# Inference - Time Intervention: Eliciting Truthful Answers From a Language Model

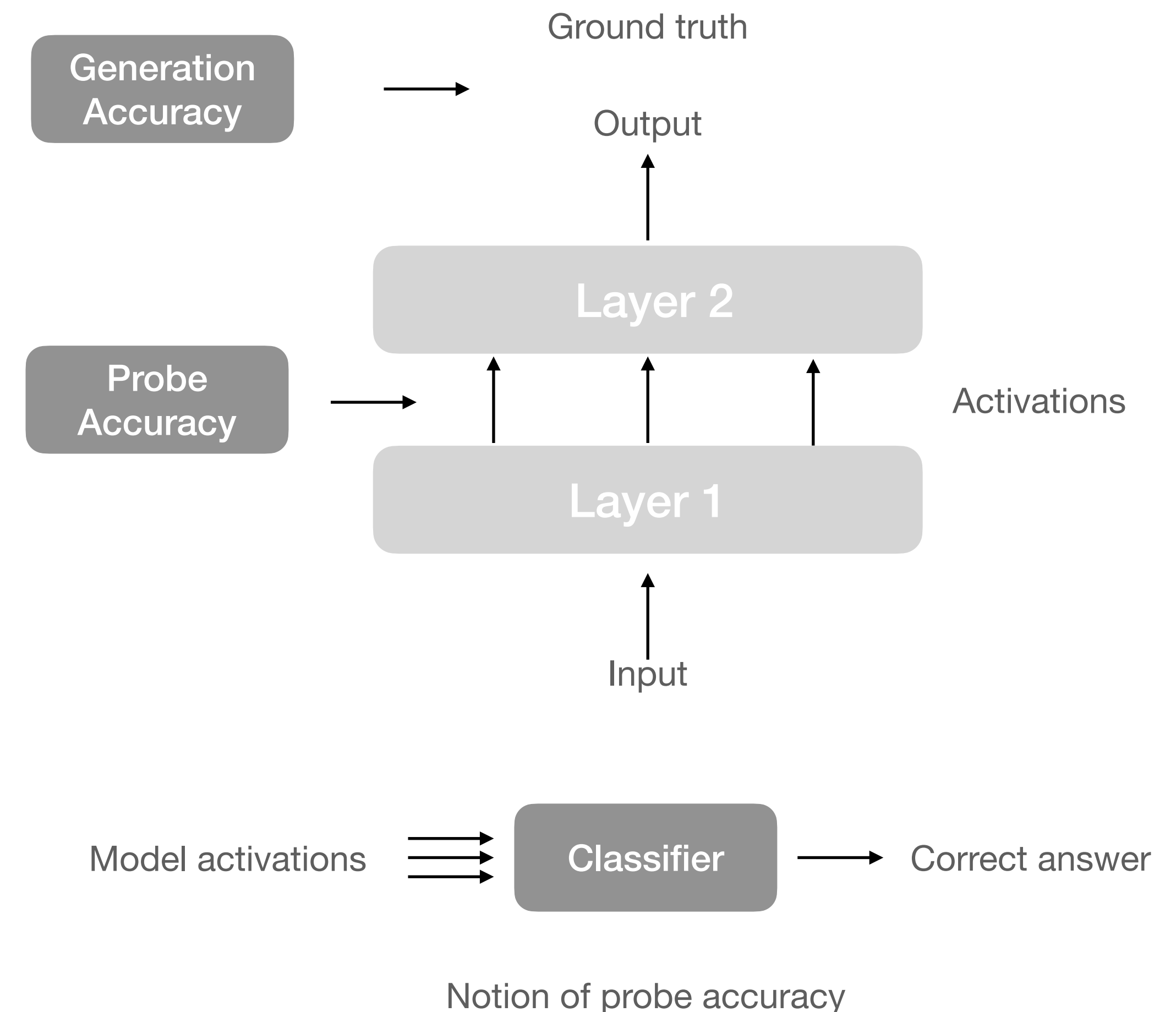**Kenneth Li, Oam Patel, Fernanda Viegas, Hanspeter Pfister, Martin Wattenberg**

Abhibha

# Motivation

- LLM's often generate wrong answers or 'hallucinations'.

- On prompting again, they sometimes provide the right answer.

- Focus on mistakes where the model 'knows' the correct answer but standard generation tactics fail to generate the right response.
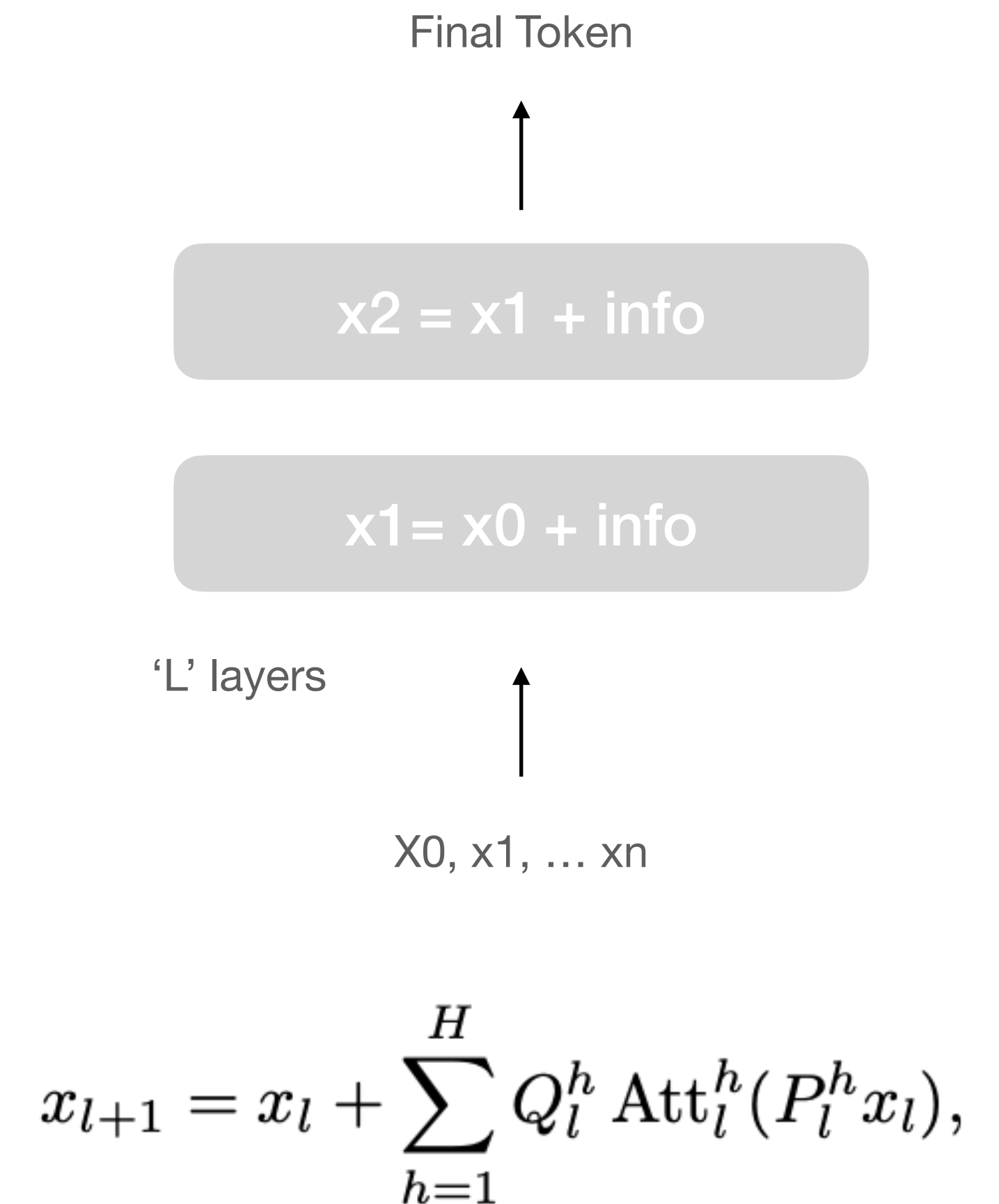
# Proposed Idea

- LM's contain latent, interpretable structure related to factuality - can be used to reduce incorrect anwers

- Inference Time Intervention

  - Identify attention heads with high probing accuracy for truthfulness

  - Shift activations during inference along truth-correlated directions.



Notion of probe accuracy

# Background
## Multi-head attention (MHA)
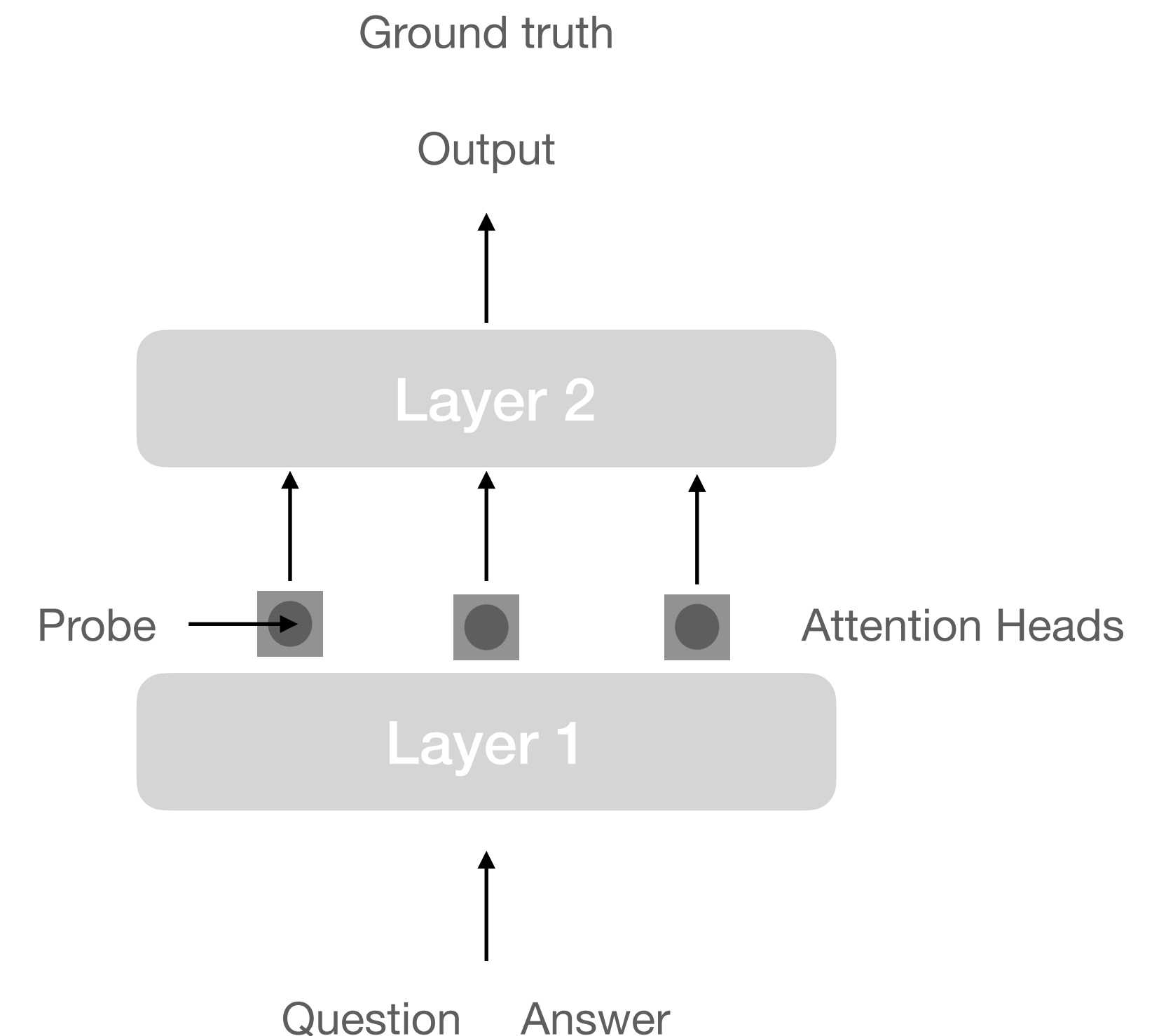
- Tokens embedded into high dimensional space: x0 = R (D x H)

- residual stream = x0 , x1, … xn

- H seperate linear operations

- P -> maps stream activation to D - dim space R(D x DH)

- Q -> maps it back  R(DH x D)

Final Token

x2 = x1 + info

x1= x0 + info

'L' layers

X0, x1, … xn

$$x_{l+1} = x_l + \sum_{h=1}^{H} Q_l^h \, \mathrm{Att}_l^h (P_l^h x_l),$$
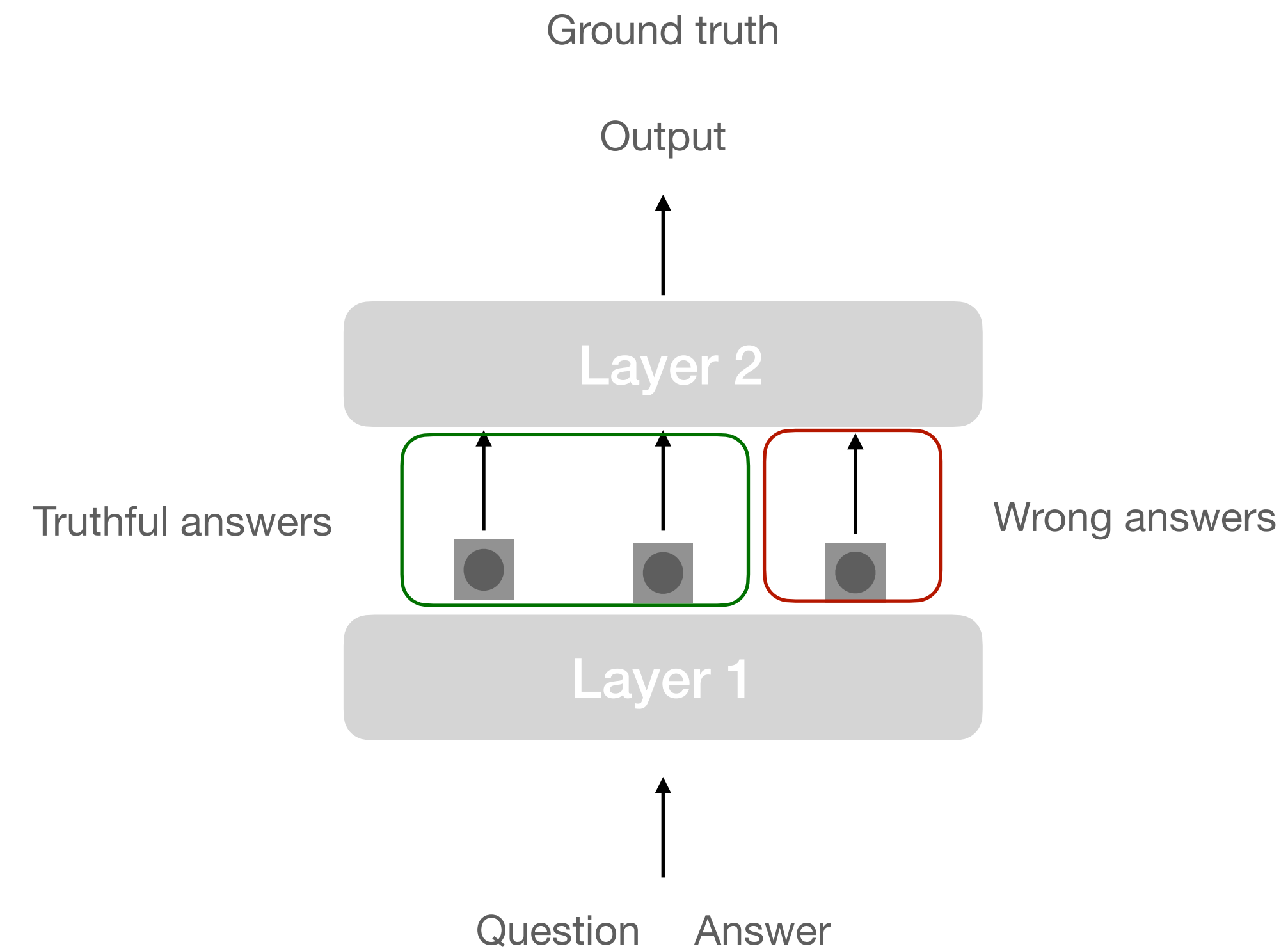
# Workflow
## Probing Truthfulness

- Investigate whether there are vectors in the activation space that correspond to 'truthfulness' or factuality

- For each input

  - Concatenate Q + A

  - Take out head activations to collect 'probing dataset' for each head in each layer

  - Randomly split dataset and train binary classifier.

- 18th head, 14th layer -> Val accuracy of 83.3%

Ground truth

Output

Layer 2

Probe → ▢ ▢ ▢ Attention Heads

Layer 1

Question  Answer

# Workflow
## 'Truthful' Direction

- Mean of the activations associated with truthful and false answers is calculated

- Mass mean shift: Truthful direction determined by vector pointing from mean activation of wrong answers to mean activation of truthful answers.

- Probe Weight Direction:

  - Direction found by linear probing.

  - Equivalent to doing a gradient descent on the head activation to maximize its probability of being predicted as truthful.

- Contrast Consistent Search (CCS) direction

  - Does not require labeled inputs.

  - Train CCS on TruthfulQA by sampling one truthful and one false answer for each question.

# Workflow

## Inference Time Intervention

- During inference Shift activations in the 'truthful' direction.

- Only consider top K heads that appear strongly related to truthfulness.

- Estimate std dev along truthful direction σ using activations from both training and val sets

- Parameters:

  - K - #heads intervened

  - Alpha - Strength of Intervention

  - Determined using standard hyperparameter sweep.

Strength

$$x_{l+1} = x_l + \sum_{h=1}^{H} Q_l^h \left( \text{Att}_l^h(P_l^h x_l) + \alpha \sigma_l^h \theta_l^h \right).$$

0 for non selected heads

# Dataset

- Truthful QA: Operationalise concept of 'truth'

- Adversarially constructed dataset to measure truthfulness of answers generated by LM's

- Reorganise dataset to generate 918 QA pairs each with binary truthfulness label.

- Tracks

  - MCQ

  - Generation

- Note: Goal to focus on specific aspect of truth-telling: avoiding common human misconception.

# Metrics

- Metric used is: true * informative (product of scalar truthful and informative scores)

- Captures #questions answered truthfully and penalises the model from replying 'I have no comment'

- Human annotation is expensive -> use 2 finetuned GPT3- 13B models to classify each answer as true or false or informative or not

# Experiments
## Baseline

- Supervised fine tuning

  - Supervised pre-training on QA pairs and pretraining on Open Web Text.

- Few shot prompting: Do 50 shot prompting on Truthful QA

- Instruction IFT: Study IFTed models from llama -> Alpaca and Vicuna.

# Results

Comparison with baselines that utilize 5% of TruthfulQA to make LLaMA-7B more truthful.

CE is the pre-training loss

KL is the KL divergence between next-token distributions pre- and post-intervention.

MC acc % represents the percentage of instances where the model's confidence in its answer aligns with the actual correctness of the answer.

| | True*Info (%) | True (%) | MC acc. (%) | CE | KL |
|---|---|---|---|---|---|
| Baseline | 30.5 | 31.6 | 25.7 | 2.16 | 0.0 |
| Supervised Finetuning | 36.1 | 47.1 | 24.2 | 2.10 | 0.01 |
| Few-shot Prompting | 49.5 | 49.5 | **32.5** | - | - |
| Baseline + ITI | 43.5 | 49.1 | 25.9 | 2.48 | 0.40 |
| Few-shot Prompting + ITI | **51.4** | **53.5** | **32.5** | - | - |

# Results

It can be applied on top of few-shot prompting or instruction fine-tuning at the cost of a relatively low increases of CE loss and KL divergence.

|  | True*Info (%) | True (%) | MC acc. (%) | CE | KL |
|---|---|---|---|---|---|
| Alpaca | 32.5 | 32.7 | 27.8 | 2.56 | 0.0 |
| Alpaca + ITI | 65.1 | 66.6 | 31.9 | 2.92 | 0.61 |
| Vicuna | 51.5 | 55.6 | 33.3 | 2.63 | 0.0 |
| Vicuna + ITI | 74.0 | 88.6 | 38.9 | 3.36 | 1.41 |

Table 2: Comparison with instruction finetuned baselines using 2-fold cross-validation.

We observe that mass mean shift performs the best and also has a better tolerance for stronger intervention strength.

|  | $\alpha$ | True*Info (%) | True (%) | MC acc. (%) | CE | KL |
|---|---|---|---|---|---|---|
| Baseline | - | 30.5 | 31.6 | 25.7 | 2.16 | 0.0 |
| random direction | 20 | 31.2 | 32.3 | 25.8 | 2.19 | 0.02 |
| CCS direction | 5 | 33.4 | 34.7 | 26.2 | 2.21 | 0.06 |
| ITI: Probe weight direction | 15 | 34.8 | 36.3 | 27.0 | 2.21 | 0.06 |
| ITI: Mass mean shift | 20 | **42.3** | **45.1** | **28.8** | 2.41 | 0.27 |

Table 3: Comparison with different intervention directions and their respective optimal $\alpha$'s. Results are from 2-fold cross-validation, a different protocol from Table 1.

# Results

TruthfulQA is split into 38 subcategories, including misconceptions, stereotypes, history, Mandela effect, and others.

Plot the true*informative scores of all subcategories with 10 or more questions compared to the baseline without intervention.

We observe that ITI increases truthfulness across most types of questions. There is no one category that seems responsible for the overall increase, and we see no clear pattern as to which categories show the biggest effect.

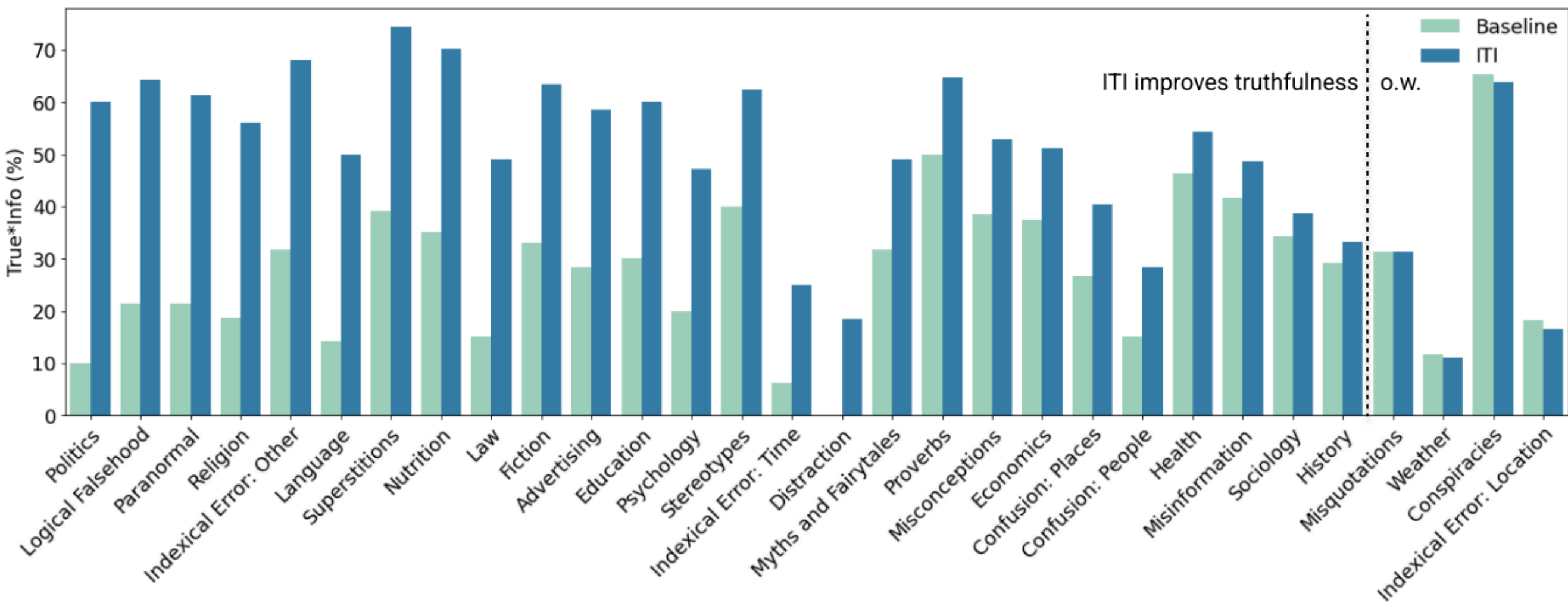## 5.1 Results Across TruthfulQA Categories



Figure 5: True*informative scores split across subcategories on LLaMA-7B, sorted by the difference between baseline and ITI. Subcategories with less than 10 questions are not shown.

# Conclusion and Future Work

- **Introduction of ITI**: ITI enhances language model factuality by adjusting activations towards "truthful" directions, outperforming existing methods on TruthfulQA.

- **Future Research Directions**: ITI's performance on diverse, real-world datasets needs exploration, with a focus on balancing truthfulness and helpfulness.

- **Understanding Representations**: Potential exists for unsupervised identification of "truthful" directions, with an interest in understanding complex attributes like "truth" in model representations.

# My take

- What is the most useful thing that I learned from this paper? Why?

  - I learned that the activation space is interpretable in some way and we can control generation at that level as well.

- What is the one thing I'd do differently if I were the author doing this project?

  - Incorporate more human - evaluation as it is more nuanced and can capture more subtleties in the model output. Eg: The authors talk about 'truthfulness' in one sense only.

- Does the paper help me think about my own project/next project?

  - If prompt engineering doesn't work and we encounter some serious hallucinations in the future, I could try this framework. Also, it makes me think, can we control the generation of vision-based models as well in a similar way?

# Thanks for listening! Questions?