

## Instruction

Describe the fear of flying.

## LLM



## Explanation Target

The fear of flying, also known as aviophobia or aerophobia, is a ...

The fear of flying, also known as aviophobia or aerophobia, is a ...

or Some reference text

B = LLM Base response

A = Aspect

R = Reference text



embeddings

Emb·t

### 1 CONCEPT EXTRACTION

Describe the fear of flying.

Describe, fear, flying



GPT-4o mini

### 2 CONCEPT REPLACEMENT

Mention, concern, journey

**Replacement Strategy**  
 $r$  = remove  
 $n$  = neutral  
 $a$  = antonym

### 3 CONCEPT IMPORTANCE ESTIMATION: $\phi(\text{fear})$

#### Coalitions

$S_0 = \emptyset$

$S_1 = \{\text{describe}\}$

$S_2 = \{\text{flying}\}$

$S_3 = \{\text{describe, flying}\}$

#### S

Mention the concern of journey.

Describe the concern of journey.

Mention the concern of flying.

Describe the concern of flying.

#### SU{fear}

Mention the fear of journey.

Describe the fear of journey.

Mention the fear of flying.

Describe the fear of flying.

#### Value Function

(repeat for N concepts)

$\cos(\text{Emb} \cdot \text{LLM}(\text{S}), \text{Emb} \cdot \text{t})$

$\cos(\text{Emb} \cdot \text{LLM}(\text{SU}\{\text{fear}\}), \text{Emb} \cdot \text{t})$

Attribution  $\phi(\text{fear})$



### ConceptX Explanation

Describe the fear of flying.

Describe the fear of flying.

Describe the fear of flying.

if t = LLM Base response

if t = Aspect

if t = Reference text