

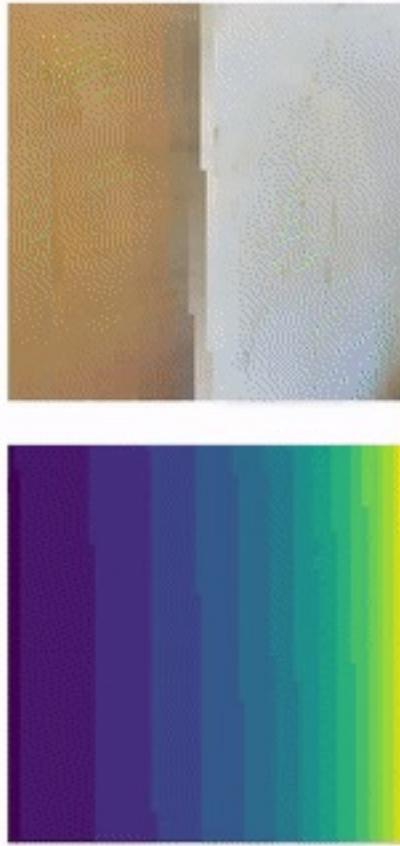
# Visualizing, evaluating and transferring reasoning patterns in VQA,

Christian Wolf  
September 29<sup>th</sup>, 2021



*The group in Feb. 2020: Corentin Kervadec, Steeven Janny, Edward Beeching, Fabien Baradel, Théo Jaunet, Quentin Possamai.*

# Topological neural maps



[Beeching, Dibangoye, Simonin, Wolf, ECCV 2020]

# Self-supervision for Deep-RL

1<sup>st</sup> ranked at Multi-ON Challenge (CVPR 2021)



Rank	Team	Progress	PPL	Success	SPL
1	Lyon	67	44	55	35
2	SGoLAM	64	38	52	32
3	VIMP	57	36	41	26

**1<sup>st</sup> place: Team Lyon,** Pierre Marza, Laetitia Matignon, Olivier Simonin, Christian Wolf, Learning mapping and spatial reasoning with auxiliary tasks.

# Navigation in real environments



**LABS**  
NAVER LABS EUROPE



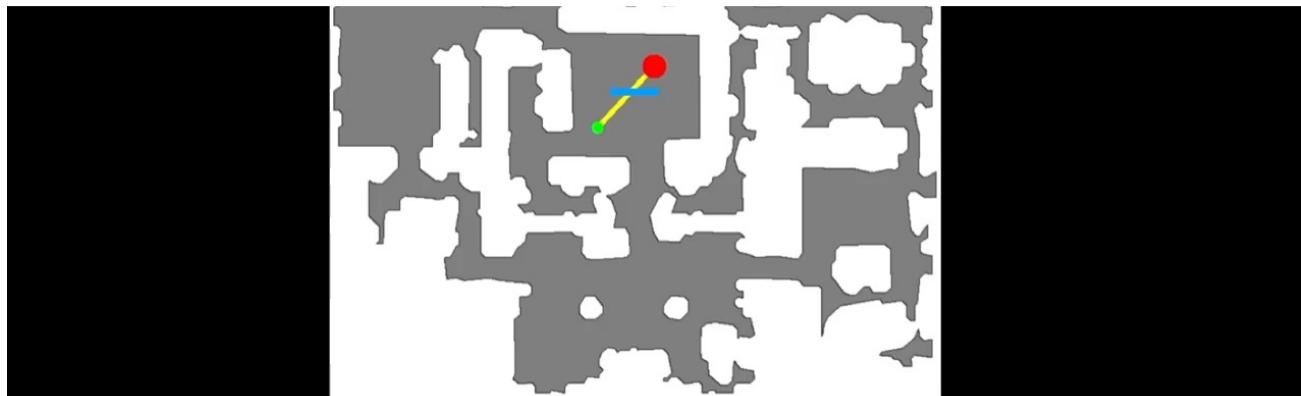
Assem  
Sadek

Guillaume  
Bono

Boris  
Chidlovskii

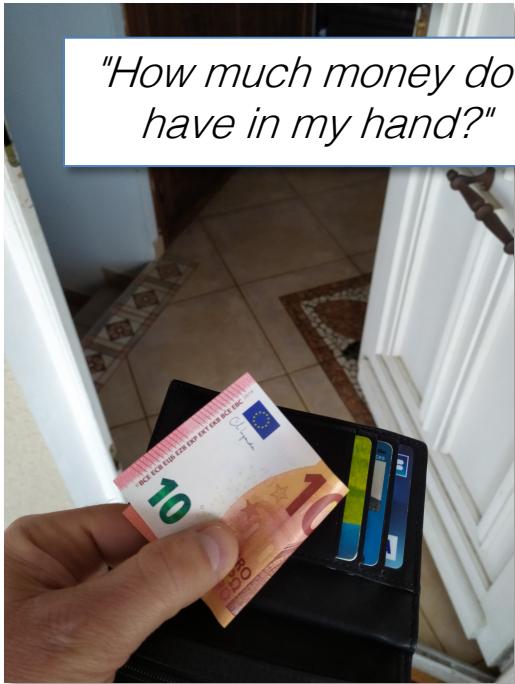
Christian  
Wolf

# Navigation in real environments



# Vision and Language Reasoning

*"How much money do I have in my hand?"*



*"What is in this jar?"*



*"Did I leave the door open?"*



*"Did I leave the lights on?"*



Corentin  
Kervadec



Grigory  
Antipov



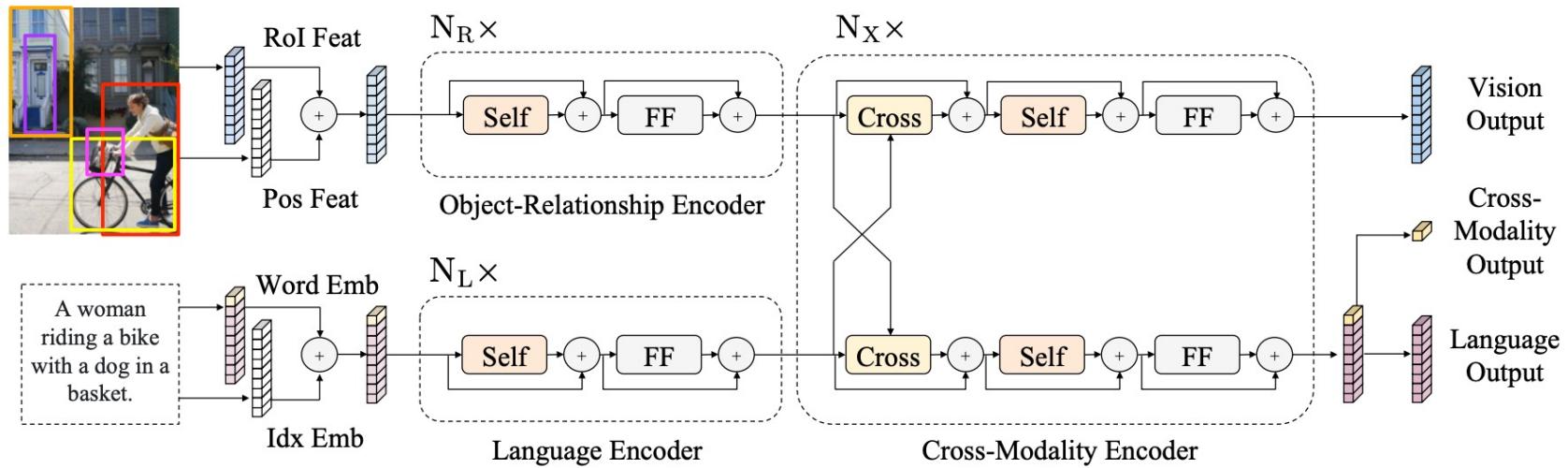
Moez  
Baccouche



Christian  
Wolf

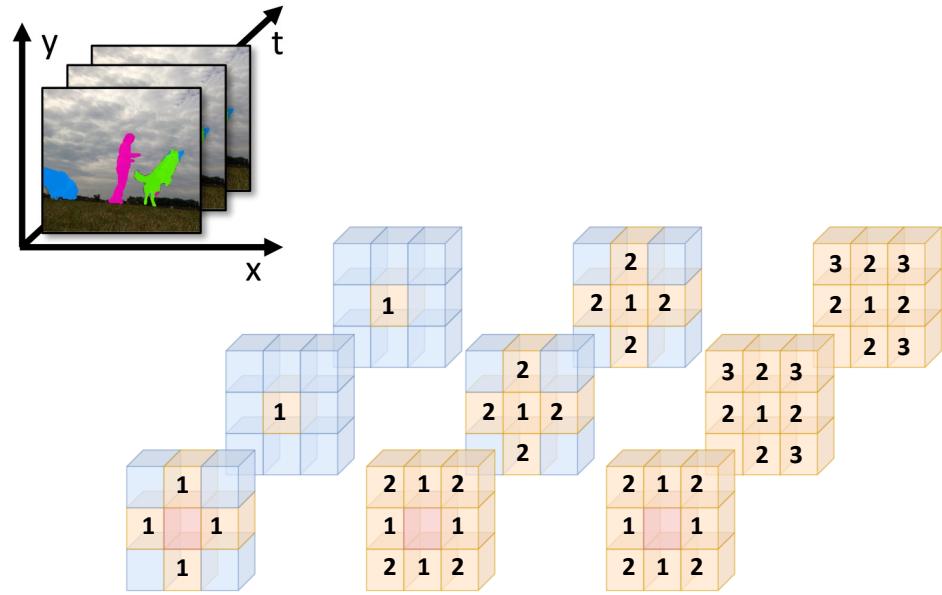
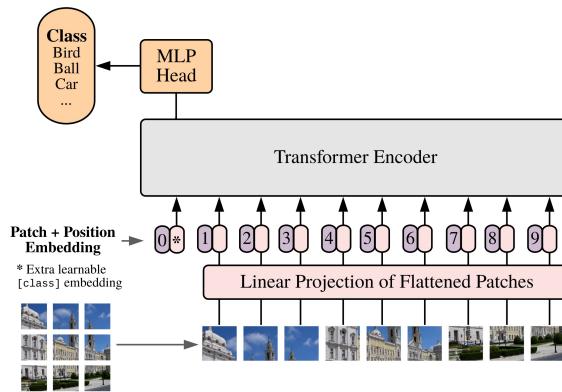
# LXMERT

A vision and language encoder with self-attention and cross-attention.



Tan, H. and Bansal, M. LXMERT: Learning cross-modality encoder representations from transformers. EMNLP-IJCNLP 2019.

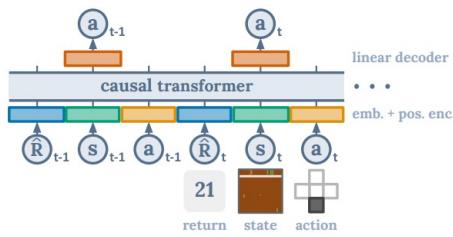
# Transformers for images and videos



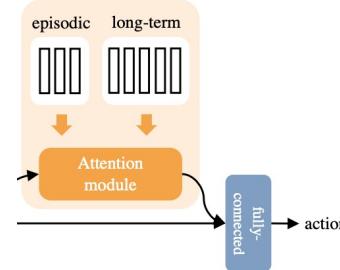
A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, N. Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, ICLR 2021.

[Duke, Ahmed, Wolf, Arabi, Taylor, CVPR 2021]

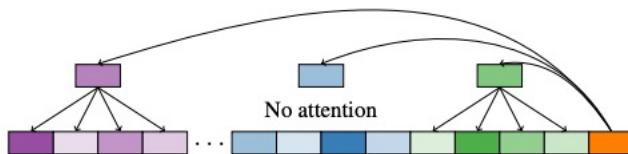
# Transformers for robotics



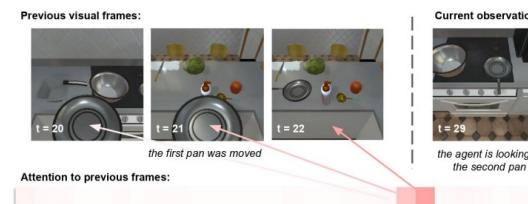
L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, I. Mordatch.  
Decision Transformer:  
Reinforcement Learning via  
Sequence Modeling arxiv 6/2021



L. Mezghani, S. Sukhbaatar, T. Lavril, O. Maksymets, D. Batra, Piotr B., K. Alahari, Memory-Augmented Reinforcement Learning for Image-Goal Navigation, arxiv 1/2021

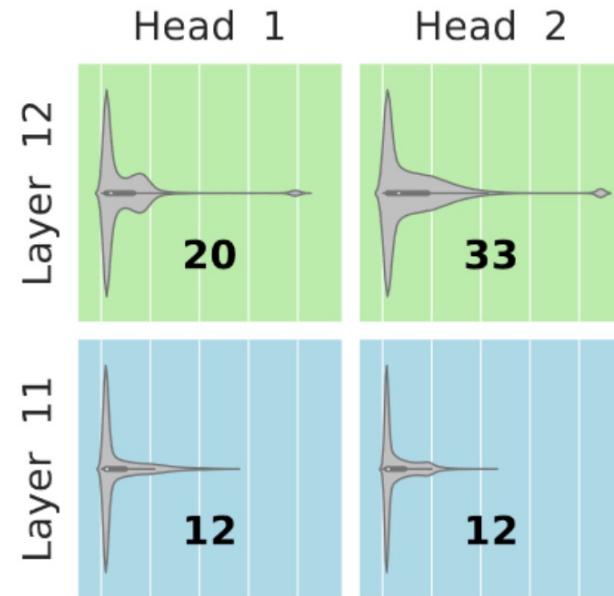
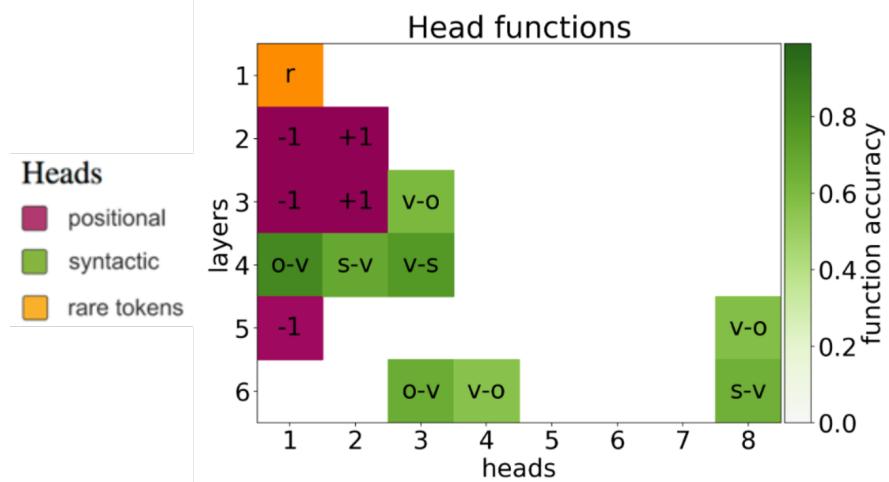


A.K. Lampinen, S.C.Y. Chan, A. Banino, F. Hill, Towards mental time travel: a hierarchical memory for reinforcement learning agents., arxiv 5/2021



A. Pashevich, C. Schmid, C. Sun,  
Episodic Transformer for Vision-and-Language Navigation, ICCV 2021.

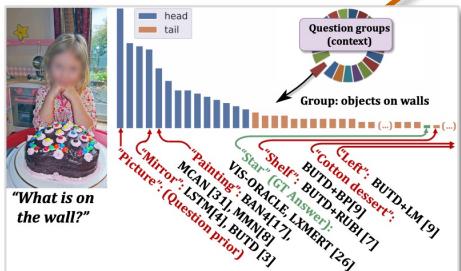
# Visualization of attention maps



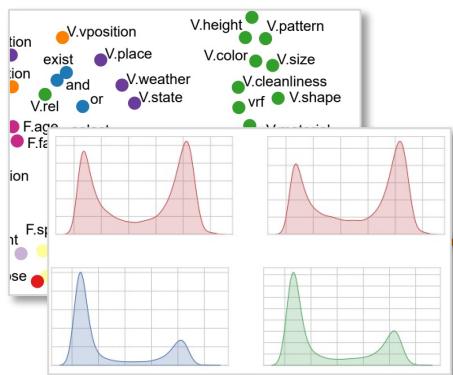
Voita, E., Talbot, D., Moiseev, F., Sennrich, R., and Titov, I., *Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned.* ACL 2019.

Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Pavlović, M., Sandve, G. K., Greiff, V., D. Kreil, M. Kopp, G. Klambauer, J. Brandstetter, S. Hochreiter. *Hopfield networks is all you need.* Pre-print arXiv:2008.02217, 2020.

How can we evaluate biases  
in learning? (CVPR 2021a)

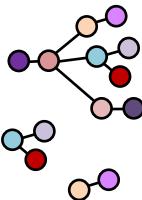


How can we visualize and  
transfer reasoning? (CVPR  
2021b, VIS 2021)



# VQA

Can we supervise reasoning  
programs? (NeurIPS 2021)



Can we weakly supervise word-  
object alignment? (ECAI 2020)

Can we ground object  
detection through language  
(under preparation)



Corentin  
Kervadec



Grigory  
Antipov

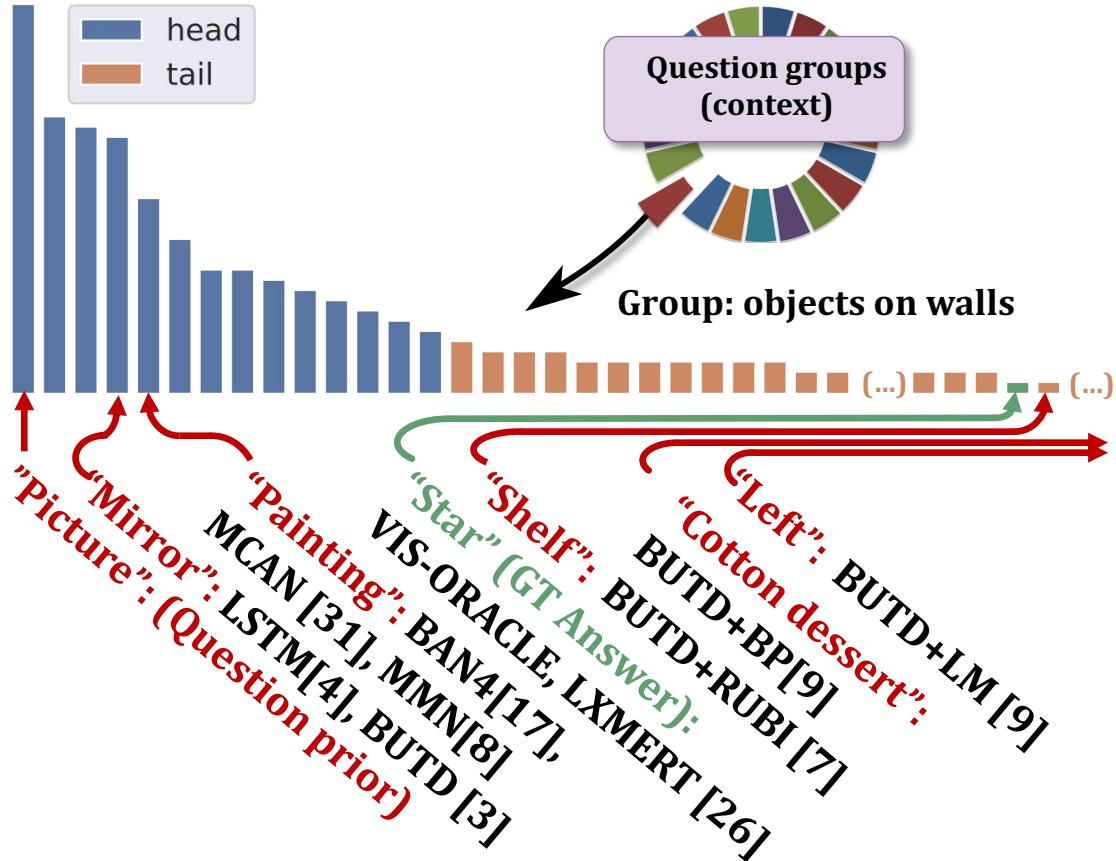


Moez  
Baccouche

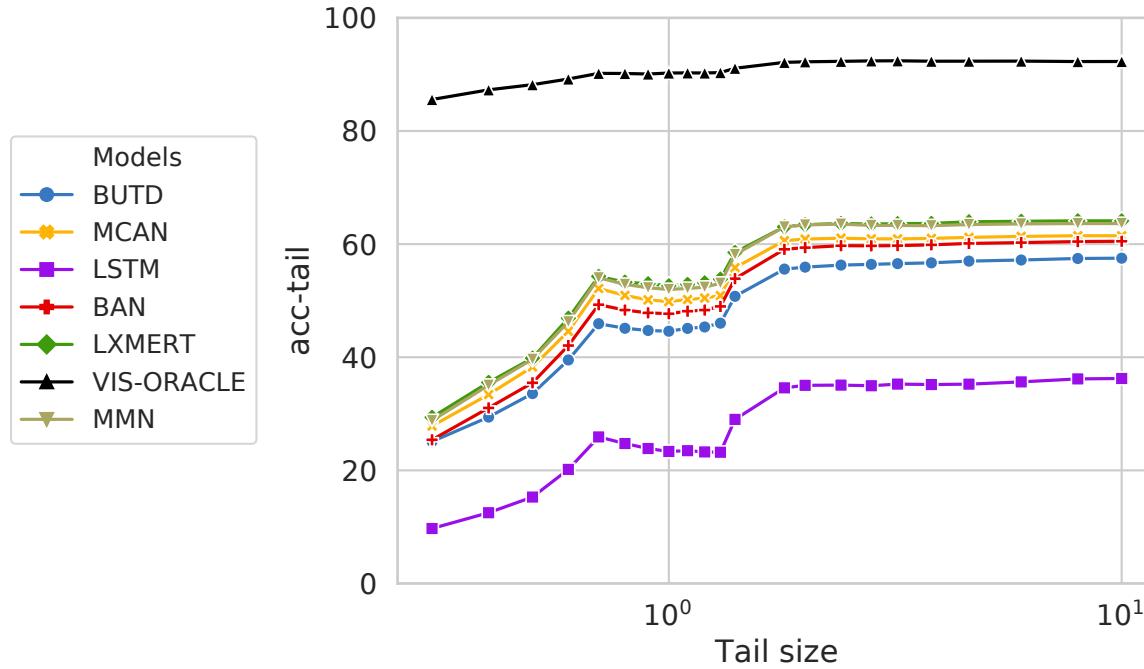


Christian  
Wolf

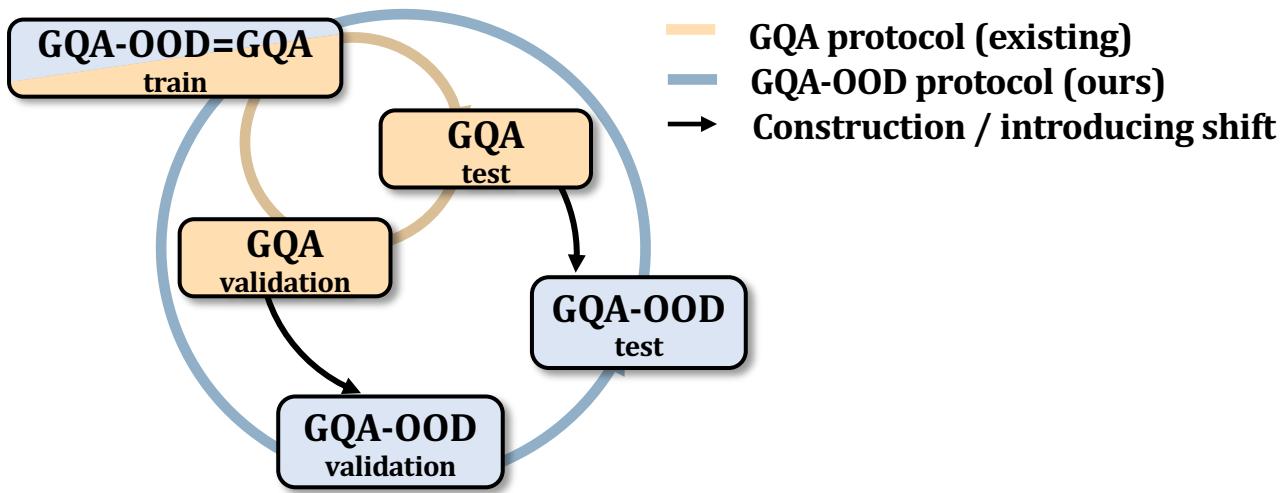
# Roses Are Red, Violets Are Blue... but Should VQA Expect Them To?



# Reasoning vs. bias exploitation



# OOD split



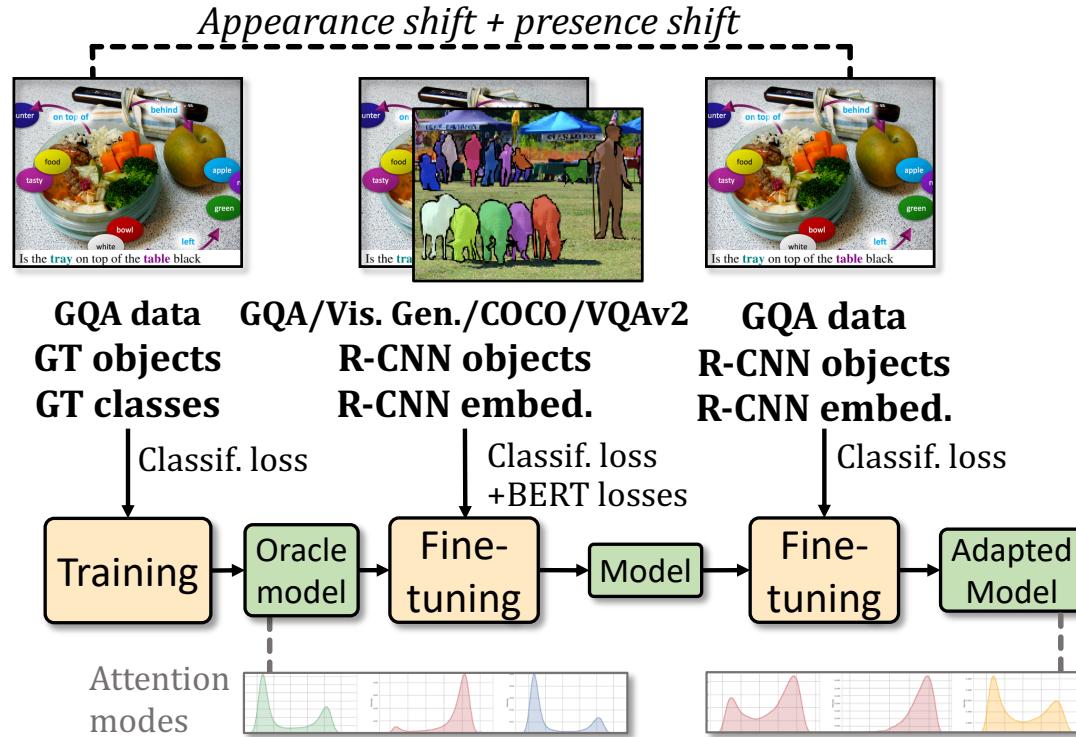
# Calibrating the metric

Model	Baseline benchm. Tot. Acc.	Proposed benchmark (Acc-tail)		
		$\alpha=1.2$	$\alpha=0.5$	$\alpha=0.3$
BUTD [3] + <i>bal</i>	$60.7 \pm 0.4$	$45.4 \pm 0.3$	$33.8 \pm 0.5$	$24.6 \pm 0.5$
BUTD [3] + <i>all</i>	$59.8 \pm 0.1$	$41.9 \pm 0.1$	$29.5 \pm 0.3$	$18.3 \pm 0.6$
$\Delta$ (relative):	-1.4%	-7.7%	-12.9%	-25.7%

# Results

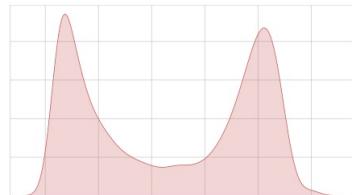
Model	VQA2 overall	GQA overall	GQA dist.	VQA-CP2 overall	GQA-OOD acc-tail
Q. Prior	32.1	27.0	55.6	8.8	17.8
LSTM [4]	43.0	39.1	3.6	22.1	24.0
BUTD [3]	63.5	51.6 $\pm$ 0.3	1.8	40.1	42.1 $\pm$ 0.9
MCAN [29]	<b>66.1</b>	56.3 $\pm$ 0.2	1.6	<b>42.5</b>	46.5 $\pm$ 0.5
BAN4 [18]	65.9	54.7 $\pm$ 0.4	1.6	40.7	47.2 $\pm$ 0.5
MMN [8]	-	<b>59.6</b>	1.8	-	48.0
LXMERT [24]	<b>69.9</b>	<b>59.6</b>	<b>1.5</b>	-	<b>49.8</b>
BUTD [3]	<b>63.5</b>	51.6 $\pm$ 0.3	1.8	40.1	42.1 $\pm$ 0.9
+RUBi+QB	-	<b>51.9</b> $\pm$ 1.1	<b>1.7</b>	<b>47.6</b> $\pm$ 3.7	42.1 $\pm$ 1.0
+RUBi [7]	61.2	43.6 $\pm$ 2.0	1.9	44.2	35.7 $\pm$ 2.3
+LM [9]	56.4	39.7 $\pm$ 0.7	2.1	52.0	32.2 $\pm$ 1.2
+BP [9]	63.2	39.6 $\pm$ 0.3	2.2	39.9	30.8 $\pm$ 1.0

# Oracle transfer

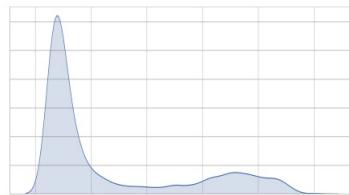


[Kervadec, Jaunet, Antipov, Baccouche, Vuillemot, Wolf,  
CVPR 2021b]

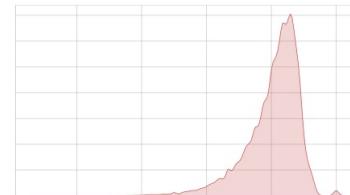
# K-numbers and attention modes



(a) Bimorph



(b) Dirac

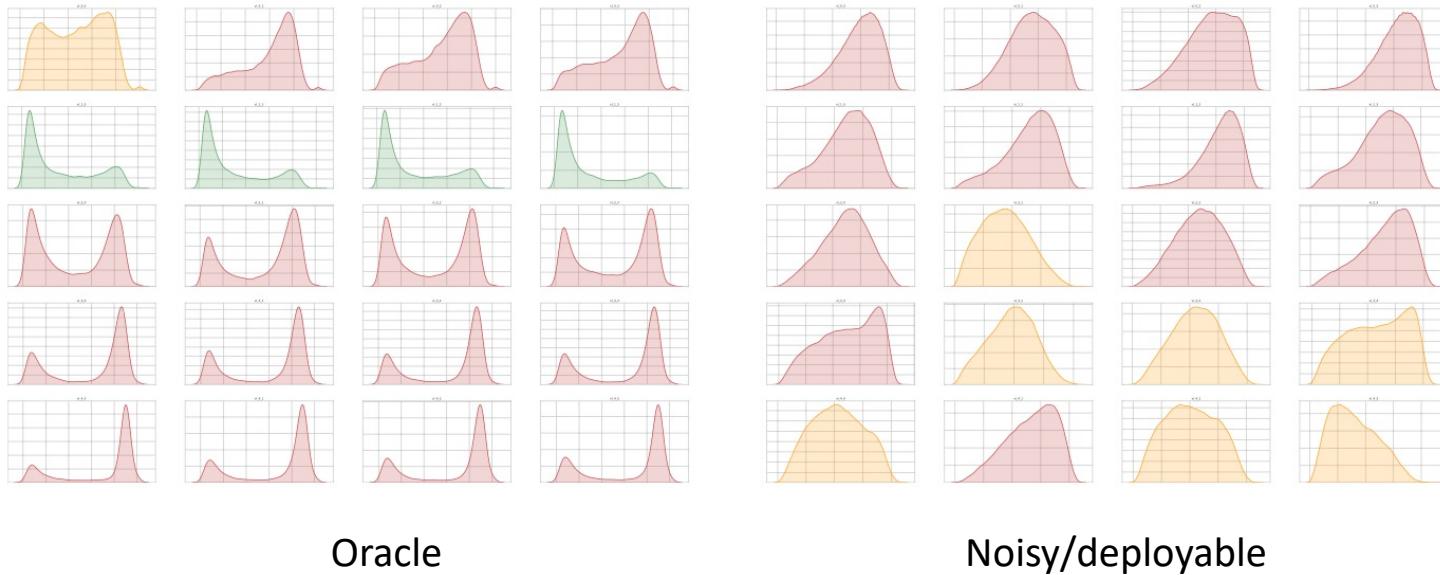


(c) Uniform

Distribution of the number  $k$  of tokens required to reach 90% of the attention energy, on GQA-val.

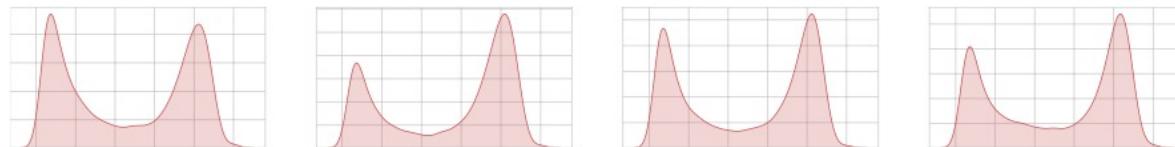
X-axis: ratio of the tokens  $k$  w.r.t. the total number of tokens.

# Analyzing reasoning patterns

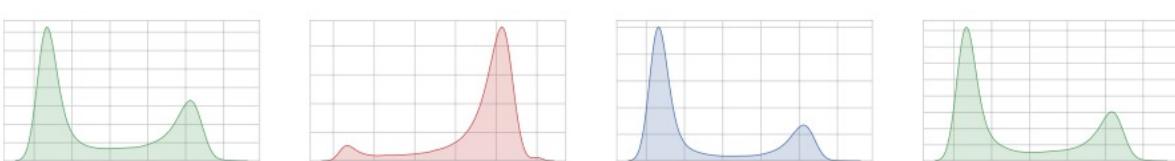


# Attention and language tasks (1)

(a)  
overall



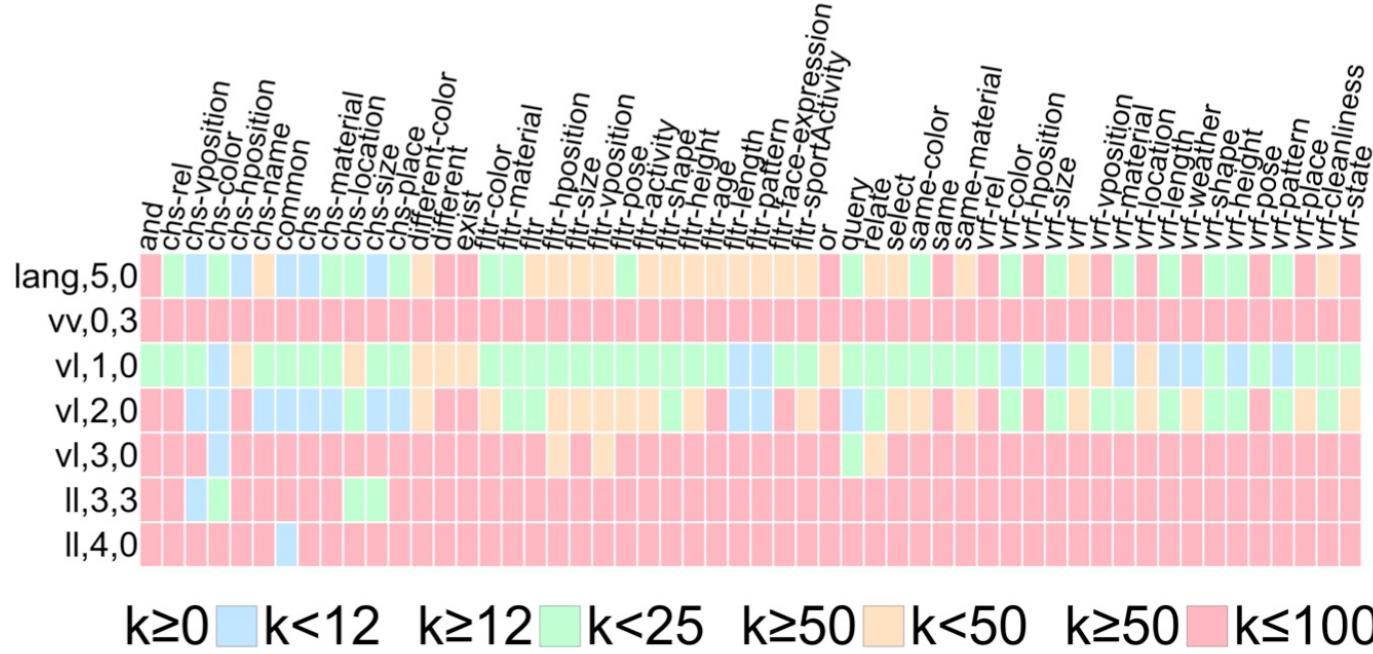
(b)  
choose  
color



Samples of all functions vs. samples with questions involving “choose color”.  
Oracle model, heads of third layer of  $T_X^{L \leftarrow V}$

Task dependence of attention head behavior:  
Activation of the 1st, 2nd and 4th head, desactivation of the 3rd.

# Attention and language tasks (2)

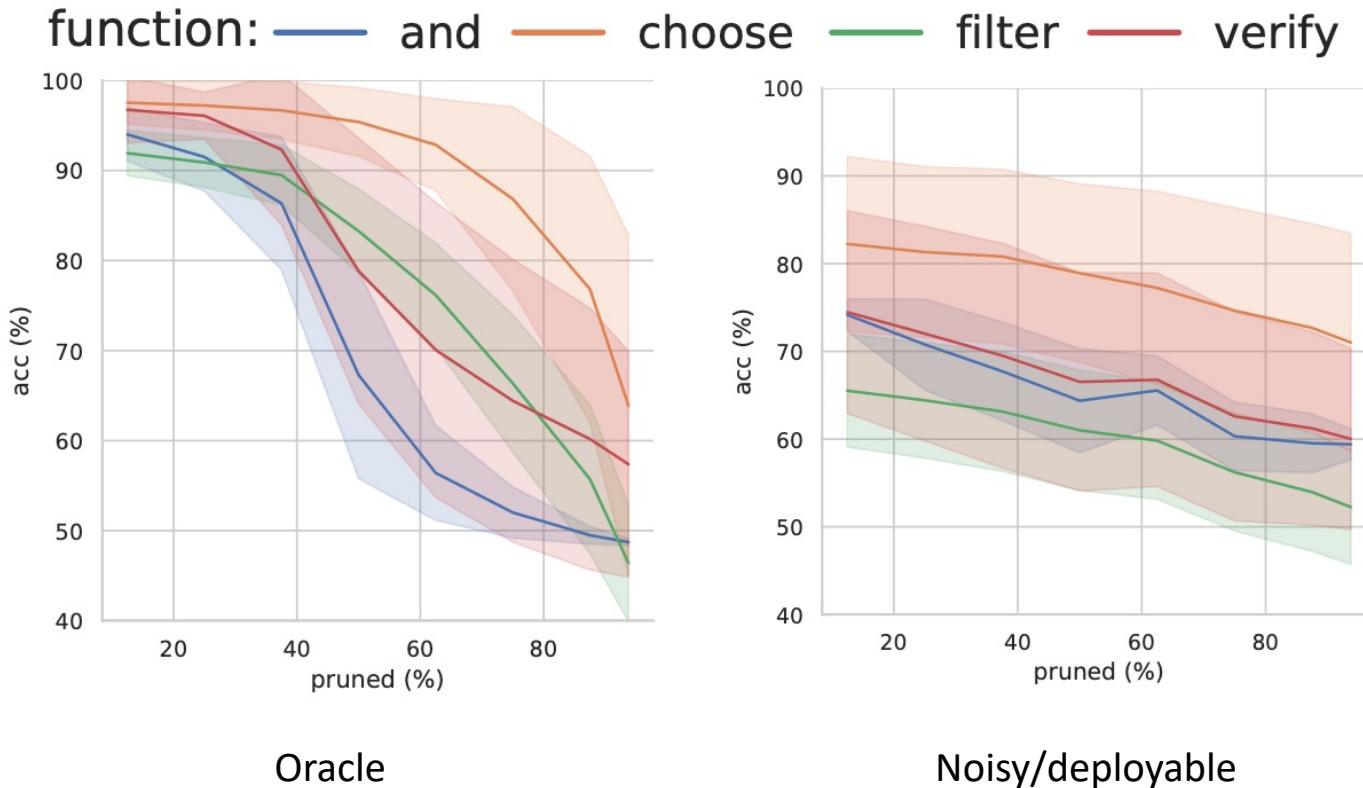


Attention heads behave differently depending on the function.  
A given function causes different attention modes for different heads.

(Oracle model, median k-numbers)

# Head pruning

Pruned attentions	n/a	L	V	$L \leftarrow V$	$V \leftarrow L$
Accuracy	91.5	37.9	91.4	52.8	68.1



# How transferable are Reasoning Patterns in VQA?

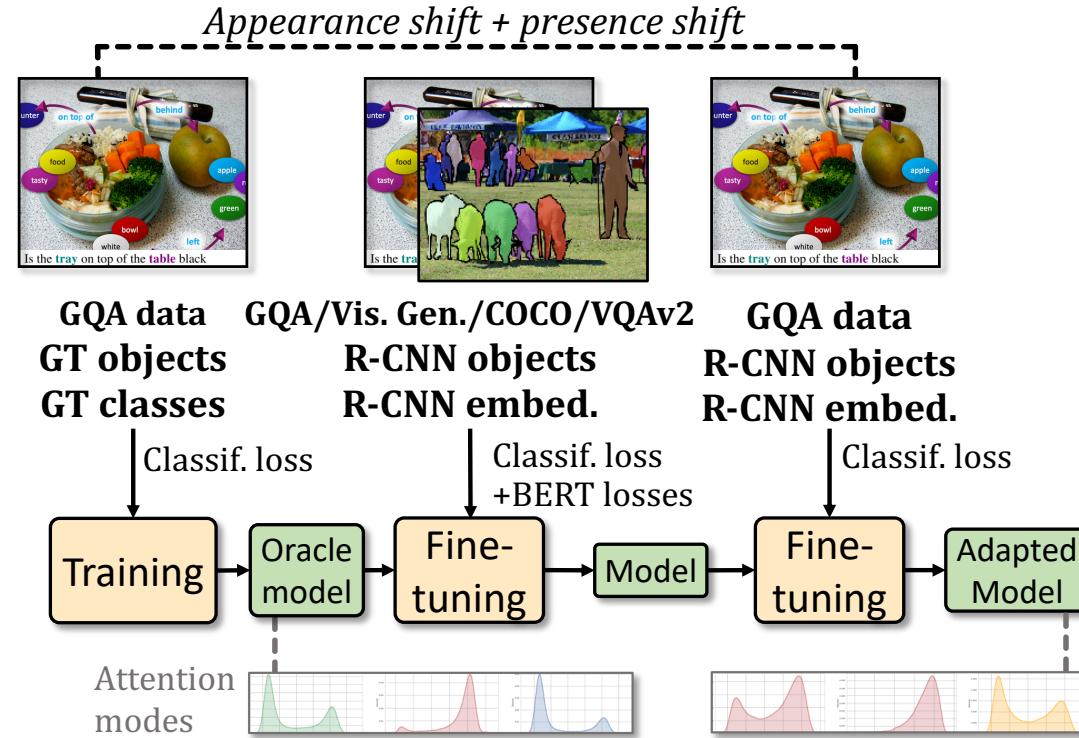
CVPR 2021

Corentin Kervadec Théo Jaunet Grigory Antipov  
Moez Baccouche Romain Vuillemot Christian Wolf

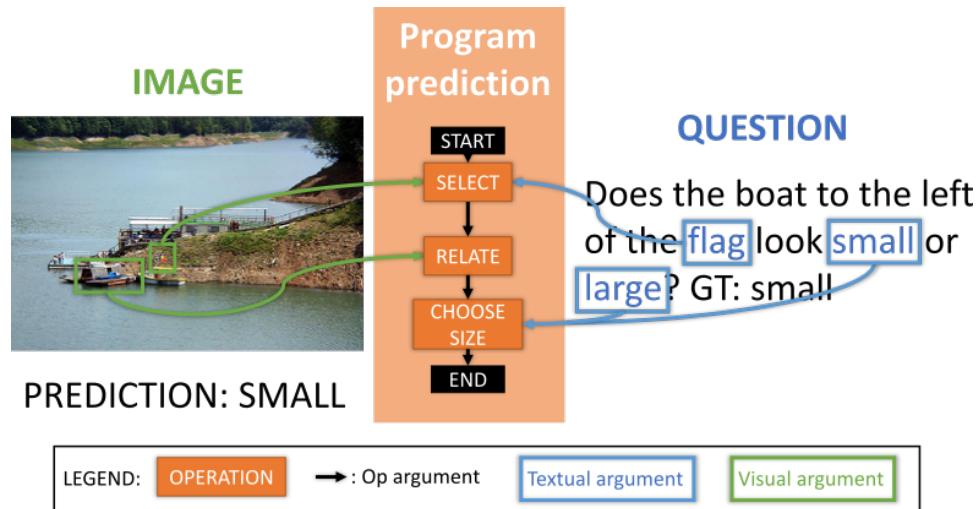
LIRIS, INSA-Lyon, Ecole Centrale de Lyon, Orange

<https://visqa.liris.cnrs.fr>

# Are reasoning patterns transferrable?

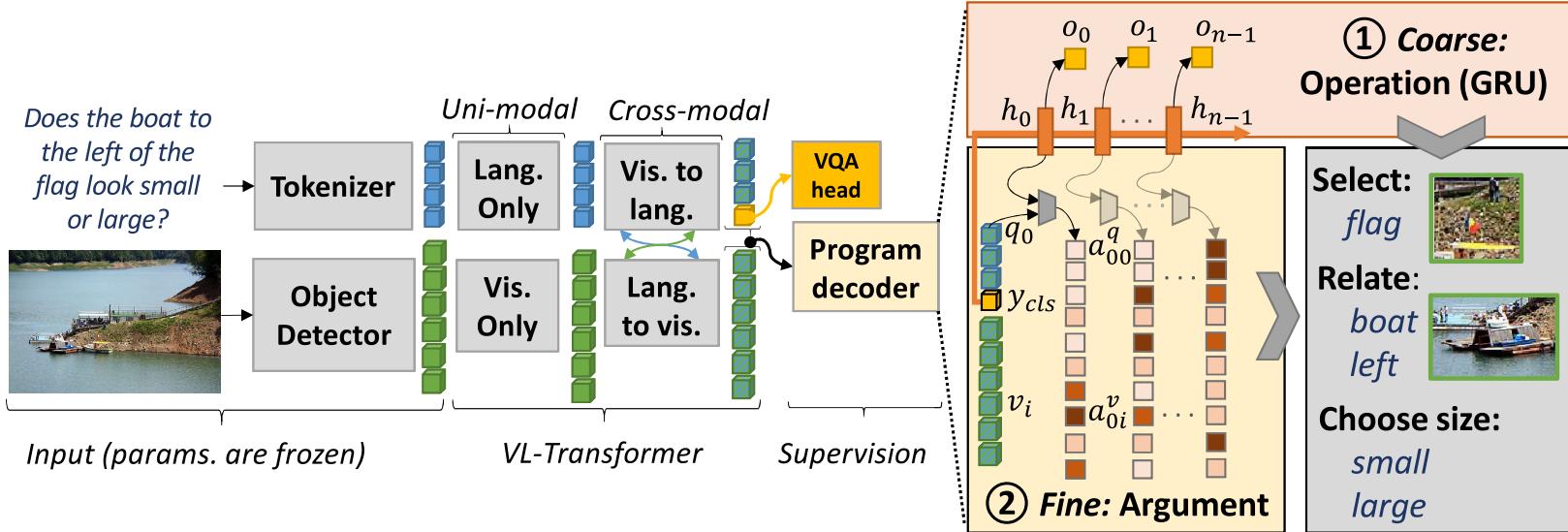


# GT reasoning programs



[Kervadec\*, Wolf\*, Antipov, Baccouche,  
Nadri, NeurIPS 2021]

# Program prediction



[Kervadec\*, Wolf\*, Antipov, Baccouche,  
Nadri, NeurIPS 2021]

# Oracle + program supervision

Model	Oracle transf.	Prog. sup.	GQA-OOD [20]		test-dev	GQA [17]			AUC <sup>†</sup> prog.
			acc-tail	acc-head		binary*	open*	test-std	
scratch	(a) Baseline		42.9	49.5	52.4	-	-	-	/
	(b) Oracle transfer	✓	$48.2 \pm 0.3$	$54.6 \pm 1.1$	$57.0 \pm 0.3$	74.5	42.1	57.3	/
	(c) Ours	✓	$48.8 \pm 0.1$	$56.1 \pm 0.3$	$57.8 \pm 0.2$	<b>75.4</b>	<b>43.0</b>	<b>58.2</b>	97.1
+ Lxmert	(d) Baseline		47.5	55.2	58.5	-	-	-	/
	(e) Oracle transfer	✓	47.1	54.8	58.4	77.1	42.6	58.8	/
	(f) Ours	✓	$48.0 \pm 0.6$	$56.6 \pm 0.6$	$59.3 \pm 0.3$	<b>77.3</b>	<b>44.1</b>	<b>59.7</b>	96.4

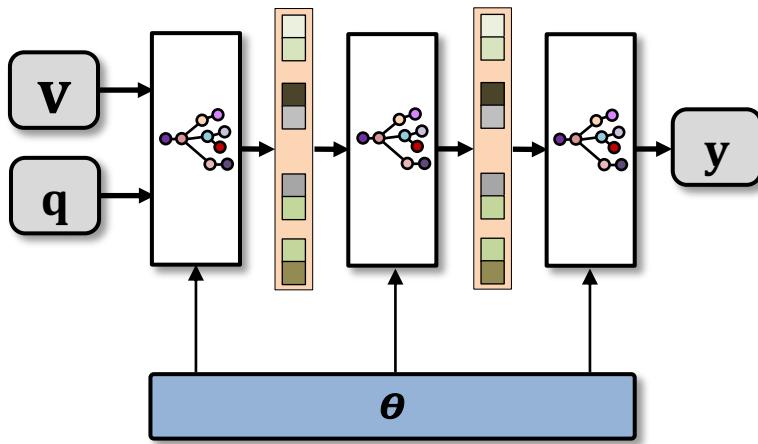
Table 1: Impact of program supervision on *Oracle transfer* [23] for vision-language transformers. LXMERT [36] pre-training is done on the GQA unbalanced training set. We report scores on GQA [17] (*test-dev* and *test-std*) and GQA-OOD (*test*). \* binary and open scores are computed on the test-std; <sup>†</sup> we evaluate visual argument prediction by computing AUC@0.66 on GQA-val.

# Comparison with SOTA

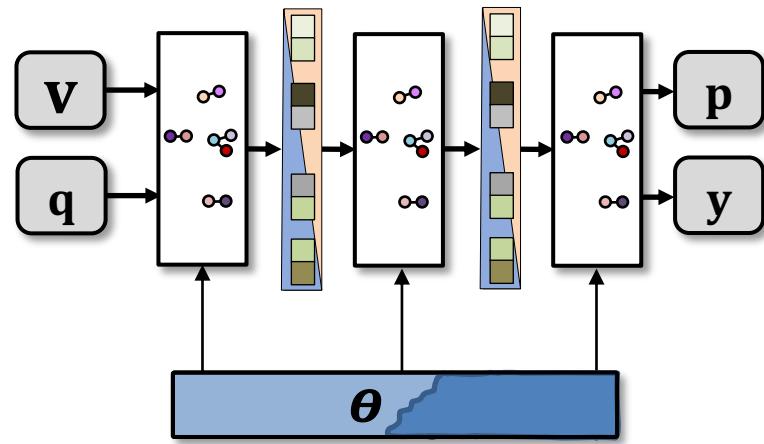
Method	Visual feats.	Additional supervision	Training data (M)		GQA-OOD [18]		GQA [15]		
			Img	Sent	acc-tail	acc-head	bin.	open	all
BAN4 [22]	RCNN [2]	-	≈ 0.1	≈ 1	47.2	51.9	76.0	40.4	57.1
MCAN [37]	RCNN [2]	-	≈ 0.1	≈ 1	46.5	53.4	75.9	42.2	58.0
Oracle transfer [21]	RCNN [2]	-	≈ 0.18	≈ 1	48.3	55.5	75.2	44.1	58.7
MMN [6]	RCNN [2]	Program	≈ 0.1	≈ 15	48.0	55.5	78.9	44.9	60.8
LXMERT [30]	RCNN [2]	-	≈ 0.18	≈ 9	<b>49.8</b>	57.7	77.8	45.0	60.3
<b>Ours</b>	VinVL [39]	Program	≈ 0.1	≈ 15	49.1	<b>59.7</b>	80.1	48.0	63.0
NSM [14]	SG [14]	Scene graph	≈ 0.1	≈ 1	-	-	78.9	<b>49.3</b>	63.2
OSCAR+VinVL [39]	VinVL [39]	-	≈ 5.7	≈ 9	-	-	<b>82.3</b>	48.8	<b>64.7</b>

Compact version: embedding dim=128, 4 heads  
 26M params instead of 212M (LXMERT: dim=768, 8 heads).

# Sample complexity



Classical training, CE loss on answers **y**

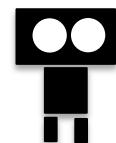


CE loss on **y** + program prediction **p**

- Learned knowledge of the reasoning processes
- Latent variables necessary for reasoning over multiple hops (used by reasoning processes)

- Decomposition of the underlying (unknown) reasoning function

# Outlook: language as Universal knowledge source



*I am looking for a pillow and I see a glimpse of a bedroom there, where pillows are frequently found. Let's check at this place first*



Pierre  
Marza



Laetitia  
Matignon



Olivier  
Simonin



Christian  
Wolf

# Conclusion

- Vision-Language reasoning as
  - A goal in its own right
  - A mean to transfer knowledge to downstream tasks
- Short cuts and Clever-Hans effects are ubiquitous
- Out-of-distribution testing is important
- Cleaner data leads to easier learning of proper reasoning
- Challenges
  - Find correct learning signals
  - Create inductive biases
  - Identify and collect data sources (world models?)