbol{y}, \boldsymbol{\hat{y}}) = \mathbb{E}_{i \sim \mathrm{Uniform}} \left[ \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_i \right] - \left( \alpha \mathbb{E}_{i \sim \boldsymbol{\hat{y}}} \left[ \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_i \right] + (1 - \alpha) \nabla_{\boldsymbol{\theta}} \boldsymbol{z}_y \right)$$

Like cross-entropy, without bias amplification

Negative cross-entropy with uniform targets

![](_page_57_Picture_0.jpeg)

Article 2 Results

![](_page_58_Picture_0.jpeg)

## Forgetting vs. generalization

- Train ViT-S/32 models on CIFAR-10 with perturbing batch sizes  $m \in \{2^k \mid k \in \{0, \dots, 9\}\}$
- Discretize model outputs with varying thresholds

![](_page_58_Figure_4.jpeg)

There exist thresholds where generalization correlates more strongly with forgetting than flatness of the loss surface.

![](_page_59_Picture_0.jpeg)

#### OBF vs. standard SAM variants

Model	Method	Perturb -	ImageNet-				Charmanaa	
			V1	Real	V2	R	Sketch	Snarpness
ViT-S/32	AdamW	None	$69.29{\scriptstyle \pm 0.26}$	$75.31{\scriptstyle \pm 0.28}$	$55.48{\scriptstyle \pm 0.58}$	$19.02{\scriptstyle \pm 0.47}$	$16.38{\scriptstyle \pm 0.34}$	165.6±15.2
	SAM <sup>[1]</sup>	Standard OBF	$72.77{\scriptstyle\pm0.06}\\74.49{\scriptstyle\pm0.04}$	$78.89{\scriptstyle \pm 0.05} \\ 81.31{\scriptstyle \pm 0.05}$	58.81±0.33 61.13±0.18	21.63±0.23 <b>25.31</b> ±0.41	$\begin{array}{c} 19.68 {\scriptstyle \pm 0.50} \\ \textbf{22.58} {\scriptstyle \pm 0.13} \end{array}$	$\frac{14.9_{\pm 1.1}}{3.9_{\pm 1.4}}$
	GSAM <sup>[2]</sup>	Standard OBF	73.41±0.05 74.41±0.12	79.48±0.08 <b>81.41</b> ±0.11	59.94±0.15 <b>61.08</b> ±0.18	$22.18{\scriptstyle \pm 0.15} \\ 25.15{\scriptstyle \pm 0.23}$	$20.28{\scriptstyle \pm 0.15} \\ 22.24{\scriptstyle \pm 0.07}$	11.6±1.2 <b>3.1</b> ±0.7
	ASAM <sup>[3]</sup>	Standard OBF	74.45±0.11 <b>74.73</b> ±0.19	$81.23{\scriptstyle \pm 0.11} \\ 81.24{\scriptstyle \pm 0.25}$	${}^{60.78 \pm 0.25}_{60.95 \pm 0.28}$	$24.07{\scriptstyle\pm 0.12}\\24.65{\scriptstyle\pm 0.26}$	$21.68{\scriptstyle \pm 0.23}\\22.40{\scriptstyle \pm 0.10}$	$6.5{\scriptstyle \pm 0.4}\ 30.3{\scriptstyle \pm 11.6}$
ResNet-50	SGD	None	$76.86{\scriptstyle \pm 0.07}$	83.28±0.11	$65.00{\scriptstyle \pm 0.14}$	$20.29{\scriptstyle \pm 0.36}$	$20.53{\scriptstyle \pm 0.46}$	230.4±42.7
	SAM <sup>[1]</sup>	Standard OBF	77.49±0.06	$\begin{array}{c} 83.78 \scriptstyle \pm 0.05 \\ 84.01 \scriptstyle \pm 0.03 \end{array}$	$\begin{array}{c} 65.26 \scriptstyle \pm 0.21 \\ 65.70 \scriptstyle \pm 0.45 \end{array}$	$21.08{\scriptstyle \pm 0.16} \\ 21.63{\scriptstyle \pm 0.18}$	$21.18{\scriptstyle \pm 0.32} \\ 22.17{\scriptstyle \pm 0.26}$	170.1±18.9 164.4±25.0
	GSAM <sup>[2]</sup>	Standard OBF	$77.43{\scriptstyle \pm 0.12} \\ 77.66{\scriptstyle \pm 0.08}$	$83.79{\scriptstyle \pm 0.19} \\ 84.09{\scriptstyle \pm 0.07}$	$\begin{array}{c} 65.37 \scriptstyle \pm 0.26 \\ 66.01 \scriptstyle \pm 0.09 \end{array}$	$21.37{\scriptstyle\pm0.21}\\21.76{\scriptstyle\pm0.23}$	$21.52{\scriptstyle \pm 0.56}\\22.26{\scriptstyle \pm 0.47}$	171.0±16.8 161.4±10.9
	ASAM <sup>[3]</sup>	Standard OBF	77.30±0.02	84.07±0.03 <b>84.66</b> ±0.05	65.55±0.16 66.55±0.15	21.71±0.02 23.84±0.12	21.75±0.15 <b>24.21</b> ±0.42	<b>33.6</b> ± <b>2.99</b> 39.1±1.28

[1] Foret et al. Sharpness-aware minimization for efficiently improving generalization. ICLR 2021.

[2] Zhuang et al. Surrogate gap minimization improves sharpness-aware training. ICLR 2022.

[3] Kwon et al. ASAM: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. ICML 2021.

![](_page_60_Picture_0.jpeg)

### Article 2 Conclusion

![](_page_61_Picture_0.jpeg)

## Key takeaways

Targeted forgetting and relearning can be a powerful driver of generalization

- Gradient-based approaches are effective at targeting undesirable biases
- Perturbed forgetting protects the global model state against suboptimal forgetting steps

SAM's training dynamics are more important than loss surface flatness

• Is the pursuit of flatter minima misleading?

![](_page_62_Picture_0.jpeg)

![](_page_62_Figure_1.jpeg)

## Selective Attention

### Selective attention in humans

- Human working memory is extremely limited
  - It can can hold at most ~4 "chunks" of information at a time [1]
- Selective attention is critical in managing these limited mental resources
  - We prioritize stimuli that is most relevant, and ignore the rest <sup>[2]</sup>
  - Humans struggle to simultaneously attend to separable features (color, shape, orientation) <sup>[3]</sup>, or novel objects at the same location <sup>[4]</sup>

[1] Cowan. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. Behavioral and brain sciences, 2001.

[2] Treisman. Strategies and models of selective attention. Psychological review, 1969.

[4] Rock and Gutman. *The effect of inattention on form perception*. Journal of Experimental Psychology: Human Perception and Performance, 1981.

<sup>[3]</sup> Treisman and Gelade. A feature-integration theory of attention. Cognitive psychology, 1980.

### Selective attention in humans

- Human working memory is extremely limited
  - It can can hold at most ~4 "chunks" of information at a time [1]
- Selective attention is critical in managing these limited mental resources
  - We prioritize stimuli that is most relevant, and ignore the rest <sup>[2]</sup>
  - Humans struggle to simultaneously attend to separable features (color, shape, orientation) <sup>[3]</sup>, or novel objects at the same location <sup>[4]</sup>

A preference for learning separable concepts eases learning compositional representations of our knowledge, which is essential to act effectively with limited cognitive resources in the world

[1] Cowan. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. Behavioral and brain sciences, 2001.

[2] Treisman. Strategies and models of selective attention. Psychological review, 1969.

[3] Treisman and Gelade. A feature-integration theory of attention. Cognitive psychology, 1980.

[4] Rock and Gutman. *The effect of inattention on form perception*. Journal of Experimental Psychology: Human Perception and Performance, 1981.

![](_page_65_Picture_0.jpeg)

## Selective attention in machine learning

- The attention mechanism<sup>[1]</sup> is closely related to selective attention
  Refined and popularized by the Transformer<sup>[2]</sup> architecture
- In transformers, selective attention is implicitly performed for hierarchical construction of internal representations
- But when explicitly learning representations for future downstream tasks, a prior for separately-attendable concepts is not automatically granted
  - Deprives downstream tasks easy access to task-relevant aspects of the data's structure

[1] Bahdanau et al. Neural machine translation by jointly learning to align and translate. ICLR 2015.

[2] Vaswani et al. Attention is all you need. NeurIPS 2017.

![](_page_66_Picture_0.jpeg)

## Article 3

#### SPARO: Selective Attention for Robust and Compositional Transformer Encodings for Vision

![](_page_66_Picture_3.jpeg)

European Conference on Computer Vision (ECCV), 2024

![](_page_67_Picture_0.jpeg)

![](_page_67_Picture_1.jpeg)

![](_page_68_Picture_0.jpeg)

![](_page_68_Picture_1.jpeg)

![](_page_69_Picture_0.jpeg)

![](_page_69_Figure_1.jpeg)

![](_page_69_Picture_2.jpeg)

Describe the building architecture.

![](_page_70_Picture_0.jpeg)

![](_page_70_Picture_1.jpeg)

![](_page_70_Picture_2.jpeg)

Describe the building architecture.

![](_page_71_Picture_0.jpeg)

![](_page_71_Picture_1.jpeg)

![](_page_71_Picture_2.jpeg)

![](_page_71_Picture_3.jpeg)

Describe the building architecture.

Are people arriving or leaving?








Describe the building architecture.









Describe the building architecture.











Describe the building architecture.











Describe the building architecture.











Describe the building architecture.





## SPARO (Separate-head Attention Read-Out)

SPARO is a read-out mechanism for transformers that explicitly structures encodings as a collection of separately-attended concepts.





# Article 3 SPARO with CLIP



## CLIP<sup>[1]</sup> (Contrastive Language-Image Pre-training)



(Figure taken from [1])



## SPARO with CLIP



SPARO imposes a prior that both modalities share the same compositional world with the same attendable concepts



#### Article 3 CLIP Experiments



#### Zero-shot and linear probe accuracies

Train	Model	ImageNet-					Object	Sugar
		V1	V2	Sketch	Α	R	Net	Crepe
CC3M	$\mathrm{CLIP}^{16}(\mathcal{C})$	0.141	0.122	0.068	0.033	0.177	0.080	0.611
	$\mathcal{C}\mathrm{+GAP}$	0.156	0.134	0.069	0.033	0.187	0.081	0.616
	$\mathcal{C}{+}\mathrm{Sparo}$	0.170	0.140	0.088	0.035	0.221	0.098	0.625
CC12M	$\operatorname{CLIP}^{16}(\mathcal{C})$	0.361	0.311	0.249	0.091	0.467	0.218	0.697
	$\mathcal{C}\mathrm{+GAP}$	0.382	0.330	0.262	0.101	0.501	0.241	0.695
	$\mathcal{C}{+}\mathrm{Sparo}$	0.406	0.350	0.298	0.113	0.559	0.268	0.723
CC15M	$\mathrm{CLIP}^{16}$ (C)	0.384	0.337	0.268	0.105	0.503	0.238	0.699
	$\mathcal{C}\mathrm{+GAP}$	0.399	0.343	0.287	0.114	0.531	0.252	0.701
	$\mathcal{C}+\mathrm{Sparo}$	0.437	0.378	0.317	0.145	0.579	0.279	0.730
L400M	$\mathrm{CLIP}^{32}$ (C)	0.617	0.531	0.482	0.202	0.719	0.423	0.748
	$\mathcal{C}\mathrm{+GAP}$	0.623	0.537	0.492	0.212	0.725	0.440	0.732
	$\mathcal{C}{+}\mathrm{Sparo}$	0.635	0.552	0.507	0.231	0.747	0.459	0.770

NT- J-1	ImageNet linear probe								
Model	Train:	CC3M	CC12M	CC15M	L400M				
CLIP	$(\mathcal{C})$	0.469	0.630	0.646	0.743				
$\mathcal{C}+\mathrm{GA}$	P	0.504	0.649	0.664	0.747				
$\mathcal{C}+\mathrm{Sp}$	ARO	0.561	0.700	0.711	0.755				



#### Post-hoc concept selection

Pretrained SPARO representation





#### Post-hoc concept selection





#### Post-hoc concept selection



Intervening to select a subset of slots with the highest ImageNet zero-shot performance



#### Visualizations



A herd of **cattle** walking down a road being followed by a cowboy



A herd of cattle **walking** down a road being followed by a cowboy



A herd of cattle walking down a road being followed by a cowboy



Several surfboards standing in a row on the beach



Several surfboards **standing** in a row on the beach



Several surfboards standing in a row on the **beach** 



Two people riding **horses** on a rock path



Two people **riding** horses on a rock path



Two people riding horses on a rock path



A **man** sitting alone on a park bench in a park



A man sitting alone on a park bench in a park



A man sitting alone on a park bench in a **park** 



pa

87



Article 3 Conclusion



#### Key takeaways

# Attention offers a scalable solution to impart a stronger bias for compositionality

- Constructing encodings as slots produced by bottlenecks of separate attention heads encourages learning of data variations in terms of concepts represented by the heads
   In CLIP, adds an additional prior for the modalities to share the same set of concepts
- Models trained with selective attention constraints achieve better downstream generalization



# Conclusion of the presentation



#### Discussion

- Caveats in taking insights from cognitive science
  - Human learning utilizes a vast repository of prior knowledge, experience, and is intertwined with motivation, curiosity, and emotions
  - A large number of complex cognitive phenomena occur simultaneously and cognitive science models often simplify processes to make them tractable for study
  - Directly translating these simplified models into machine learning risks overlooking the nuanced nature of human cognition and the fundamental differences in optimization
- Humans do not always systematically generalize
  - Ability to perform complex, systematic logical inferences often requires thorough practice
  - Machine learning models, however, often fail at tasks human common sense succeeds at
- Systematic generalization is not sufficient in the real world
  - Models can systematically generalize on systemic biases



#### Open questions

- Iterated learning (IL)
  - Can IL help when there is no clear preference for a beneficial structure?
  - Can we dynamically adapt the learning bottleneck during training?
  - Can we distill IL's dynamics into a method that does not require a speaker-listener setup?
- Perturbed forgetting and SAM
  - Can we separate perturbed forgetting into a more general framework beyond SAM?
  - In the perturbed forgetting perspective of SAM, what really are the undesirable biases?
  - Is there a relationship between perturbed forgetting dynamics and loss surface flatness?
- SPARO
  - Can IL amplify Sparo's preference for compositionality for cleaner disentanglement?
  - Can we guide the type of concepts SPARO should prioritize learning?
  - Can Sparo offer insights for improving the internal attention mechanisms of transformers?



#### Thank you!



Aaron Courville



Aishwarya Agrawal



Alessandro Sordoni



Arian Hosseini



Christos Bac Nguyen Tsirigotis



Dzmitry Bahdanau



Eeshan Dhekane



Frederick Tung



Gabriel L. Oliveira



Hattie Zhou



Hossein Sharifi-Noghabi



Larochelle

Jose Gallego



Max Schwarzer



Michael Noukhovitch



Ranjay Krishna



Samuel Lavoie



Sarvjeet Singh Ghotra



Simon Lacoste-Julien



Yuchen Lu



