# ProxySPEX: Inference-Efficient Interpretability via Sparse Feature Interactions in LLMs

L. Butler*    A. Agarwal*    J. S. Kang*    Y. E. Erginbas    B. Yu    K. Ramchandran

BAIR — BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

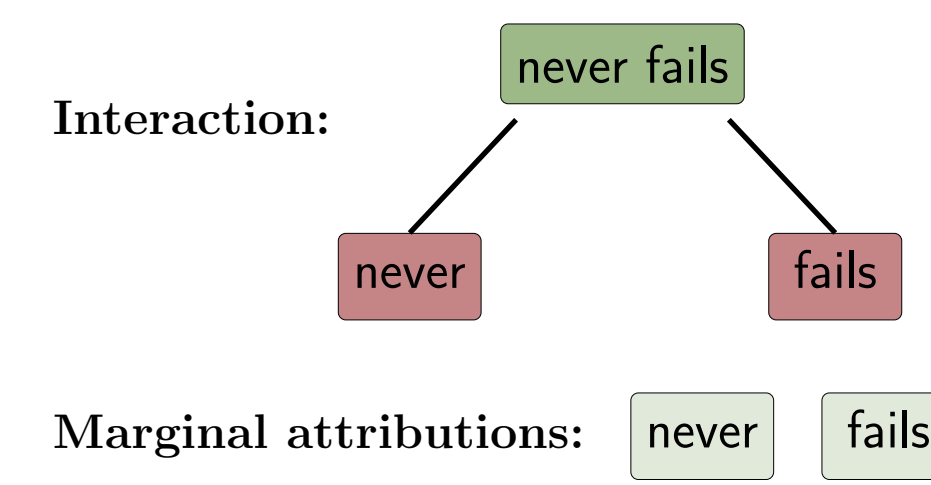Berkeley — UNIVERSITY OF CALIFORNIA

## Problem

How can we **efficiently identify** the **influential feature interactions** in LLMs?

**(a) SENTIMENT ANALYSIS**

CONTEXT
... Her acting never fails to impress. She brings depth and authenticity to every role. Her performances consistently draw the ...

PROMPT
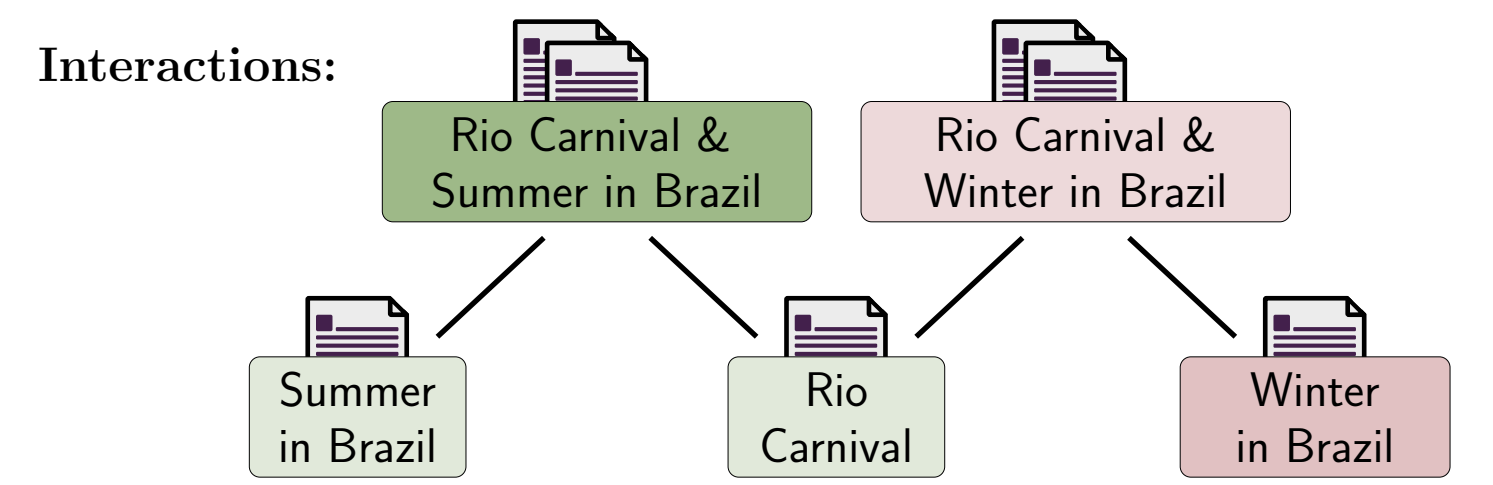*Is this a positive or negative review?*

GENERATED RESPONSE
*Positive.*

Interaction: never fails → never, fails

Marginal attributions: never, fails

**(b) RETRIEVAL AUGMENTED GENERATION**

CONTEXT
... Weather in Tokyo, Brazilian Music, Rio Carnival, Summer in Brazil, Winter in Brazil, History of Brazil, Sport in Rio ...

PROMPT
*What is the weather like during Rio Carnival?*

GENERATED RESPONSE
*Rio Carnival generally takes place during the summer season in Brazil. The weather at this time is typically hot and humid.*

Interactions: Rio Carnival & Summer in Brazil, Rio Carnival & Winter in Brazil → Summer in Brazil, Rio Carnival, Winter in Brazil
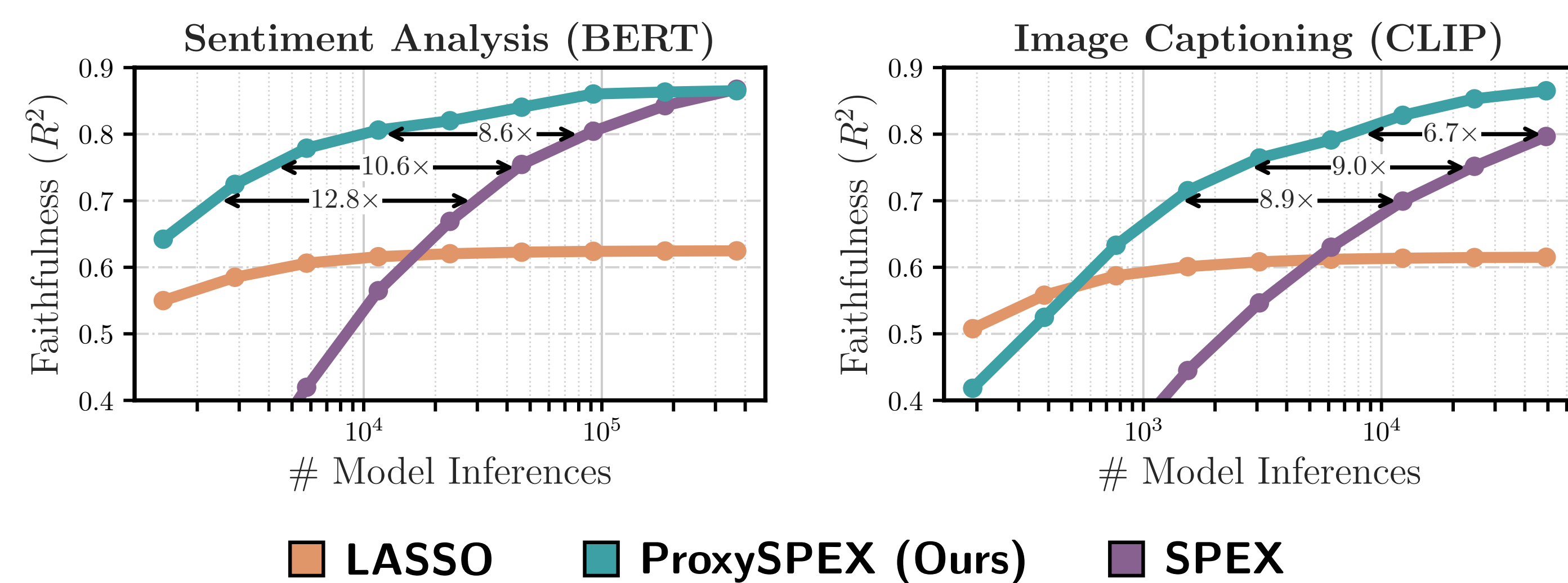
- **Examples:** double negatives in sentiment analysis tasks and multi-document understanding in question answering tasks.
- Marginal attribution approaches like SHAP/LIME scale, but don't capture important interactions.
- A prior approach (SPEX) scales, but still requires tens of thousands of model inferences, which can be prohibitive for complex models such as LLMs.

## Faithfulness at Scale

- For input $\mathbf{x}$ = "Her acting fails to impress", let $f(\mathbf{x}_S)$ be the output of the LLM under *masking pattern* $S$.
- If $S = \{1, 2, 4, 5, 6\}$, then $\mathbf{x}_S$ is "Her acting [MASK] fails to impress". This masking pattern changes the sentiment score from positive to negative.
- We aim to learn an interpretable approximate function $\hat{f}$ that is faithful to the original function $f$, measured in terms of $R^2$:

$$R^2 = 1 - \frac{\|\hat{f} - f\|^2}{\|f - \bar{f}\|^2}, \quad \text{where } \|f\|^2 = \sum_{S \subseteq [n]} f(S)^2, \bar{f} = \frac{1}{2^n} \sum_{S \subseteq [n]} f(S).$$

**Result:** ProxySPEX requires ~10× fewer inferences to achieve equally faithful explanations as SPEX.



Sentiment Analysis (BERT) / Image Captioning (CLIP) — Faithfulness ($R^2$) vs # Model Inferences. LASSO, ProxySPEX (Ours), SPEX.
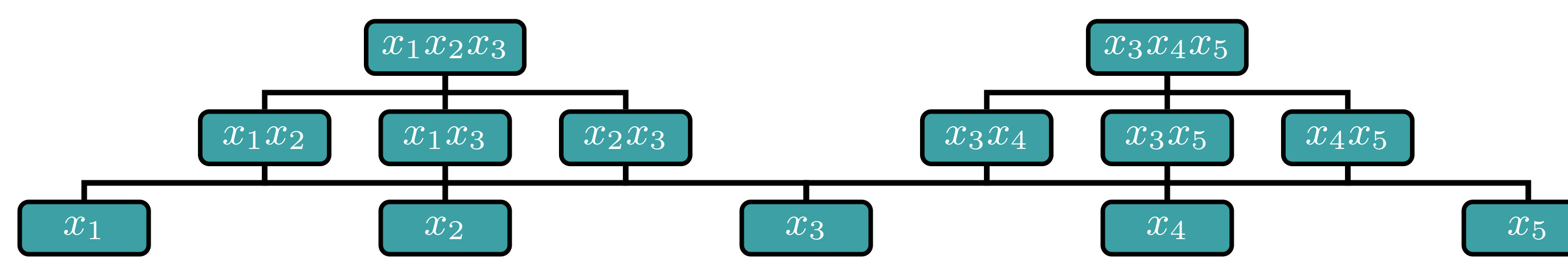
## Fourier Sparsity and Spectral Hierarchies

- Every function $f(\mathbf{x}_S)$ has a unique decomposition under the Fourier transform, expressed as:
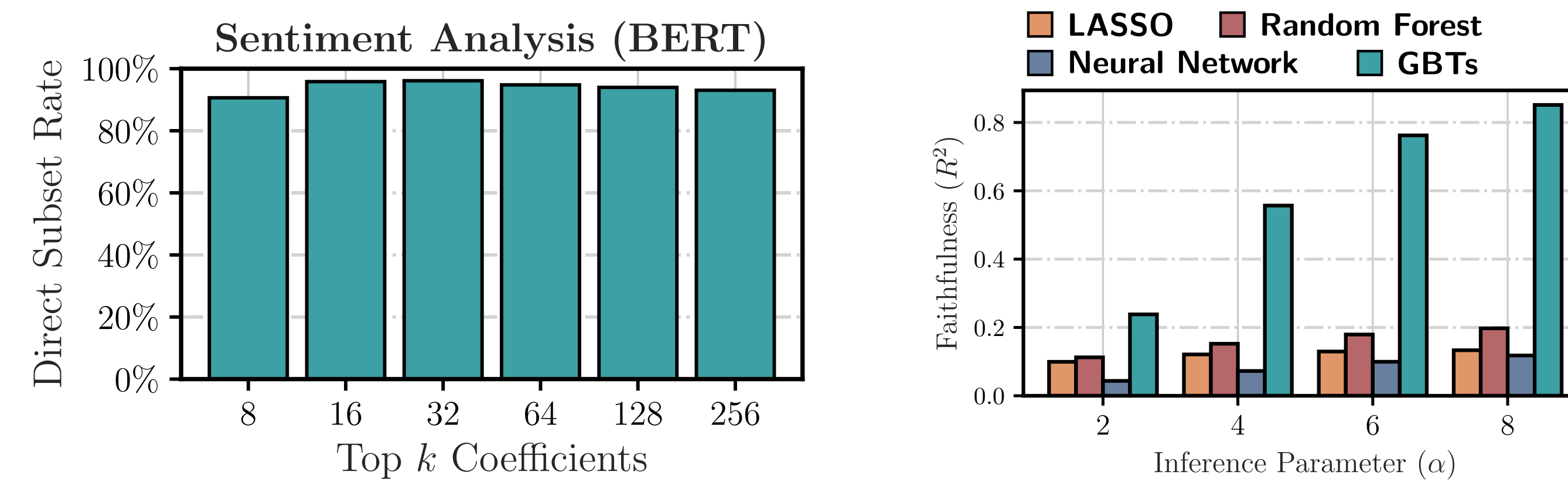
$$F(\mathbf{x}_T) = \frac{1}{2^n} \sum_{S \subseteq [n]} (-1)^{|S \cap T|} f(\mathbf{x}_S), \qquad f(\mathbf{x}_S) = \sum_{T \subseteq [n]} (-1)^{|S \cap T|} F(\mathbf{x}_T).$$

- It has been observed that $F(\mathbf{x}_T) \approx 0$ for most $T$ (**sparsity**), and most large $F(\mathbf{x}_T)$ are **low degree** such that $|T| \leq d$ for some small $d$.
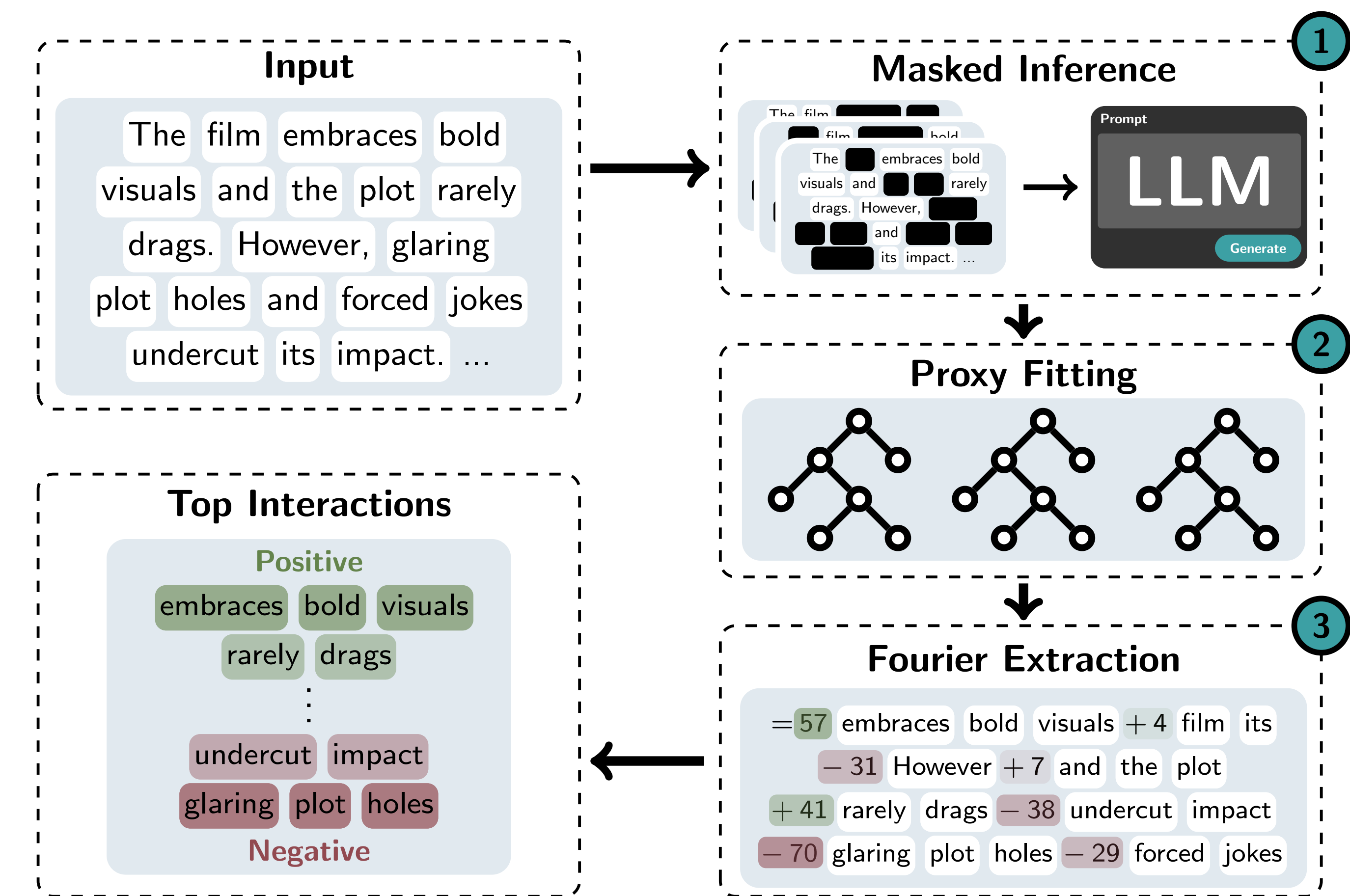
In addition to **sparse** and **low-degree**, influential interactions are **hierarchical**: higher-order interactions are accompanied by their lower-order subsets.



(Left) Direct subset rate measures the rate at which a top-$k$ interaction has a lower-order subset also contained in the top-$k$. (Right) Gradient Boosted Trees efficiently recover sparse, hierarchical interactions.
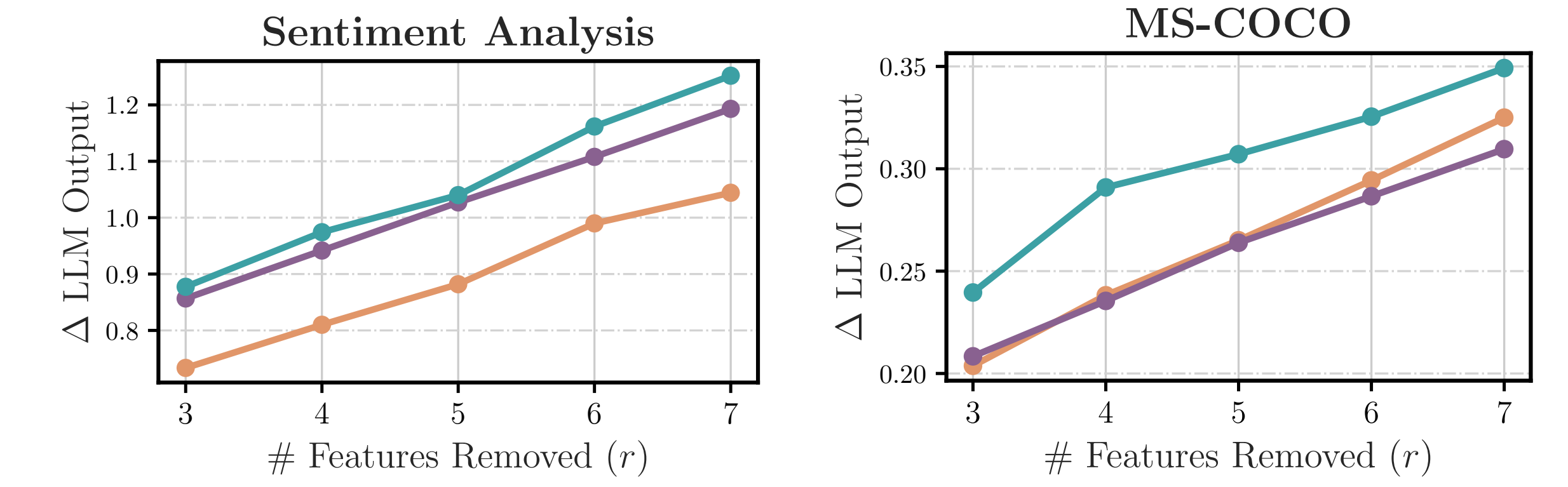
Sentiment Analysis (BERT) — Direct Subset Rate vs Top $k$ Coefficients. Faithfulness ($R^2$) vs Inference Parameter ($\alpha$). LASSO, Random Forest, Neural Network, GBTs.

## ProxySPEX Algorithm



**Input** → **Masked Inference** ① → **Proxy Fitting** ② → **Fourier Extraction** ③ → **Top Interactions**

Input: The film embraces bold visuals and the plot rarely drags. However, glaring plot holes and forced jokes undercut its impact. ...

Top Interactions:
Positive: embraces bold visuals, rarely drags
Negative: undercut impact, glaring plot holes

Fourier Extraction:
= 57 embraces bold visuals + 4 film its − 31 However + 7 and the plot + 41 rarely drags − 38 undercut impact − 70 glaring plot holes − 29 forced jokes

(1) ProxySPEX masks subsets of words and queries the LLM using this masked input. (2) It then fits Gradient Boosted Trees as a proxy model to learn the LLM's hierarchical interactions. (3) A sparse representation is extracted from the fitted GBTs, capturing influential interactions.
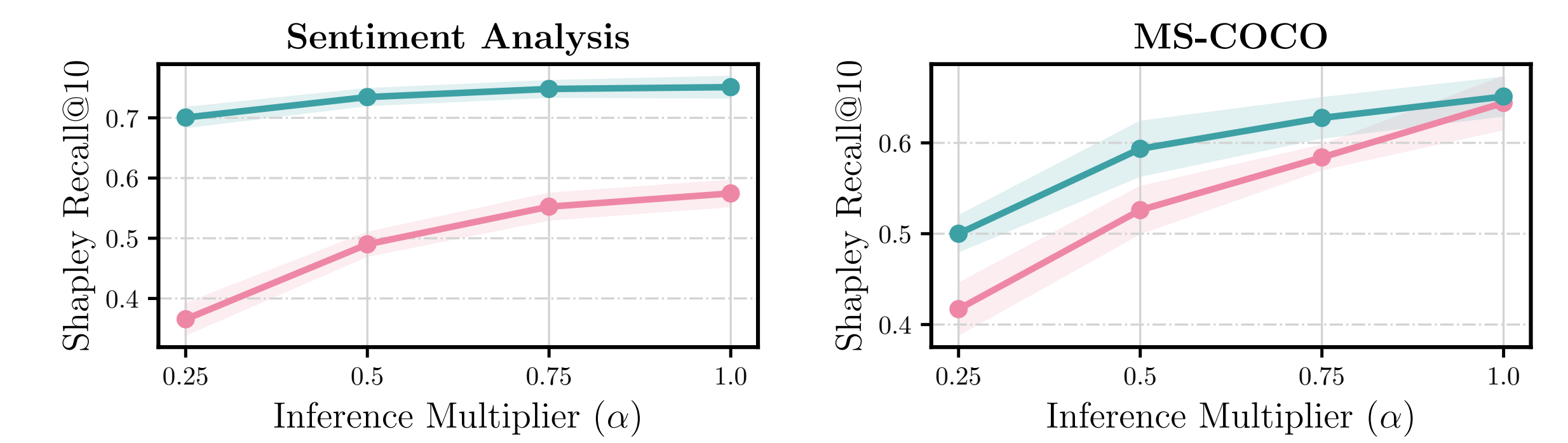
## Feature Removal



Sentiment Analysis / MS-COCO — Δ LLM Output vs # Features Removed ($r$). LASSO, ProxySPEX (Ours), SPEX.

By accounting for interactions, ProxySPEX identifies more influential features across datasets than the LASSO and SPEX.
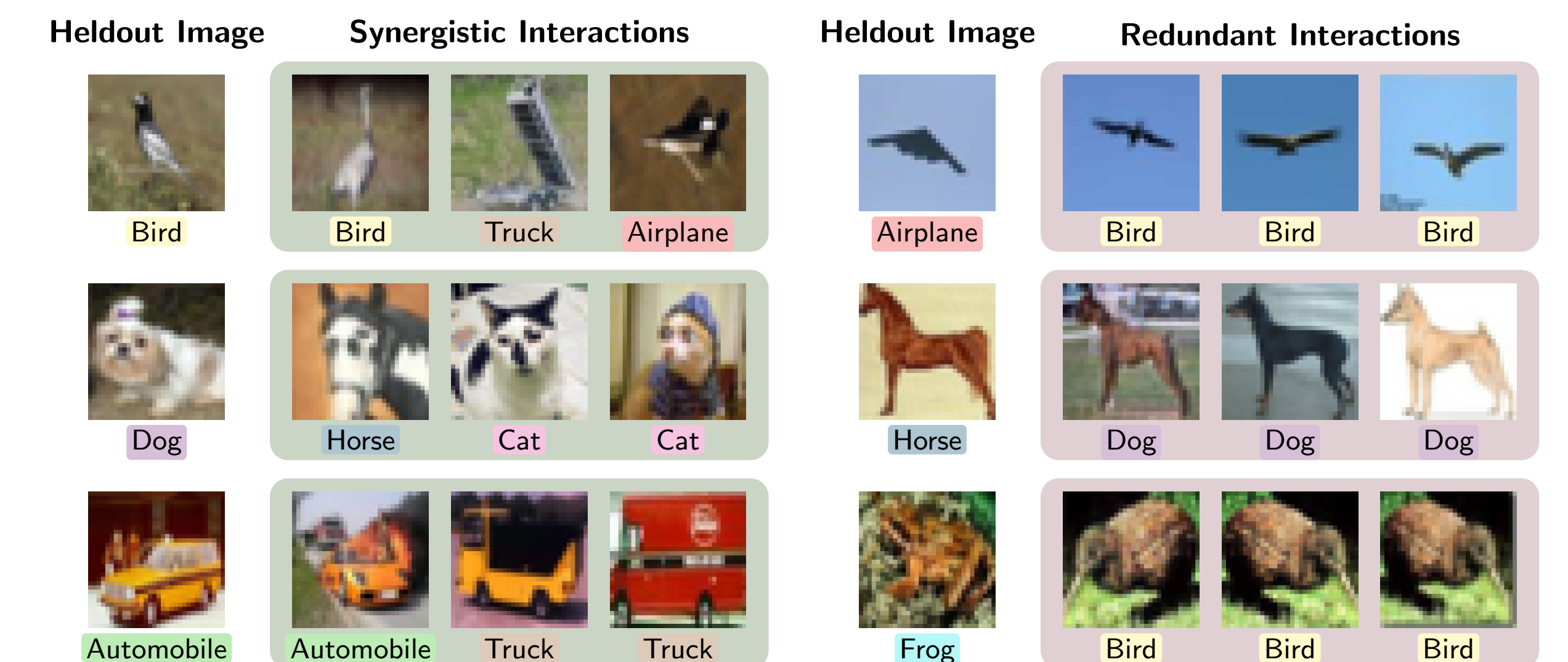
## Sample-Efficient Shapley Estimation



Sentiment Analysis / MS-COCO — Shapley Recall@10 vs Inference Multiplier ($\alpha$). KernelSHAP, ProxySPEX (Ours).

For multipliers $\alpha \in \{0.25, 0.5, 0.75, 1.0\}$, recall of the top ten Shapley values after $\alpha \cdot n \log_2(n)$ inferences. For small $\alpha$, ProxySPEX is **superior at recovering the most significant features**, while KernalSHAP outperforms as $\alpha$ increases.
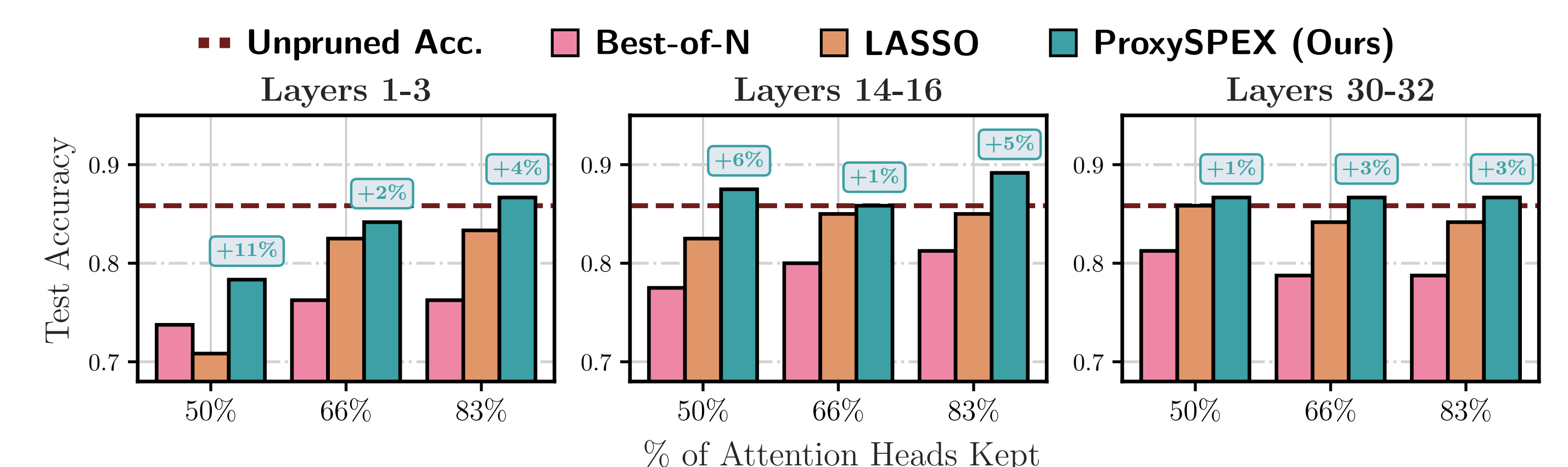
## Data Interaction Attribution

Data attribution measures how each training sample influences the prediction of a particular test point. We extend this framework to capture interactions.



Heldout Image / Synergistic Interactions / Heldout Image / Redundant Interactions

**Synergistic interactions:** data that together are more valuable together than the sum of their parts. **Redundant interactions:** Combined influence is less than the sum of the parts.

## Attention Head Interaction Attribution



Layers 1-3 / Layers 14-16 / Layers 30-32 — Test Accuracy vs % of Attention Heads Kept. Unpruned Acc., Best-of-N, LASSO, ProxySPEX (Ours).

Attention head pruning for Llama-3.1-8B-Instruct for MMLU (high-school-us-history) across different layers. Unpruned accuracy shown by dashed line.