

C. Kervadec<sup>1,2</sup> T. Jaunet<sup>2</sup> G. Antipov<sup>1</sup> M. Baccouche<sup>1</sup> R. Vuillemot<sup>3</sup> C. Wolf<sup>2</sup>

<sup>1</sup> Orange Innovation <sup>2</sup> LIRIS, INSA Lyon

<sup>3</sup> LRIS, EC Lyon

Check-out our interactive demonstration (online)

<https://reasoningpatterns.github.io>

## Reasoning vs. Biases in VQA

VQA models are notorious for their tendency to rely on **shortcuts** [2,3], preventing them to **reason**.

We claim that *shortcut learning in VQA is in part due to the visual uncertainty (image representation imperfect)*.

## Our contributions

- In-depth analysis of **reasoning patterns** at work in VQA
- Analysing attention mechanisms at work in a **VL-Transformer**
- Comparing models with **perfect-sight** vs. **noisy visual inputs**

An **oracle transfer** method:

- Transfer reasoning capabilities, learned by the **oracle**, to a **standard VQA model with noisy input**
- Improve overall performance and generalisation on GQA[1]

## Oracle vs. Standard model

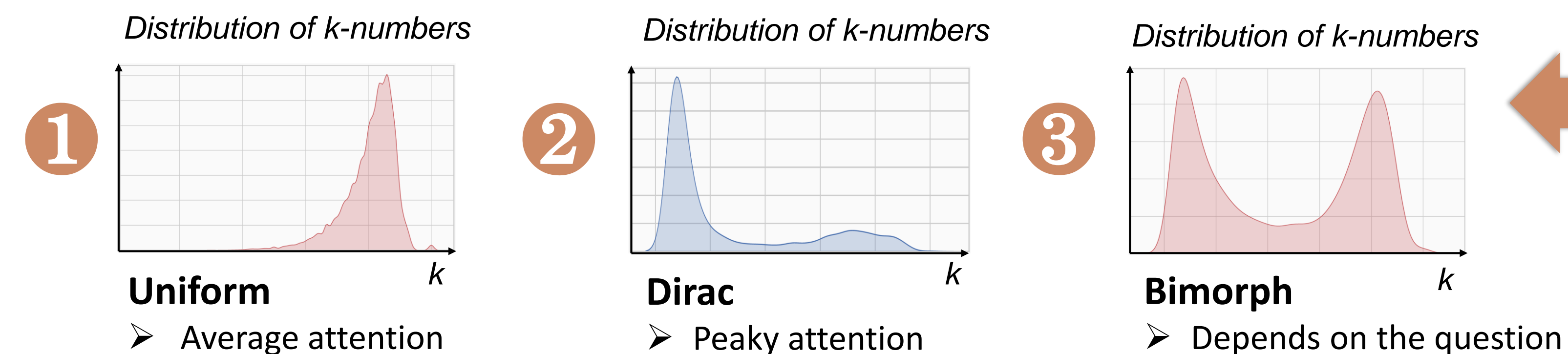
We propose to compare two settings:

- **Standard** ● *Vision is uncertain* ● *deployable*
  - Image is represented as a set of objects extracted using a pre-trained object detector
- **Oracle** ● *Vision is perfect* ● *not deployable*
  - Image is represented using human annotations

Our experiments are based on a **Vision-Language (VL)-Transformer**

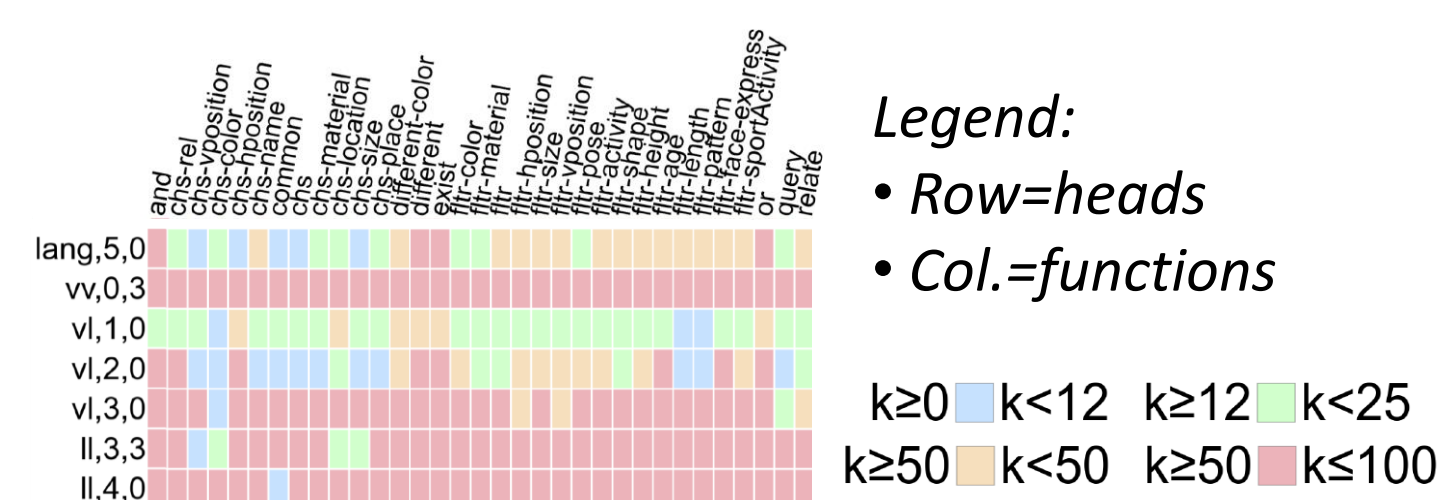
## Attention modes in VL-Transformers

We identify and plot three main attention modes in **Oracle's** attention heads:



## Oracle learned functions

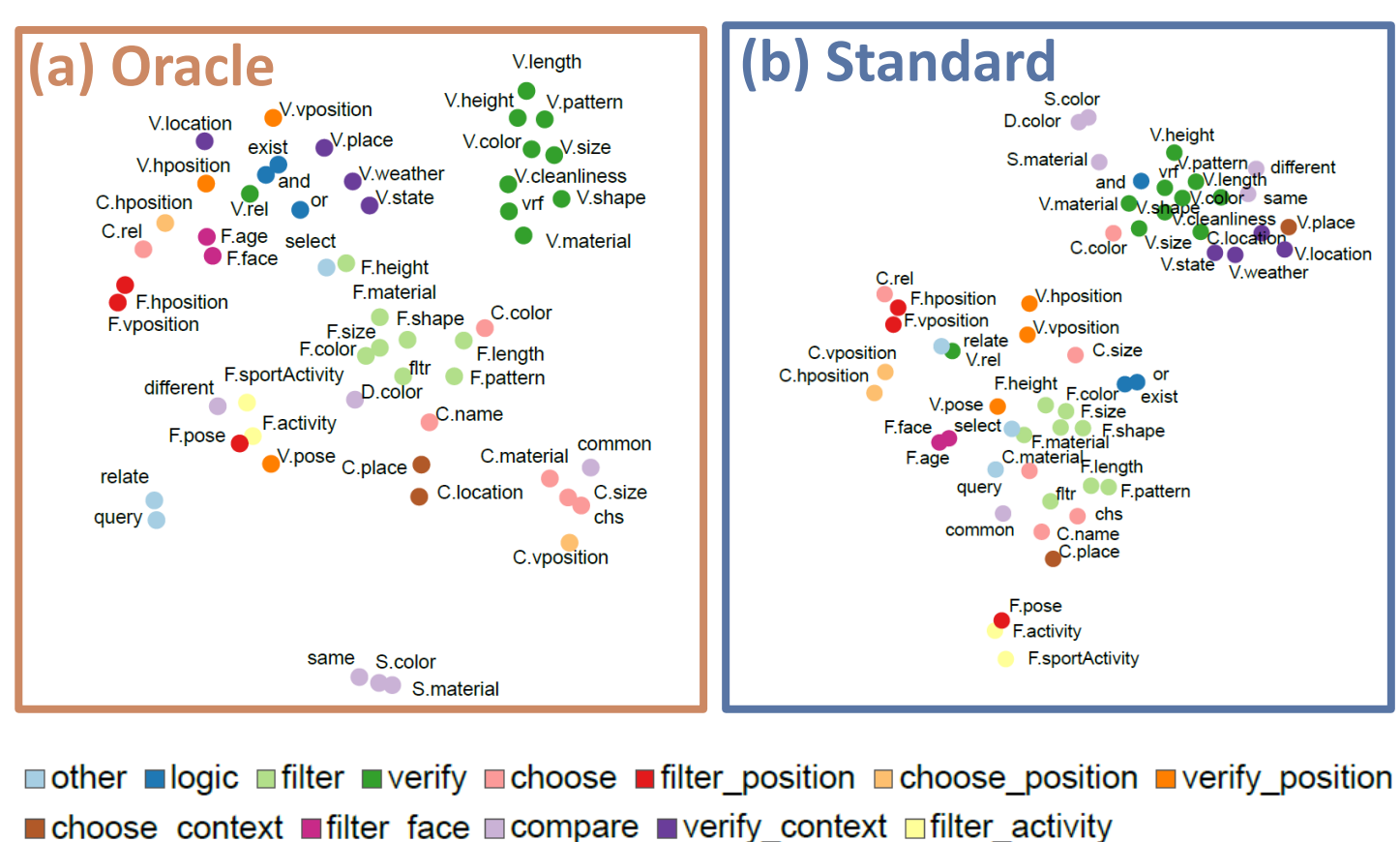
We link functions (e.g. choose color) to attention modes. Here for **Oracle**:



➤ **Oracle** adapts its attention to the task

But **Standard** did not

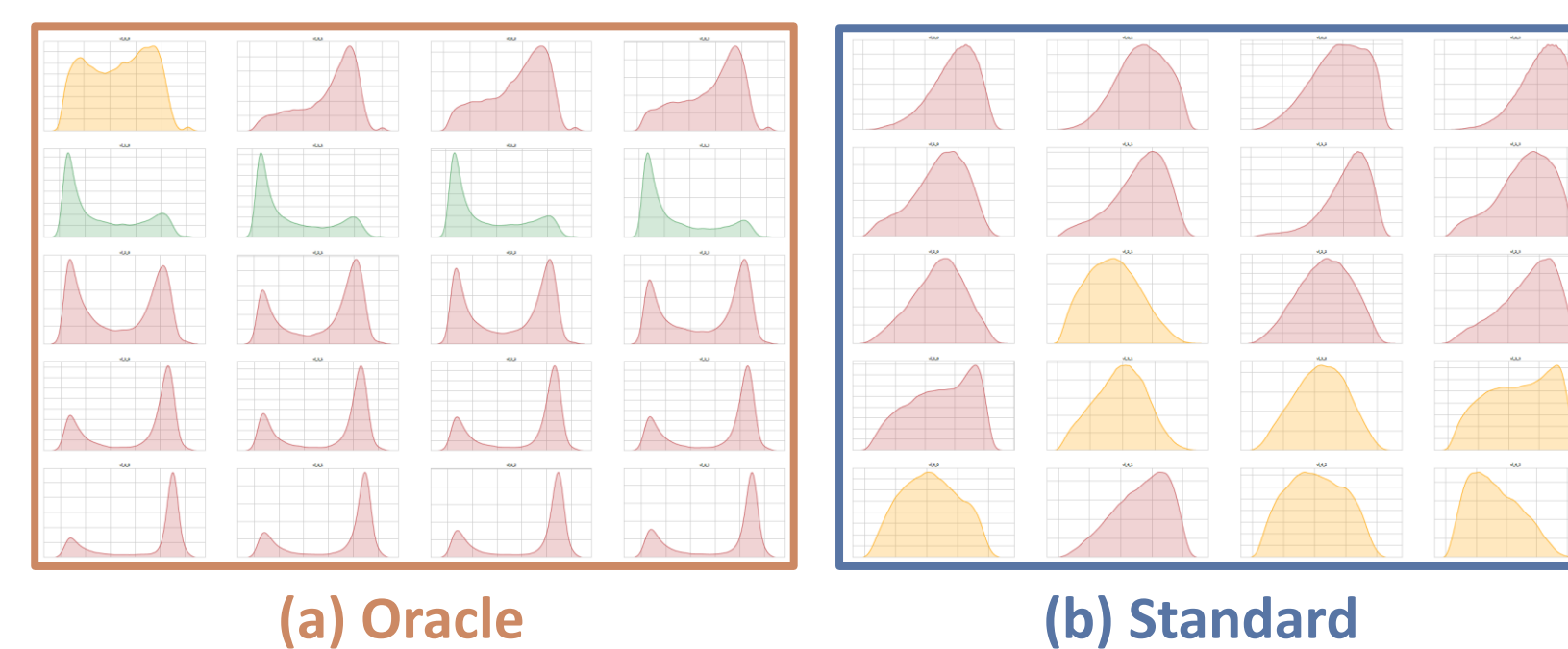
No relationships between attention modes and function in **Standard**\*:



\* t-SNE visualisation. See paper to get more details

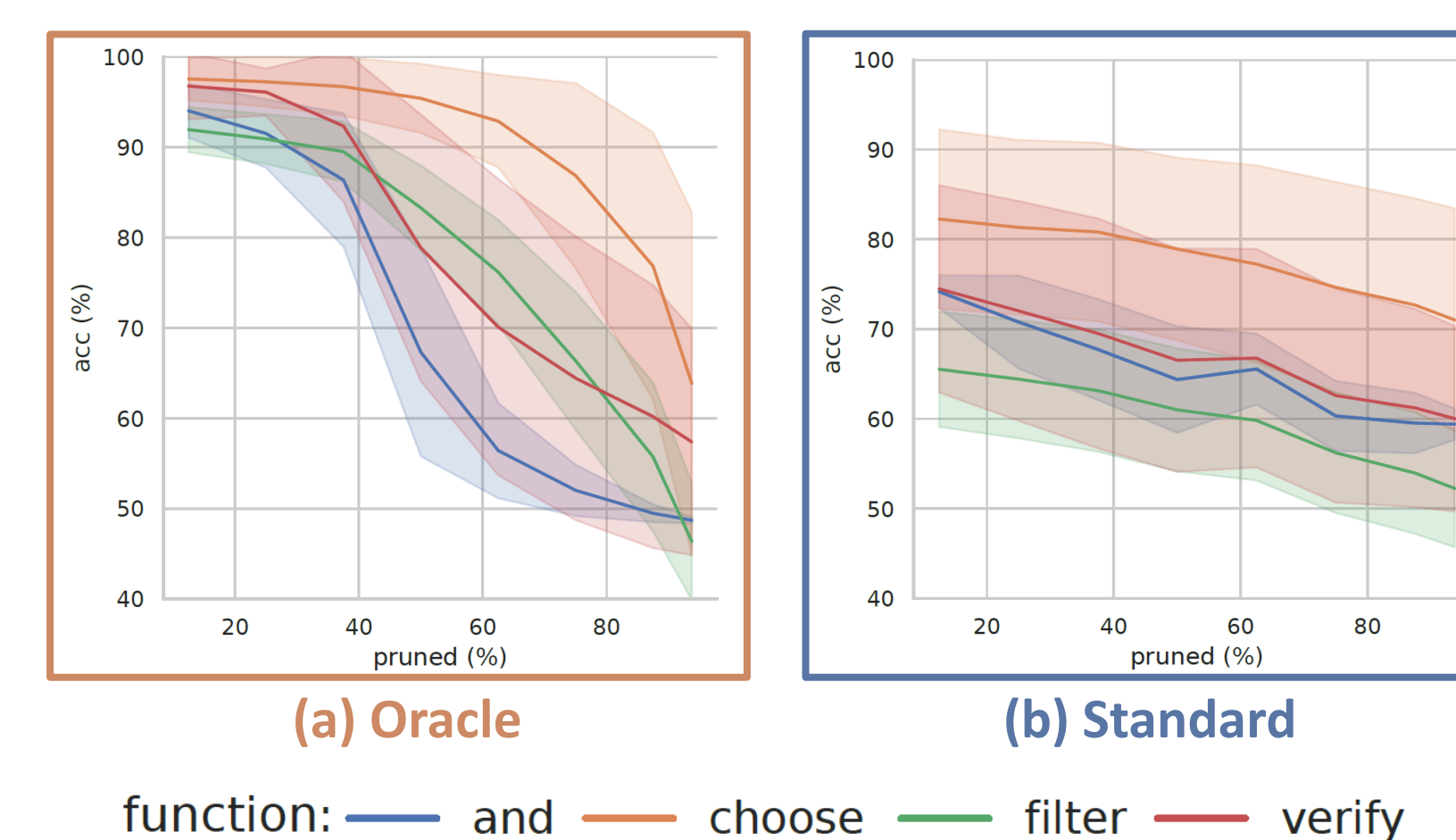
## Higher diversity in Oracle

**Oracle's** heads have more diverse attention modes: ➤ **Standard** ones are mostly uniform

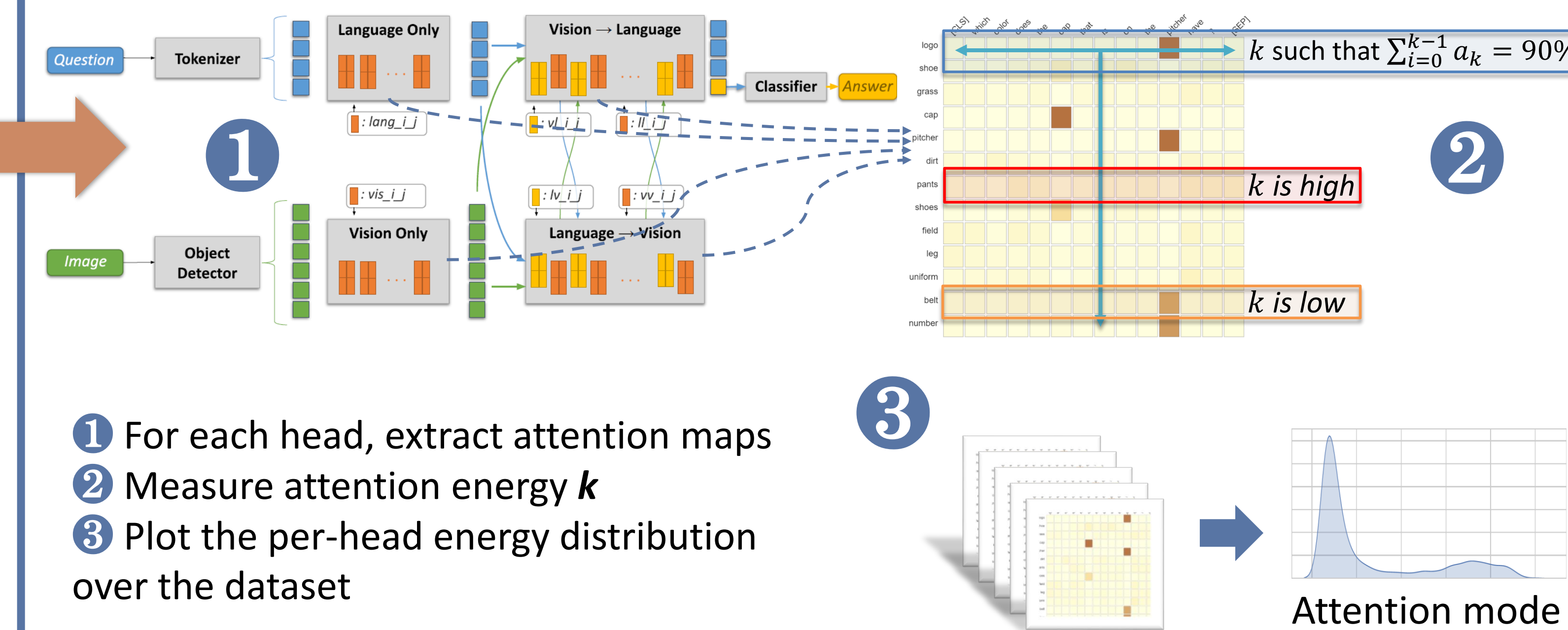


## Head pruning

Impact of pruning varying numbers of attention heads in cross-modal layers on accuracy. For **Oracle**, the impact is related to function. Result are different for **Standard**.



## Methodology: attention modes



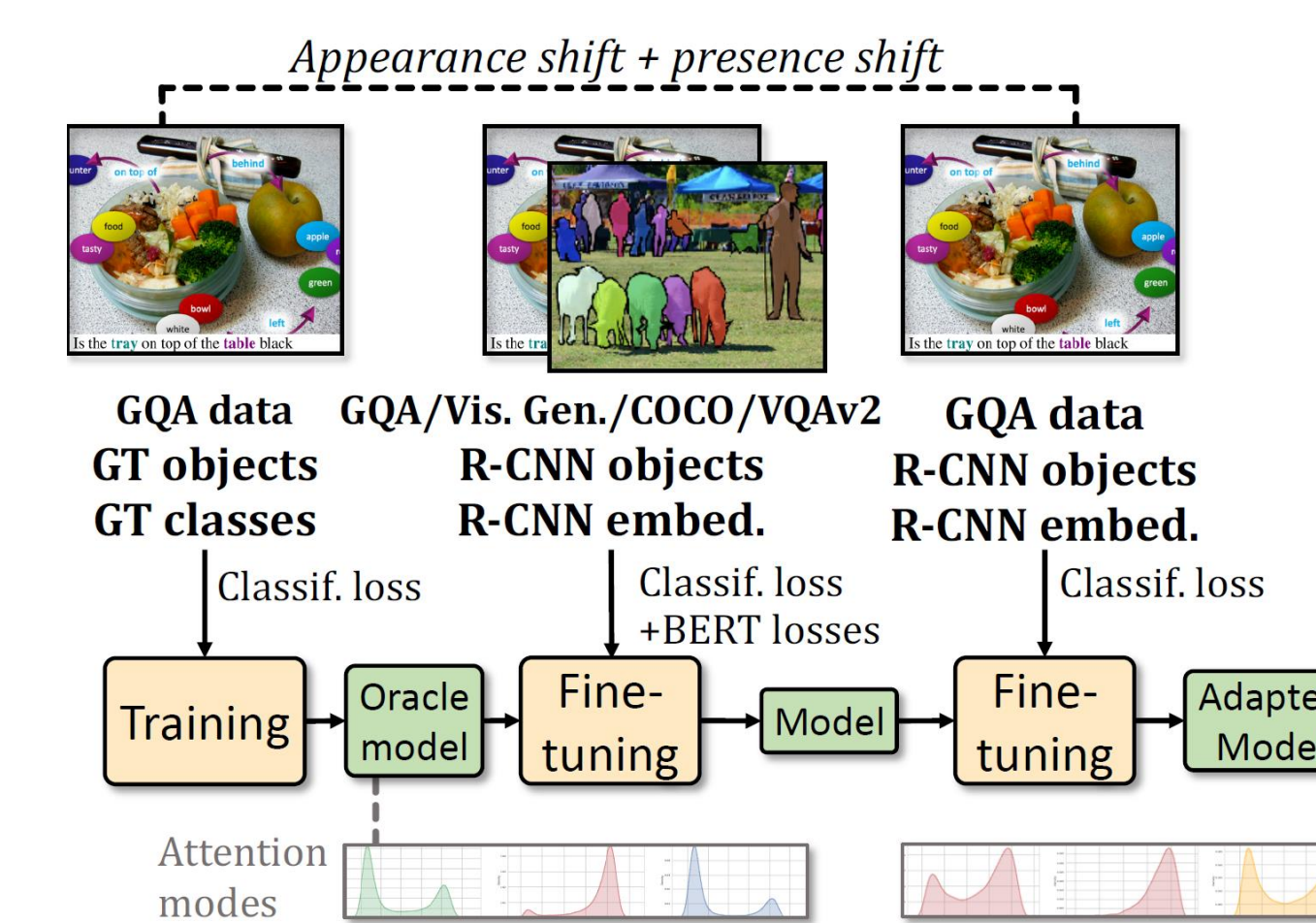
- 1 For each head, extract attention maps
- 2 Measure attention energy  $k$
- 3 Plot the per-head energy distribution over the dataset

## Oracle transfer

We observe significant differences between **Oracle** and **Standard**: we highlighted the **Oracle** ability of adapting reasoning to the task at hand.

We propose to transfer learned reasoning patterns from **Oracle** to **Standard**:

- 1 Train the oracle on perfect vision
- 2 Optionally, BERT-like pretraining
- 3 Finetune with standard (noisy) vision



Model	Pretraining		GQA-OOD [22]		GQA [19]	VQAv2 [17]
	Oracle	LXMERT/BERT	acc-tail	acc-head	overall	overall
(a) Baseline			42.9	49.5	52.4	-
(b) Ours	✓		<b>48.5</b>	<b>55.5</b>	<b>56.8</b>	-
(c) Baseline (+LXMERT/BERT)		✓	47.5	54.7	56.8	69.7
(d) Ours (+LXMERT/BERT)	✓	✓	<b>48.3</b>	<b>55.2</b>	<b>57.8</b>	<b>70.2</b>

[1] R. Geirhos, et al. Shortcut learning in deep neural networks. In Proc Nature Machine Intelligence 2020  
 [2] C. Kervadec, et al. Roses are Red, Violets are Blue... But Should VQA expect Them To? In Proc CVPR 2021  
 [3] D. Hudson, et al. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In Proc CVPR 2019