# SPEX: Scaling Feature Interaction Explanations for LLMs

J. S. Kang[1*]    L. Butler[1*]    A. Agarwal[2*]    Y. E. Erginbas[1]    R. Pedarsani[3]    B. Yu[12]    K. Ramchandran[1]

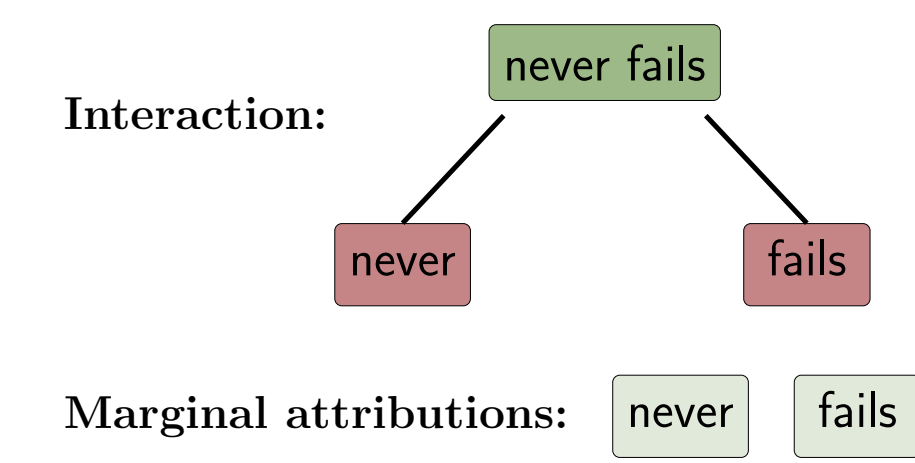[1]UC Berkeley EECS    [2]UC Berkeley Statistics    [3]UC Santa Barbara    *Equal Contribution

## Problem

LLMs identify important interactions between inputs. Can signal processing and information theory help efficiently identify these interactions using only query access to the LLM?



(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
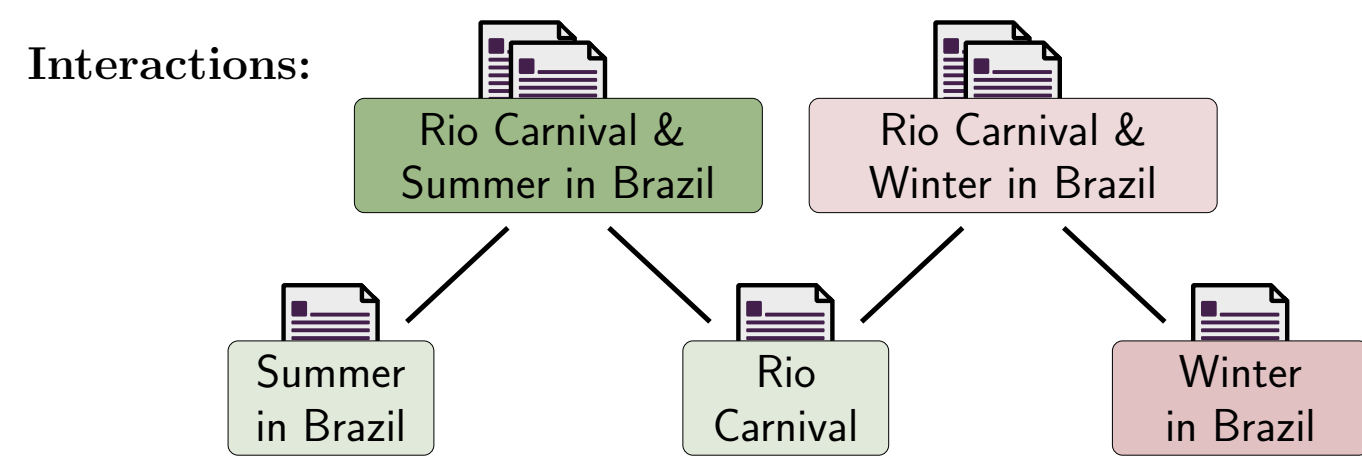
Marginal attributions: never, fails

(b) RETRIEVAL AUGMENTED GENERATION

CONTEXT
... Weather in Tokyo, Brazilian Music, Rio Carnival, Summer in Brazil, Winter in Brazil, History of Brazil, Sport in Rio ...
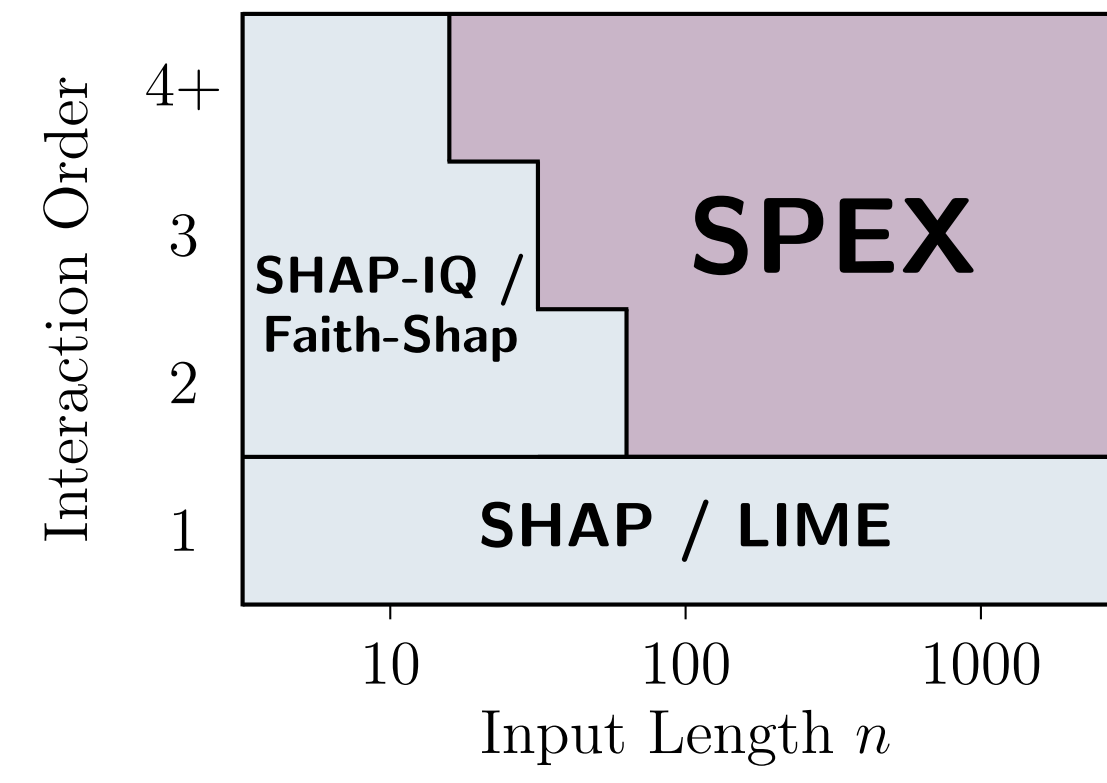
PROMPT
What is the weather like during Rio Carnival?

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

Example: Tasks can require using interactions between inputs to generate responses.

- Marginal approaches like SHAP/LIME scale, but don't capture important interactions.
- Existing interaction identification approaches are too slow to scale for practical LLM input sizes.
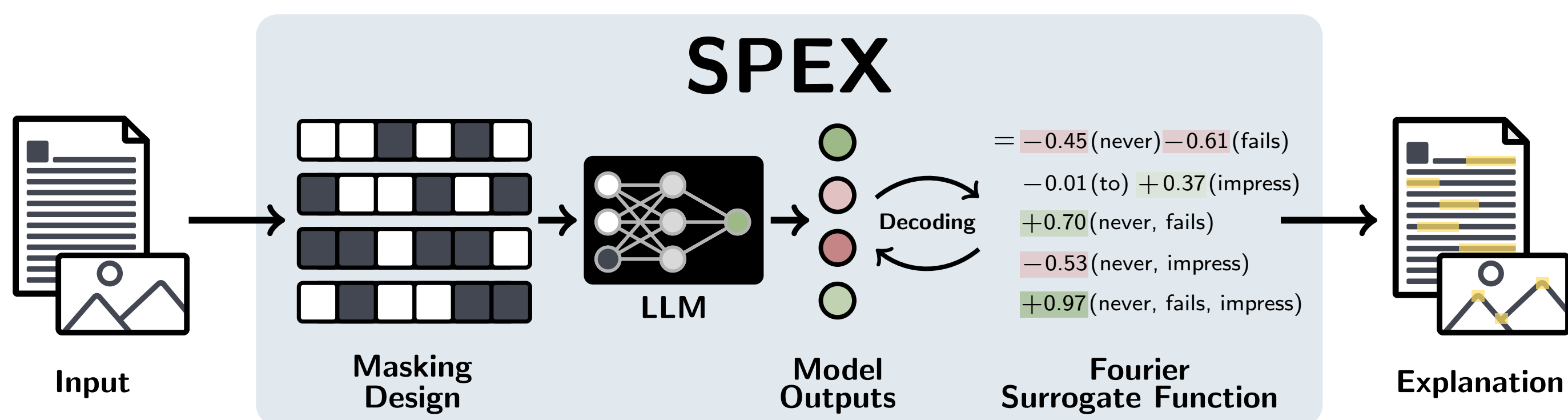- Our approach, SPEX, scales to large inputs *and* captures interactions.



## Formulation as Fourier Transform

- 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 $S = \{3\}$, then $\mathbf{x}_S$ is "Her acting [MASK] fails to impress", this masking pattern changes the score from positive to negative.
- Equivalently write $f : \mathbb{F}_2^n \to \mathbb{R}$, where $f(\mathbf{x}_S) = f(\mathbf{m})$ with $S = \{i : m_i = 1\}$. Then the Fourier transform is defined as follows:

Forward: $F(\mathbf{k}) = \frac{1}{2^n} \sum_{\mathbf{m} \in \mathbb{F}_2^n} (-1)^{\langle \mathbf{k}, \mathbf{m} \rangle} f(\mathbf{m})$    Inverse: $f(\mathbf{m}) = \sum_{\mathbf{k} \in \mathbb{F}_2^n} (-1)^{\langle \mathbf{m}, \mathbf{k} \rangle} F(\mathbf{k})$.

We find that $F(\mathbf{k}) \approx 0$ for most $\mathbf{k}$ (sparsity), and most large $F(\mathbf{k})$ are low degree such that $|\mathbf{k}| \leq d$ for some small $d$.

- SPEX exploits this sparsity using codes, to compute interactions efficiently, by computing estimates $\hat{F}(\mathbf{k})$ for a small (a-priori unknown) set of $\mathbf{k} \in \mathcal{K}$.
- Inverting our estimated $\hat{F}(\mathbf{k})$ gives us an approximate surrogate function $\hat{f}$.



SPEX

= −0.45 (never) −0.61 (fails)
  −0.01 (to) +0.37 (impress)
  +0.70 (never, fails)
  −0.53 (never, impress)
  +0.97 (never, fails, impress)

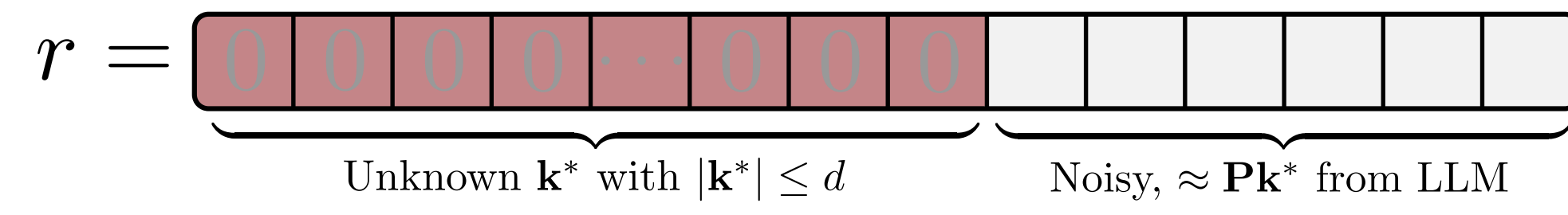Input → Masking Design → LLM → Model Outputs → Decoding → Fourier Surrogate Function → Explanation

SPEX utilizes codes to determine masking patterns. We observe the changes in model output depending on the used mask. SPEX uses message passing to learn Fourier coefficients to generate interaction-based explanations.

## Algorithm

### Step 1: Masking Design - Embedding Code Structures Through Aliasing

- We collect samples according to two matrices $\mathbf{M} \in \mathbb{F}_2^{b \times n}$ and $\mathbf{P} \in \mathbb{F}_2^{p \times n}$.
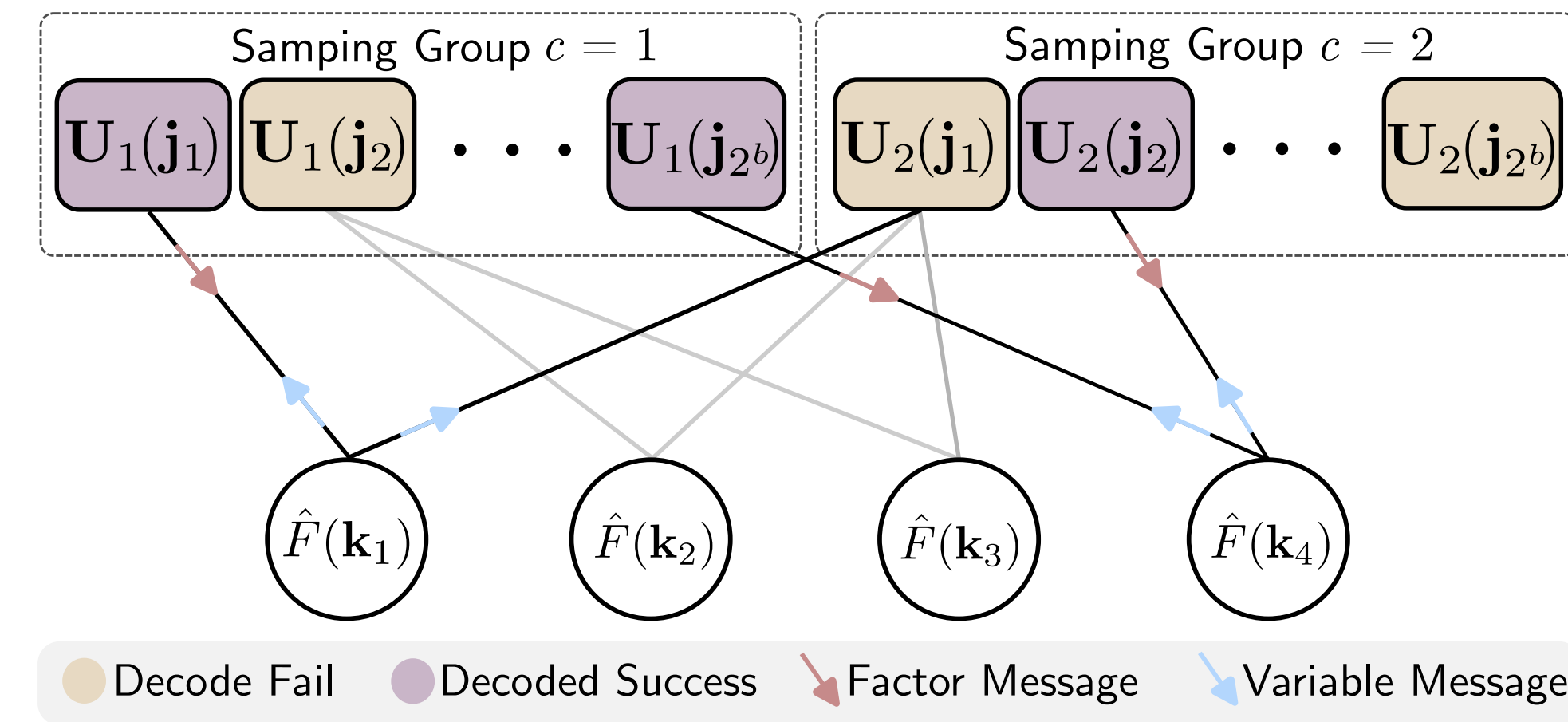
$$u_{c,i}(\boldsymbol{\ell}) = f(\mathbf{M}_c^\top \boldsymbol{\ell} + \mathbf{p}_i) \iff U_{c,i}(\mathbf{j}) = \sum_{\mathbf{k} \,:\, \mathbf{M}_c \mathbf{k} = \mathbf{j}} (-1)^{\langle \mathbf{p}_i, \mathbf{k} \rangle} F(\mathbf{k}).$$

- Depending on $\mathbf{p}_i$, the modulation $(-1)^{\langle \mathbf{p}_i, \mathbf{k} \rangle}$ changes the sign of $F(\mathbf{k})$.
- Each $U_{c,i}(\mathbf{j})$ can be seen as a noisy BPSK message containing a codeword $\mathbf{P}\mathbf{k}^*$ conveying a dominant $\mathbf{k}^*$ in the sum above.

$$r = \boxed{\phantom{xxxx}}$$

Unknown $\mathbf{k}^*$ with $|\mathbf{k}^*| \leq d$ — Noisy, $\approx \mathbf{P}\mathbf{k}^*$ from LLM

- If $\mathbf{P}$ is a parity matrix of a systematic code, we can decode $r$ to recover dominant $\mathbf{k}^*$. This can be seen as a form of *joint source channel coding*.
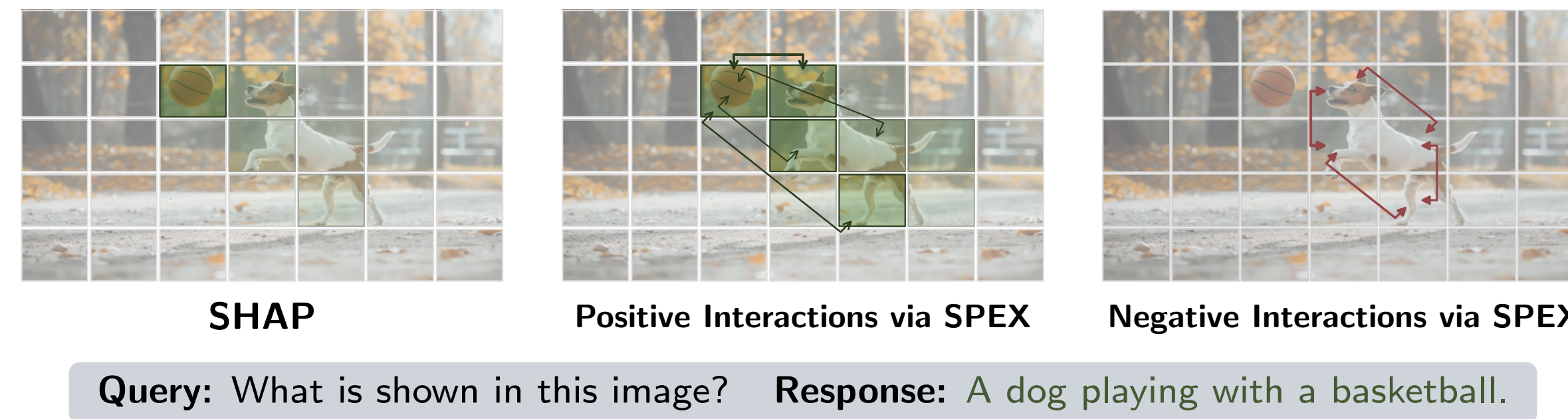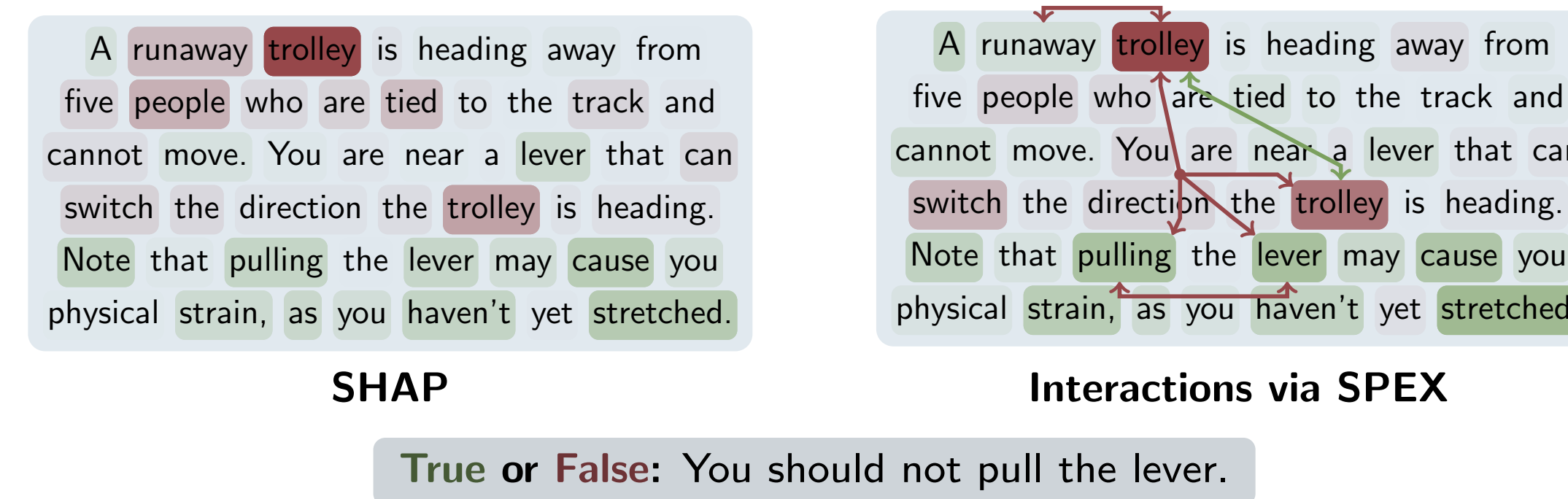
### Step 2: Message Passing - Decoding and Interference Cancellation

- Defines a bipartite graph connecting the non-zero $F(\mathbf{k})$ and $U$.
- As we recover $\hat{F}(\mathbf{k})$ and $\mathbf{k}$, we can do interference cancellation via message passing. This is inspired by sparse graph codes for robust communication.



Samping Group $c = 1$: $\mathbf{U}_1(\mathbf{j}_1)$ $\mathbf{U}_1(\mathbf{j}_2)$ $\cdots$ $\mathbf{U}_1(\mathbf{j}_{2^b})$
Samping Group $c = 2$: $\mathbf{U}_2(\mathbf{j}_1)$ $\mathbf{U}_2(\mathbf{j}_2)$ $\cdots$ $\mathbf{U}_2(\mathbf{j}_{2^b})$

$\hat{F}(\mathbf{k}_1)$ $\hat{F}(\mathbf{k}_2)$ $\hat{F}(\mathbf{k}_3)$ $\hat{F}(\mathbf{k}_4)$

Decode Fail · Decoded Success · Factor Message · Variable Message

- We can analyze the message passing with density evolution theory.

## Case Studies: Applications



A runaway trolley is heading away from five people who are tied to the track and cannot move. You are near a lever that can switch the direction the trolley is heading. Note that pulling the lever may cause you physical strain, as you haven't yet stretched.

SHAP

A runaway trolley is heading away from five people who are tied to the track and cannot move. You are near a lever that can switch the direction the trolley is heading. Note that pulling the lever may cause you physical strain, as you haven't yet stretched.

Interactions via SPEX

True or False: You should not pull the lever.



SHAP · Positive Interactions via SPEX · Negative Interactions via SPEX

Query: What is shown in this image?    Response: A dog playing with a basketball.
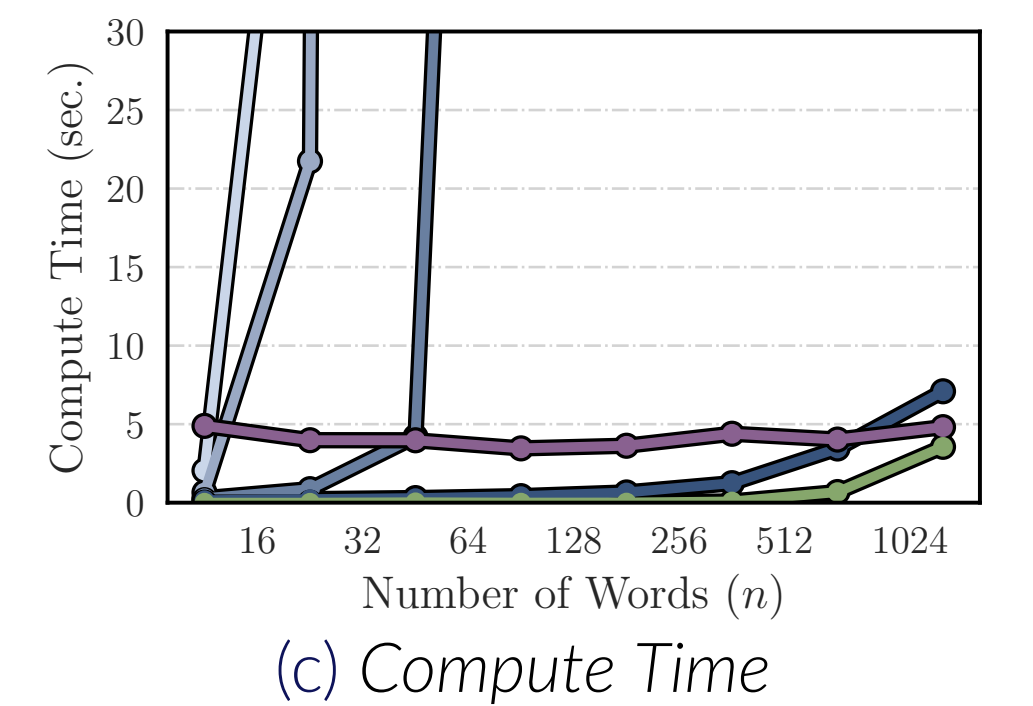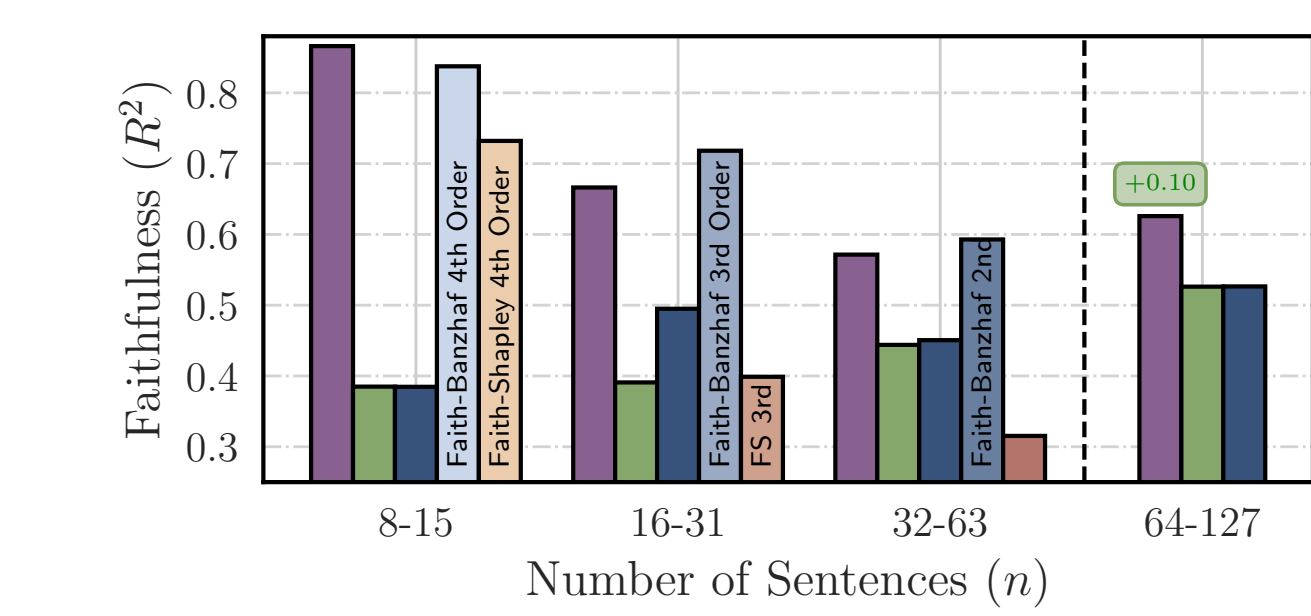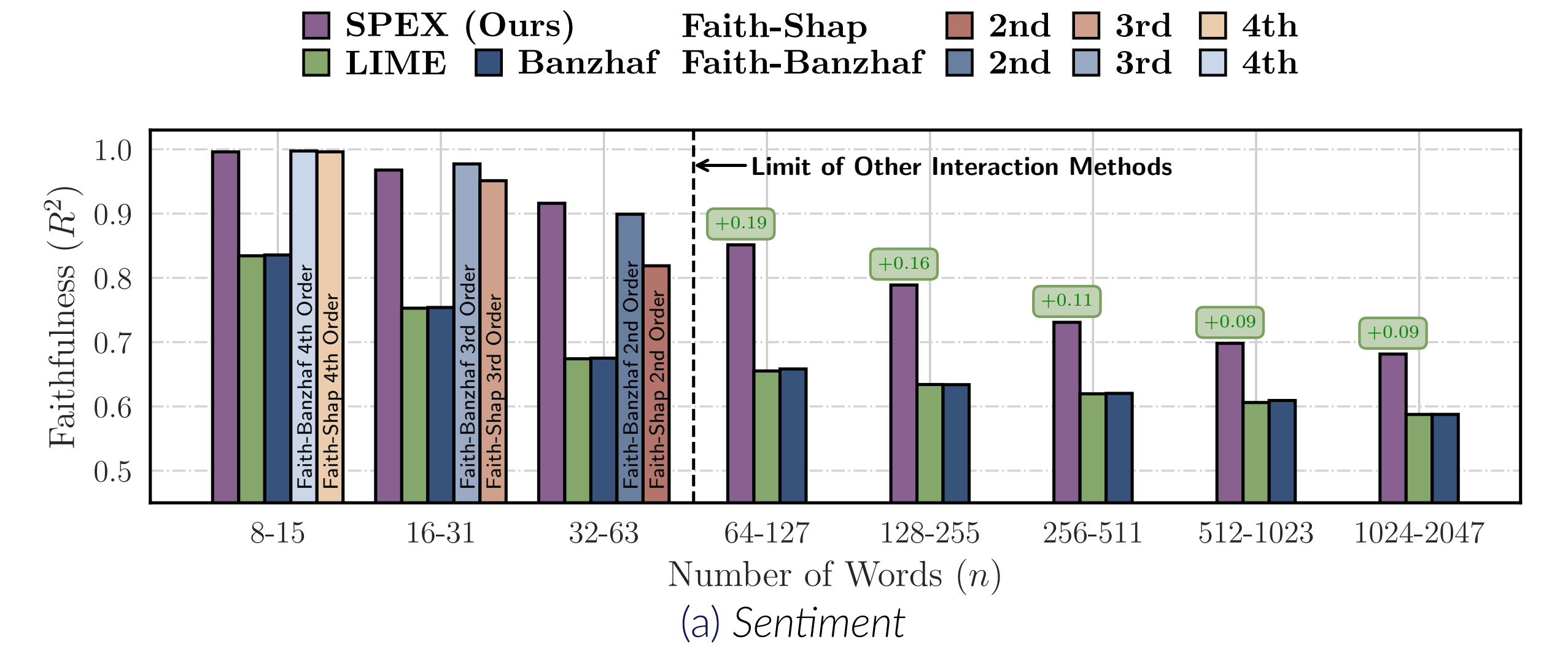
Abstract Reasoning Errors: LLMs struggle with modified versions of puzzle questions. We consider a variant of the classic trolley problem. *GPT-4o mini* incorrectly answers. We identify a strong interaction between words that commonly appear in the standard problems.
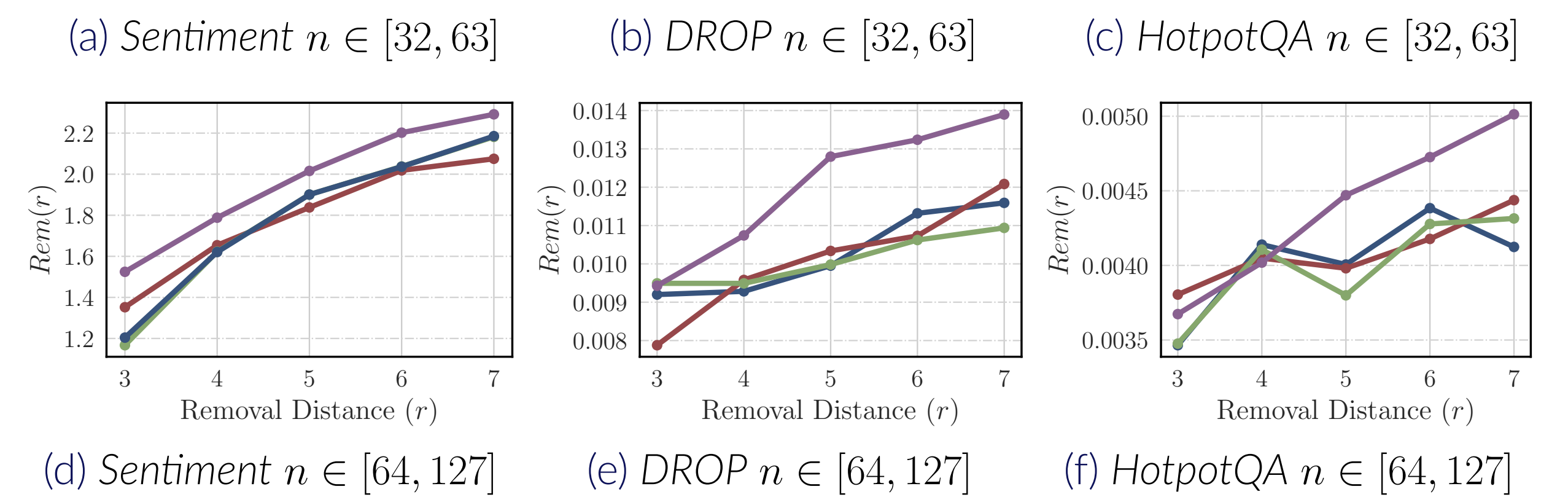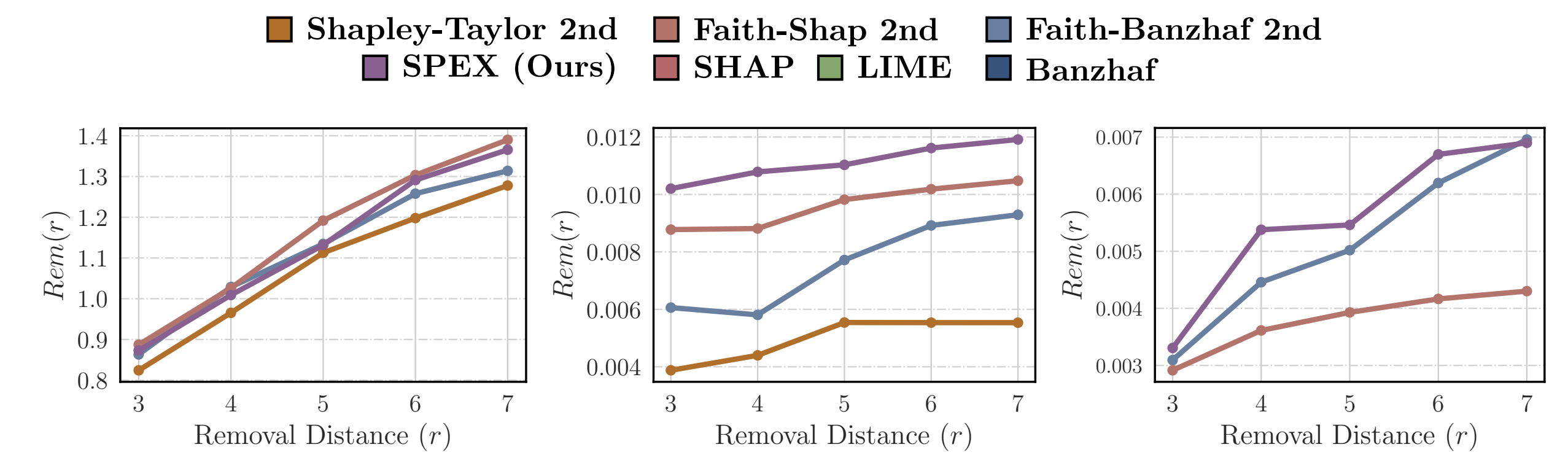
Visual Question Answer: We prompt *LLaVA-NeXT-Mistral* with *"What is shown in this image?"* for the image above. SHAP indicates the importance of image patches containing the ball and the dog. SPEX shows that the presence of *both* the dog and the basketball jointly are critical.
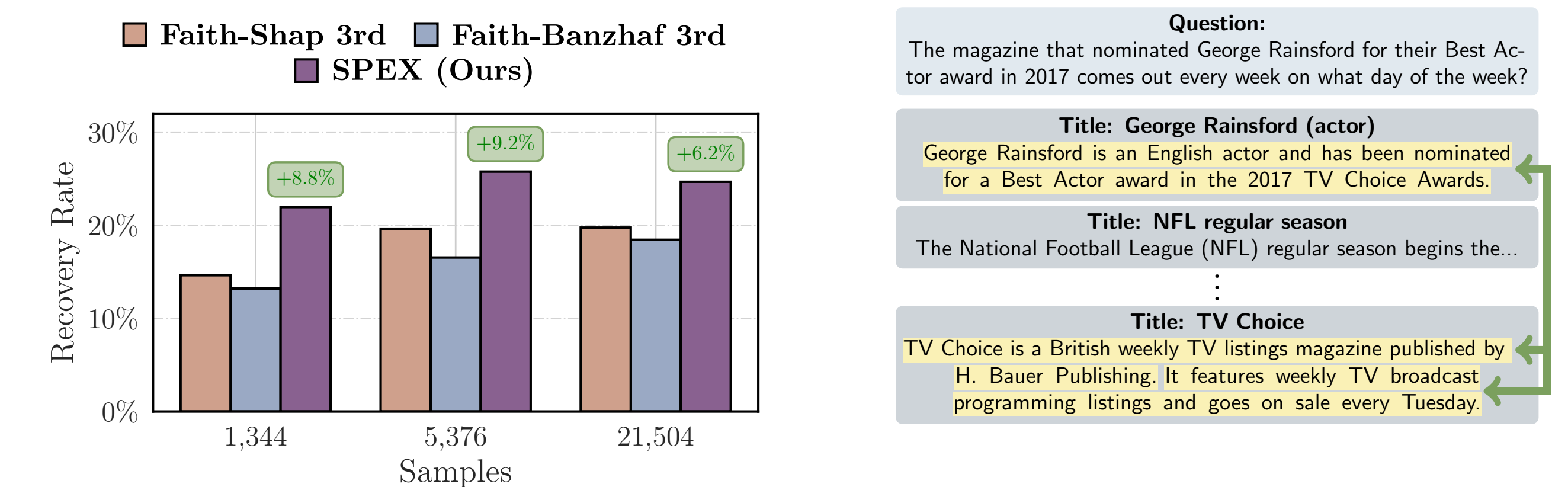
## Experiments



(a) Sentiment



(b) HotpotQA      (c) Compute Time

Faithfulness: *Faithfulness* to the real function $f$, defined in terms of $R^2$:

$$R^2 = 1 - \frac{\|\hat{f} - f\|^2}{\|f - \bar{f}\|^2}, \quad \|f\|^2 = \sum_{\mathbf{m} \in \mathbb{F}_2^n} f(\mathbf{m})^2 \quad \bar{f} = \frac{1}{2^n} \sum_{\mathbf{m} \in \mathbb{F}_2^n} f(\mathbf{m}).$$



(a) Sentiment $n \in [32, 63]$      (b) DROP $n \in [32, 63]$      (c) HotpotQA $n \in [32, 63]$

(d) Sentiment $n \in [64, 127]$      (e) DROP $n \in [64, 127]$      (f) HotpotQA $n \in [64, 127]$

Top-$r$ Removal: We identify the top $r$ influential features to model output:

$$\text{Rem}(r) = \frac{|f(\mathbf{1}) - f(\mathbf{m}^*)|}{|f(\mathbf{1})|}, \quad \mathbf{m}^* = \arg\max_{|\mathbf{m}| = n - r} |\hat{f}(\mathbf{1}) - \hat{f}(\mathbf{m})|.$$



Question:
The magazine that nominated George Rainsford for their Best Actor award in 2017 comes out every week on what day of the week?

Title: George Rainsford (actor)
George Rainsford is an English actor and has been nominated for a Best Actor award in the 2017 TV Choice Awards.

Title: NFL regular season
The National Football League (NFL) regular season begins the...

Title: TV Choice
TV Choice is a British weekly TV listings magazine published by H. Bauer Publishing. It features weekly TV broadcast programming listings and goes on sale every Tuesday.

(Left) Recovery of Human-labeled interactions in *HotpotQA*. (Right) Example interaction.

Recovery Rate@$r$: Let $S_r^* \subseteq [n]$ denote human-annotated sentence. Let $S_i$ denote feature indices of the $i^{th}$ most important interaction.

$$\text{Recovery@}r = \frac{1}{r} \sum_{i=1}^{r} \frac{|S_r^* \cap S_i|}{|S_i|}.$$