# Stochastic Parakeets , Peacocks , and Penguins : Irrelevant Context Hallucinations Reveal Class-Based (Mis)Generalization

Anonymous ACL submission

# Abstract

The widespread success of LLMs on NLP benchmarks has been accompanied by concerns that LLMs function primarily as stochastic parrots that reproduce texts similar to what they saw during pre-training, often erroneously. But what is the nature of their errors, and do these errors exhibit any regularities? In this work, 007 we examine irrelevant context hallucinations, in which models integrate misleading contextual cues into their predictions. Through be-011 havioral analysis, we show that these errors result from a structured yet flawed mechanism 012 that we term *class-based* (*mis*)generalization, in which models combine abstract class cues with features extracted from the query or context to derive answers. Furthermore, mechanistic interpretability experiments on Llama-3, 017 Mistral, and Pythia across 39 factual recall relation types reveal that this behavior is reflected in the model's internal computations: (i) abstract class representations are constructed in lower layers before being refined into specific answers in higher layers, (ii) feature selection is governed by two competing circuits - one prioritizing direct query-based reasoning, the other incorporating contextual cues - whose 027 relative influences determine the final output. Our findings provide a more nuanced perspective on the stochastic parrot argument: through form-based training, LLMs can exhibit sensitivity to structured abstractions, albeit in unreliable ways.<sup>1</sup>

# 1 Introduction

The remarkable success of LLMs on various NLP benchmarks has been accompanied by concerns that they function primarily as "stochastic parrots" that operate by "haphazardly stitching together sequences of linguistic forms" using statistical cooccurrences in pre-training data (Bender et al., 2021). This view is supported by evidence that



Figure 1: Examples demonstrating class-based (mis)generalization with Llama-3 (8B).

LLMs can reproduce training artifacts, exploit spurious correlations, and fail when faced with distribution shifts, among other issues (Carlini et al., 2021; Zhou et al., 2024; Dziri et al., 2023; Wu et al., 2024c; Mirzadeh et al., 2024).

In this work, we argue that more deeply examining model errors can reveal insights into LLM behaviors and generalization capabilities. In particular, we examine a specific and underexplored type of error — **irrelevant context hallucinations** — to investigate the mechanisms through which LLMs integrate contextual information into their predictions. We introduce a controlled experimental setting where LLMs receive irrelevant contextual information alongside a query (Figure 1).

By artificially controlling the context and query pairing, this setup allows us to explore how LLMs behave in situations they were unlikely to have encountered in pre-training. Additionally, by focusing on incorrect answers, we sidestep the data leakage issue, which is primarily concerned with memorization of correct answers (Balloccu et al.,

<sup>&</sup>lt;sup>1</sup>Code will be released upon acceptance.

106

107

108

109

110

111

112

113

114

063

064

2024; Xu et al., 2024). Thus, these controls reduce the possibility that answers stem purely from pattern matching from pre-training data, allowing us to better isolate prediction shifts driven by added context.

Through qualitative analysis of irrelevant context hallucinations, as demonstrated in Figure 1, we hypothesize that these errors exhibit structured regularities. We posit that LLMs exhibit a structured but flawed mechanism, which we term the class-based (mis)generalization hypothesis. Specifically, LLMs can leverage abstract class cues (e.g., "language"), use them to select features in the prompt (e.g., selecting the *country* feature of "Honda", instead of the year feature), and combine these abstract classes with the selected features to produce an answer (e.g., "Language" + "Japan"  $\rightarrow$ "Japanese"). This hypothesis suggests that LLMs can generalize in a systematic and structured manner in this setting, but as we will show, their reliance on these abstractions is often flawed. In some cases, it leads to correct answers via an incorrect computation (e.g., "Portuguese" in Figure 1), while in others, it results in hallucinations (e.g., "Japanese", "Norwegian" in Figure 1).

To validate our hypothesis, we conduct a behavioral analysis of how irrelevant context influences model predictions on Llama-3, Mistral and Pythia. Specifically, we perform annotations on 500 data points and show that 70% of observed shifts pattern with our class-based generalization hypothesis. Moreover, statistical analyses confirm that this phenomenon is systematic rather than due to chance or being query-dependent.

We provide further evidence of this generalization mechanism via mechanistic interpretability experiments which probe the model's internal computations across Transformer layers (Vaswani et al., 2017). Our findings reveal two key mechanisms that further support our hypothesis: (i) LLMs make hierarchical class-to-instance predictions; i.e., they construct abstract class representations (e.g., "languages") before refining them to more specific answers (e.g., "Japanese"). (ii) Feature selection is governed by competing circuits: we identify one pathway that prioritizes direct query-based reasoning and another that incorporates contextual cues. Their relative strength determines the final output. Attention knockout experiments show that ablating key heads involved in the context-based pathway can flip model predictions (e.g., flipping "Japanese" to "French"), further confirming this competitive

interaction. These findings support the class component of our hypothesis and illustrate how models select features to combine with the abstract class.

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

Crucially, our findings do not suggest that LLMs possess strong generalization abilities, nor do they contradict the theoretical basis of the stochastic parrot argument that LLMs are incapable of true language understanding from form-based training alone (Bender and Koller, 2020). Instead, to stretch the metaphor, we provide evidence of LLMs as *stochastic parakeets, peacocks, and penguins* which can leverage abstract class structures in ways that are neither purely memorized nor necessarily reliable. These abstractions result from next-token prediction during pre-training and extend beyond simple ontological hierarchies (e.g., superset-subset relationships), shaping the model's internal feature selection process.

In summary, our main contributions are:

- We introduce a novel setting that isolates how LLMs integrate irrelevant contexts, distinguishing generalization from memorization.
- We provide empirical evidence that LLMs exhibit class-based (mis)generalization, demonstrating sensitivity to abstract class structures beyond statistical co-occurrences.
- We uncover the internal computational mechanisms of class-based generalization, revealing competing circuits and hierarchical class representations.
- We propose a behavioral analysis framework that moves beyond accuracy-based evaluation, emphasizing the importance of understanding LLMs' internal mechanisms.

# 2 Related Work

**LLM Evaluation** Traditional NLP evaluation prioritizes test set performance but often overlooks how models arrive at their final answers. For LLMs trained on Internet-scale data, distinguishing genuine generalization from memorization or spurious correlations is challenging, especially with potential data leakage (Dziri et al., 2023; Wu et al., 2024c; Zhou et al., 2024; Balloccu et al., 2024; Xu et al., 2024). Prior work addresses this through data extraction (Carlini et al., 2021), statistical control (Min et al., 2022), adversarial perturbations (Mirzadeh et al., 2024), and error analysis (Dziri et al., 2023). In contrast, we take a behaviorfocused approach in a controlled setting, reducing the likelihood of pure pattern matching. Addi-



Figure 2: Class-based generalization framework: feature selection and combination.

tionally, by analyzing errors — often ignored by
accuracy-based metrics — we gain deeper insights
into model mechanisms.

Irrelevant Context Hallucinations Hallucina-168 tions in text generation have been studied in the 169 absence of context (McKenna et al., 2023; Kang 170 et al., 2024; Meng et al., 2022) and in cases with relevant context (Cao et al., 2020, 2022a; Maynez 172 et al., 2020; Lee et al., 2018; Adlakha et al., 2024; 173 Chuang et al., 2024; Petroni et al., 2020; Li et al., 174 2023). We focus on irrelevant context hallucina-175 tions, where extraneous context influences predic-176 tions. Unlike prior work on evaluating or miti-177 gating such errors (Cao et al., 2022b; Cuconasu 178 et al., 2024; Wu et al., 2024a; Petroni et al., 2020; 179 Li et al., 2023; Shi et al., 2023; Mirzadeh et al., 180 2024), we explain their underlying class-based 181 (mis)generalization mechanisms, conceptually and mechanistically. 183

Mechanistic Interpretability Mechanistic inter-184 pretability methods (Olah, 2022; Nanda, 2023) reverse-engineer LLMs via vocabulary projection (Belrose et al., 2023; Geva et al., 2022; nostalgebraist, 2020) and computational interventions (Ghandeharioun et al., 2024; Stolfo et al., 2023; Finlayson et al., 2021). Extending prior work (Merullo 190 et al., 2024; Wu et al., 2024b; Lv et al., 2024; Geva et al., 2022), we use these techniques to uncover how LLMs are influenced by irrelevant contexts. 193 By linking behavioral analysis with internal mech-194 anisms, we provide a mechanistic perspective on 195 irrelevant context hallucinations.

# **3** Framework and Hypothesis

197

198

200

In this section, we describe the abstract framework illustrated in Figure 2 and present our class-based (mis)generalization hypothesis. **Setting** Consider a query Q representing the question of interest and an irrelevant context C prepended to it. Let  $A_Q$  denote the model's answer to Q alone and  $A_{C+Q}$  denote the answer with the added context (contextual answer). We define **context features** and **query features** as sets of properties or attributes of the entities in C and Q, respectively. For example, for context features in Figure 2, the feature names are {*Type*, *Country*, *Year*, *etc.*}, with corresponding feature values {*Model, Japan, 1972, etc.*}. The **class** of  $A_{C+Q}$  derived from Q (e.g., "languages") determines which features the model should prioritize (e.g., derive "Japanese" based on the "Country: Japan" feature).

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

233

235

236

237

238

239

**Class-Based** (**Mis**)**Generalization Hypothesis** Given the setup C+Q, when the context influences the model predictions, instead of relying solely on the query, we hypothesize a structured mechanism by which the model integrates contextual information into its predictions. Specifically, we propose **class-based generalization**, where language models process context in two steps: they first *derive* an abstract class (e.g., "languages") and then *select* and *combine* relevant features from C or Q (e.g., "Japan" or "France") to generate an answer (e.g., "Japanese' or "French").

Let **query-based candidates** be answers derived from query features combined with the class (e.g., "French") and **context-based candidates** be answers derived from context features (e.g., "Japanese"). If  $A_{C+Q}$  is the query-based candidate, we define the case as **query-dominant**; otherwise, it is **context-dominant**. These terms emphasize the final outcome rather than the intermediate steps. See Table 1 for examples.

A special case arises when a token of the expected class is already present in the prompt (Appendix B), making the model more likely to copy it directly (Jiang et al., 2024).

C + Q	$A_Q$	$C_{\text{cand.}}$	$Q_{ ext{cand.}}$	$A_{C+Q}$	Case
C: Honda Civic, produced by Honda. Q: The original language of A Secret was	French	Japanese	French, English	Japanese	Context- dominant
C: City of Boroondara is in Melbourne. Q: Prime Minister of Malaysia is a legal term in	Malaysia	Australia	<b>Malaysia</b> , Malaysian	Malaysia	Query- dominant

Table 1: Examples of context- and query-dominant categorizations with context- and query-based candidates.

#### 241 242

243

244 245

246 247

2

2

23

25

253

254

20

257

25

20:

26

262

263

264 265

26

2

270 271

272

273

2



# 4 Dataset and Experimental Design

Models & Datasets We evaluate three pretrained LMs — Llama-3.1 (8B) (AI@Meta, 2024), Mistral v0.3 (7B) (Jiang et al., 2023), and Pythia (6.9B-deduped) (Biderman et al., 2023) — using their base versions to assess raw model behavior. We use the ParaRel dataset (Elazar et al., 2021), which consists of 39 factual QA subdatasets. Dataset statistics are provided in Table 12 in the Appendix A. Both Q and C are sourced from these datasets. Experiments are run on two RTX8000 GPUs.

**Experimental Setup** We compare two conditions: 1) **Q-only**, where each Q is formatted using a predefined template and a subject-relation-object (s, r, o) triplet from ParaRel, resulting in 27.6K queries. The model's top-1 prediction is  $A_Q$ . 2) **C+Q**, where each Q is prepended with context demonstrations from other subdatasets spanning various relation types, introducing controlled contextual variation. We randomly sample 100 examples per subdatasets<sup>2</sup>, generating 3,900 context variations per query, totaling 106M examples. The model's top-1 prediction is  $A_{C+Q}$ .

**Context- and Query-Based Candidates** To make the definitions from Sec. 3 precise, we define a **context-based candidate**  $x \in C_{\text{cand.}}$  to be a candidate among the top three<sup>3</sup> predictions under C + Q but not among the top ten<sup>4</sup> predictions under Q. A **query-based candidate**  $x \in Q_{\text{cand.}}$  appears in the top predictions under both conditions. Note that  $A_{C+Q}^{\text{top3}} = C_{\text{cand.}} \cup Q_{\text{cand.}}$ . See Table 1 for examples.

 $A_{C+Q}^{\text{top3}} := \{ x \mid x \in \text{top 3 candidates under } C + Q \}$  (1)

$$A_Q^{\text{top10}} := \{ x \mid x \in \text{top 10 candidates under } Q \}$$
(2)

$$C_{\text{cand.}} := \{ x \mid x \in A_{C+Q}^{\text{top3}} \text{ and } x \notin A_Q^{\text{top10}} \}$$
(3)

$$Q_{\text{cand.}} := \{ x \mid x \in A_{C+Q}^{\text{top3}} \text{ and } x \in A_Q^{\text{top10}} \}$$

$$\tag{4}$$

# 5 Behavioral Analysis of Contextual Answers

276

277

278

279

281

282

283

286

291

292

294

295

296

297

299

300

301

302

303

We now investigate how irrelevant context influences model predictions, verifying our class-based (mis)generalization hypothesis through textuallevel behavioral analysis. Specifically, we examine: (1) whether irrelevant context causes behavioral changes (Sec. 5.1), (2) whether the influence of irrelevant context aligns with our hypothesis (Sec. 5.2). (3) whether the observed correlation between irrelevant context and context-based candidates is statistically significant (Sec. 5.3).

Case	Top-3 Candidates	Llama	Mistral	Pythia
No in- fluence	1. All query-based $(C_{\text{cand.}} = \emptyset)$	47.9%	48.0%	39.3%
Query- dominant	2. Mix, top-1 is query- based	27.9%	25.7%	27.2%
Context-	3. Mix, top-1 is context- based	15.1%	17.0%	19.2%
dominant <sup>-</sup>	4. All context-based $(Q_{\text{cand.}} = \emptyset)$	10.1%	10.3%	14.3%

Table 2: Breakdown of samples according to the composition of  $A_{C+Q}^{\text{top-3}}$ , based on 106M datapoints. Detailed results can be found in Table 7 in the Appendix.

# 5.1 Behavioral Changes Induced by Irrelevant Context

In this section, we investigate whether adding irrelevant context leads to behavioral changes in model predictions. From an accuracy perspective, we observe a slight decrease: across 39 subdatasets, the accuracy for Llama-3 drops from 47.2% to 43.1%, and for Mistral, from 38.2% to 35.3% (Table 6 in Appendix). While these changes are modest, accuracy alone does not provide a complete picture of changes in model predictions. To address this gap, we measure the answer change rate ( $\Delta$  Rate) after adding the irrelevant context:  $\Delta$ Rate =  $\frac{|A_{C+Q}\neq A_Q|}{\#$  datapoints}. For Llama-3, 38.3% of responses changed after adding irrelevant context, while for Mistral, nearly half of the datapoints

<sup>&</sup>lt;sup>2</sup>P264 has only 53 examples, so we include all of them.

<sup>&</sup>lt;sup>3</sup>A threshold of three ensures that classified context-based candidates are strongly influenced by context.

 $<sup>^{4}</sup>$ A threshold of 10 ensures that classified context-based candidates are not plausible answers under Q alone.

Query Type	Ctx. Type	Context Demonstration + Query and Answer
Language	Person/ Music	<b>Prompt:</b> Amilcare Ponchielli plays opera. The original language of A Hunting Accident was Answer: $A_{C+Q} =$ Italian, $A_Q =$ English
	Make/ Model	<b>Prompt:</b> Toyota Alphard, produced by Toyota. The original language of A Hunting Accident was Answer: $A_{C+Q} =$ Japanese, $A_Q =$ English
Place	Person/ Religion	<b>Prompt:</b> Indo-Greek Kingdom is follower of Buddhism. Alpha Island is a part of the continent of <b>Answer:</b> $A_{C+Q} = Asia$ , $A_Q = Alpha$
Flace	Place	<b>Prompt:</b> Council of States of Switzerland is a legal term in Switzerland. Alpha Island is a part of the continent of Answer: $A_{C+Q}$ = Europe, $A_Q$ = Alpha

Table 3: Examples of context-based candidates across different query and context types.

(48.0%) experience a shift in predictions (Table 6 in Appendix).

305

307

308

310

311

313

314

315

316

317

320

321

325

326

328

332

334

336

We further examine the cases under the C + Q condition based on the composition of  $(A_{C+Q}^{\text{top-3}})$ (Table 2). Roughly 48%<sup>5</sup> of samples are unaffected by the irrelevant context for Llama and Mistral (case 1), meaning all top-3 candidates are querybased). However, when predictions are influenced by the added context (cases 2, 3 and 4), about half of these instances (49.5% for Llama, 52.5% for Mistral) become context-dominant. These results demonstrate the influence of irrelevant contexts, even if the overall accuracy is little changed.

# 5.2 Human Annotation of Context-Based Candidates

Next, we examine whether these behavioral changes pattern with our class-based generalization hypothesis. To do so, we annotate context-based candidates, which capture the shifts induced by irrelevant context. We assess whether each answer explicitly integrates *identifiable features* from the context *and* combines them with the *expected class* indicated by the query. Annotation procedure and examples are provided in Appendix D.

We perform this annotation on a randomly sampled set of 500 context-based candidates across different subdatasets. Our results reveal that 81.6% of the responses incorporate features from the provided context, 84.4% belong to the correct class, and 71.0% satisfy both criteria – combining identifiable context features with the correct abstract class. This finding provides strong evidence for our hypothesis as a majority of these samples can be explained by the hypothesis. Table 3 provides illustrative examples of the model's output adapting to contextual cues.

338

339

340

341

342

343

344

345

346

347

349

351

352

353

354

355

356

357

358

359

361

362

363

364

365

366

367

368

369

371

# 5.3 Statistical Validation of Contextual Influence

Next, we investigate whether the correlation between irrelevant context and context-based candidates is statistically significant. To quantify the dependence between a context C (e.g., *Honda*) and its associated context-based candidate  $C_{cand.}$  (e.g., Japanese), we compute the pointwise mutual information (PMI) between them. Specifically, we sample 100 distinct contexts from various subdatasets. Each context is paired with 100 different queries belonging to the same expected class (e.g., languages, places, etc.), resulting in 10,000 instances per class. Since context-based candidates are determined independently of the queries, each context  $C_i$  is paired with its corresponding candidate  $C_{\text{cand.,i}}$ , regardless of the 100 queries. This yields 100 pairs of  $(C_i, C_{\text{cand. }i})$  per class, such as (Honda, Japanese) for languages, and (Honda, Japan) for places. The mean PMI across the 100 pairs of each class is computed as:

$$\mu_{\text{observed}} = \frac{1}{100} \sum_{i=1}^{100} \text{PMI}(C_i, C_{\text{cand.},i}) \quad (5)$$

$$PMI(C_i, C_{cand,i}) = \log \frac{P(C_i, C_{cand,i})}{P(C_i)P(C_{cand,i})}.$$
 (6)

In this formula,  $P(C_i) = 1/100$ , since we have 100 distinct contexts.  $P(C_{\text{cand.},i})$  is estimated based on its frequency among all 10,000 generated answers  $A_{C+Q}$  for the given expected class. Similarly,  $P(C_i, C_{\text{cand.},i})$  is computed from its cooccurrence within these samples. Across all models and expected classes, the mean PMI is approximately 4, suggesting a strong association between contexts and their corresponding candidates.

<sup>&</sup>lt;sup>5</sup>Notably, due to the conservative choice of 10 for  $A_Q^{\text{top-10}}$ , some answers in case 1 might also be context-based but already appear in the top-10 predictions under the Q condition. Therefore, we exclude these cases from further analyses.



Figure 3: Logit attribution (C+Q condition) along residual stream  $(R_{T,l}^1, R_{T,l}^2)$  reveals the construction of abstract class representation in the lower layers, with competition between  $Q_{\text{cand.}}$  (dashed) and  $C_{\text{cand.}}$  (dotted) in the mid to higher layers. The example token in parenthesis correspond to Table 1. Additional results are in Appendix F.

L16	L17	L18	L19	L20	L21	L22	L23	L24	L25	L27	L28	L29	L30	L31
languages	languages	/is	languages	languages	languages	languages	languages	English						
/is	/is	languages	/is	English	English	English	English	English	English	English	English	English	English	Japanese

Table 4: Logit lens on Llama-3 showing top-1 predictions shifting from abstract concepts (e.g., 'languages') to concrete instances (e.g., 'English' or 'Japanese') across layers. The first and second row correspond to  $R_{T,l}^1$ , and the second row is  $R_{T,l}^2$ , respectively. See Appendix F.2 for the corresponding prompt and associated probabilities.

To formally assess statistical dependence, we perform a one-sample t-test against the null hypothesis  $\mathbb{E}[\text{PMI}(C_i, C_{\text{cand.},i})] = 0)$  (which would indicate independence). With a *p*-value of 0.001, we reject the null hypothesis, concluding that *C* and *C*<sub>cand.</sub> exhibit significant dependence. (See Table 8 in the Appendix for full results.)

# 6 Mechanistic Analysis of Contextual Answers

We next investigate whether the models' internal computations reflect the class-based generalization that we observed above. In Sec. 6.1, we use logit attribution to show that models construct **abstract class representations**, supporting the class component of our hypothesis. In Sec. 6.2 and Sec. 6.3, we apply activation patching and attention knockout to reveal that feature selection in our hypothesis arise from **competition between circuits**, where distinct query-based pathways (computing  $Q_{cand.}$ ) and context-based pathways (computing  $C_{cand.}$ ) compete to determine the final answer. These findings provide mechanistic evidence for our hypothesis.

**Data** We randomly draw 1,000 context-dominant and 1,000 query-dominant datapoints from case 2 and case 3 in Table 2 as these cases have both query- and context-based candidates.

# 6.1 Logit Attribution

373

375

387

397

400

**Method** To explore how models build answers across layers, we apply logit attribution (nostalge-

braist, 2020) to trace predictions across layers by projecting hidden states onto the vocabulary space. Given a prompt with T tokens and a model with L layers, we extract hidden states at the last token position  $h_{T,j} \in \mathbf{R}^d$ , where  $j \in \{1, ..., L\}$  and d is the hidden size. These are projected onto the vocabulary space using Unembed(LayerNorm $(h_{T,j})) \in$  $\mathbf{R}^{|V|}$ , where the Unembed matrix corresponds to the transpose of the input embedding weights. Models maintain a residual stream for each token i, which accumulates information as it passes through each layer. At each layer, two key transformations occur: attention update  $(A_{i,l})$  and MLP update  $(M_{i,l})$ . Mathematically, the updates follow:

$$A_{i,l} = \operatorname{ATTN}(R_{i,l}^0) \tag{7}$$

$$R_{i,l}^1 = A_{i,l} + R_{i,l}^0 \tag{8}$$

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

$$M_{i,l} = \mathsf{MLP}(R_{i,l}^1) \tag{9}$$

$$R_{i,l}^2 = M_{i,l} + R_{i,l}^1 \tag{10}$$

where  $R_{i,l}^1$  is the residual stream after attention at layer l, and  $R_{i,l}^1$  is the final residual stream at layer l after the MLP update. (See Appendix F.)

**Findings** To understand how different tokens evolve across layers, we project the last token residual stream  $R_{T,l}^1$  (after attention) and  $R_{T,l}^2$  (after the MLP) onto the vocabulary space at each layer. Figure 3 tracks the logits for  $C_{\text{cand.}}$ ,  $Q_{\text{cand.}}$ , and class tokens under the C + Q condition. Additional results are provided in Appendix F. Figure 3 reveals a hierarchical class-to-instance process in answer generation. Early layers prioritize class token logits (solid) like "languages", suggesting that the model first constructs abstract class representations. Around the middle layers, candidate answer logits (dashed/dotted) begin to rise, refining these abstract representations into concrete answers. In Table 4, a concrete example of logit lens top-1 predictions reveals how Llama-3 shifts from abstract class to concrete instances. This pattern supports our hypothesis that **models leverage class-based information in shaping their predictions**.

Moreover, the figures highlight a competition between  $C_{\text{cand.}}$  (dashed) and  $Q_{\text{cand.}}$  (dotted), particularly in context-dominant cases (pink). In early layers, logits for  $C_{\text{cand.}}$  and  $Q_{\text{cand.}}$  form two distinct groups, regardless of dominance. Around layer 14,  $Q_{\text{cand.}}$  (dotted) in both cases begin to split, followed by  $C_{\text{cand.}}$  (dashed) in layer 17. By layer 24,  $C_{\text{cand.}}$  (dark pink) surpass  $Q_{\text{cand.}}$  (light pink) logits in context-dominant settings, marking a decisive shift in the competition. After this, the early two-group pattern reemerges but with reversed dominance --- context-based candidates prevail in context-dominant cases, and query-based candidates in query-dominant cases. By layer 29, the final prediction is fully formed, with the top logits corresponding to the final output. These observations reveal key insights: (i) existence of competition: even when the final prediction is query-dominant, context-based candidates remain actively computed across layers. (ii) critical transition (Layers 17-24): the decisive competition between query- and context-based candidates occurs primarily in this range, determining which candidate is promoted.

#### 6.2 Activation Patching

Method To understand the competition between  $C_{\text{cand.}}$  and  $Q_{\text{cand.}}$ , we investigate whether distinct context and query circuits exist within the model's internal activations. We apply activation patching (Ghandeharioun et al., 2024; Meng et al., 2022), a technique for causal intervention that selectively perturbs and restores activations to assess their contribution. We conduct three model runs: (1) Clean run: Standard forward pass with the original prompt, recording activations  $\bigcup h_{i,l}^0$ . (2) Corrupted run: Forward pass with Gaussian noise  $\epsilon \sim \mathcal{N}(0, \sigma^6)$ ) injected into context or query topic token embeddings, yielding perturbed activations

 $\bigcup h_{i,l}^1$ , and the final log-probabilities of candidates  $\log p(t|\lfloor \rfloor h_{i,l}^1)$ . (3) Restoration run: Same as the corrupted run, but iterating over all i and l, restoring each  $h_{i^*,l^*}^0$ , while keeping the rest cor-rupted. By injecting noise at context subject and object (context patching) or query subject posi-tion (query patching) and measuring the recovery of predictions, we differentiate context and query circuits, tracing how features from these tokens propagate through the model and how they con-tribute to context-based or query-based candidates. 

The **restoration effect** for each  $i^*$  and  $l^*$  is calculated as in Eq. 11, where  $t \in \{C_{\text{cand.}}, Q_{\text{cand.}}\}$ , with higher values indicate stronger contributions.

$$\text{RE}(i^*, l^*, t) = \log p(t | h^0_{i^*, l^*} \cup h^1_{-i^*, -l^*})$$

$$-\log p(t|\bigcup h_{i,l}^1) \tag{11}$$

$$\Delta_{\text{query}}(i^*, l^*) = \text{RE}(i^*, l^*, Q_{\text{cand.}}) - \frac{12}{\text{RE}(i^*, l^*, C_{\text{cand.}})}$$

$$\Delta_{\text{context}}(i^*, l^*) = \operatorname{RE}(i^*, l^*, C_{\text{cand.}}) - \operatorname{RE}(i^*, l^*, Q_{\text{cand.}})$$
(13)

In context and query circuits, we compute Eq. 13 (Figures 4a, 4c) and Eq. 12 (Figures 4b, 4d), respectively. Comparing restoration effects maps circuits responsible for context- and query-based candidates and identifies where their competition occurs. (See Appendix G for implementation details and additional results.)

**Findings** The results reveal distinct circuits for context- and query-based pathways. Figures 4a and 4c show the same context circuit aggregating information from the context subject and object in both cases, transferring it to the final token position from layer 17 onward. In contrast, Figures 4b and 4d indicate that the same query circuit for both cases integrating query subject information earlier than in the context circuit, from layer 8. The log-probability increases after layer 24 in context circuit and after layer 16 in query circuit.

Both circuits exist across context- and querydominant cases, but their relative strength determines the final prediction. In context-dominant cases, the context circuit wins, with a larger logprobability difference (max 2.28) compared to the query circuit (max 1.10). Conversely, in querydominant cases, the query circuit exerts a stronger influence (max 1.37 vs. 1.25). Notably, between layers 17 and 24, the query-dominant case shows minimal context information transfer (Figure 4c),

<sup>&</sup>lt;sup>6</sup>We follow Meng et al. (2022) to set  $\sigma = 0.3$  as three times of the empirical standard deviation of the input embeddings.



Figure 4: Left-hand plots demonstrate the context circuit, which extracts features from context and computes context-based candidates, while right-hand plots illustrate the query circuit. These circuits are the same in both context- and query-dominant cases; the difference lies in their strength, revealing the competition between context- and query-based candidates. An example is C\_REL0 = [BOS], C\_SUBJ1=' Honda Civic', C\_REL2=', produced by',C\_OBJ3=' Honda', Q\_REL4='. The original language of ',Q\_SUBJ5=' A Secret', Q\_REL6=' was'.

aligning with slower logit attribution growth (Figure 3). This confirms that both pathways exist for both cases with final predictions depending on their relative activation strength, and layers 17 to 24 are the key to promoting context-based candidates.

527

528

529

530

531

532

533

535

537

541

542

543

552

# 6.3 Flipping Model Predictions via Attention Knockout

To examine the causal role of internal competition in shaping the final output  $A_{C+Q}$ , we intervene in two key layers of the context circuit: layer 17 (where context first transfers to the last token) and layer 24 (where it is most integrated). By restricting attention to the query in the context-dominant case and to the context in the query-dominant case, we test whether predictions can be flipped (e.g., "Japanese" to "French"). See Appendix H for details and additional results. Table 5 (Llama) shows that in the context-dominant case, blocking context flow causes  $Q_{\text{cand.}}$  probabilities to surpass  $C_{\text{cand.}}$ on average, flipping 465/1000 datapoints to querybased candidates. In the query-dominant case, intervention increases Ccand. probability by 4.7 and decreases  $Q_{\text{cand.}}$  probability by 8.4, flipping 225/1000 datapoints. These results confirm the competition between  $C_{\text{cand.}}$  and  $Q_{\text{cand.}}$ , and that these two layers are the key to promoting context-based candidates.

553SummaryThese findings support the class-based554(mis)generalization hypothesis. Logit attribution555confirms that models first construct abstract class556representations before refining them into specific557answers. Activation patching reveals competing558circuits for feature selection: one favoring direct559query-based pathway and the other integrating con-

	Orig.	L17+L24		2 Ra	and.					
	Prob.	Prob.	$\Delta$	Prob.	$\Delta$					
Context-Dominant										
$C_{ ext{cand.}} Q_{ ext{cand.}}$	25.5 8.6	13.1 14.8	-12.4 + 6.2	21.0 10.8	-4.5 +2.2					
Query-Dominant										
$Q_{\text{cand.}} \\ C_{\text{cand.}}$	35.2 6.6	26.8 11.3	-8.4 +4.7	29.6 7.4	-5.7 +0.8					

Table 5: Effect of attention knockout on context- ( $C_{cand.}$ ) and query-based ( $Q_{cand.}$ ) candidate probabilities on Llama-3. "Orig." = No intervention, "2 Rand." = Average of interventions on two random layers over three runs.  $\Delta$  denotes the change from the original setting.

textual cues, with their strength shaping the final output. Notably, context circuit strengthens between layers 17 and 24, validated by the flipped predictions from attention knockout. 560

561

562

563

564

565

566

567

569

570

571

572

573

574

575

576

577

# 7 Conclusion

By analyzing the mechanism behind irrelevant context hallucinations, our study demonstrates that LLMs exhibit class-based (mis)generalization, relying on abstract class structures in a systematic yet flawed manner. Through mechanistic analysis, we show that this phenomenon arises from hierarchical class-to-instance predictions and competing circuits that mediate feature selection. These findings challenge a potential misconstrual of the stochastic parrot hypothesis that LLMs can only regurgitate surface-level patterns. Rather, we argue that they utilize class structures in ways that are neither purely memorized nor necessarily reliable.

# 8 Limitations

578

Our work has several limitations. First, our experiments are conducted in a controlled setting, 580 which helps isolate generalization from memoriza-581 tion and enables analysis at both behavioral and mechanistic levels. However, future work could improve upon this by designing setups that disentangle memorization and generalization in natu-585 rally occurring text. Second, our study is limited 586 to English-language datasets, and we only evaluate 587 models of certain sizes (around 7-8B) due to computational constraints. It remains an open question whether class-based generalization is influenced by model size. Do larger models exhibit stronger 591 correlations of this kind? Do smaller models also display class-based generalization, and if so, what 593 is the minimum size required? Third, in the mech-594 anistic interpretability section, we focus primarily on layer-wise analysis to support our main hypothesis, while attention head analysis is left for future work. Finally, while we conduct interventions, our primary goal is not to mitigate contextual halluci-599 nations. Developing mitigation methods informed by our findings and evaluating their effectiveness is an important direction for future research.

# References

606

607

610

611

612

613

614

615

616

617

620

621

623

625

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2024. Evaluating correctness and faithfulness of instructionfollowing models for question answering. *Transactions of the Association for Computational Linguistics*, 12:681–699.
- AI@Meta. 2024. Llama 3 model card.
- Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondrej Dusek. 2024. Leak, cheat, repeat: Data contamination and evaluation malpractices in closedsource LLMs. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93, St. Julian's, Malta. Association for Computational Linguistics.
- Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 610–623.

Emily M Bender and Alexander Koller. 2020. Climbing towards NLU: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198. 630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.
- Meng Cao, Yue Dong, and Jackie Cheung. 2022a. Hallucinated but factual! inspecting the factuality of hallucinations in abstractive summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3340–3354, Dublin, Ireland. Association for Computational Linguistics.
- Meng Cao, Yue Dong, Jingyi He, and Jackie Chi Kit Cheung. 2022b. Learning with rejection for abstractive text summarization. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9768–9780, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Meng Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit Cheung. 2020. Factual error correction for abstractive summarization models. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6251–6258, Online. Association for Computational Linguistics.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Yung-Sung Chuang, Linlu Qiu, Cheng-Yu Hsieh, Ranjay Krishna, Yoon Kim, and James R. Glass. 2024. Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1419–1436, Miami, Florida, USA. Association for Computational Linguistics.
- Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonellotto, and Fabrizio Silvestri. 2024. The power of noise: Redefining retrieval for rag systems. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, volume 17 of *SIGIR* 2024, page 719–729. ACM.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter

797

798

799

West, Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, Soumya Sanyal, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. 2023. Faith and fate: Limits of transformers on compositionality. In *Thirty-seventh Conference on Neural Information Processing Systems*.

687

693

697

699

700

701

703

704

707

709

710

711

712

713

714

715

716

718

719

721

722

724

725

726

727

728

729 730

731

733 734

737

738

739

740

741

742

743

- Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *Transactions of the Association for Computational Linguistics*, 9:1012–1031.
- Matthew Finlayson, Aaron Mueller, Sebastian Gehrmann, Stuart Shieber, Tal Linzen, and Yonatan Belinkov. 2021. Causal analysis of syntactic agreement mechanisms in neural language models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1828–1843, Online. Association for Computational Linguistics.
- Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2022. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 30–45, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. 2024. Patchscopes: A unifying framework for inspecting hidden representations of language models. In *Forty-first International Conference on Machine Learning*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Yibo Jiang, Goutham Rajendran, Pradeep Ravikumar, and Bryon Aragam. 2024. Do llms dream of elephants (when told not to)? latent concept association and associative memory in transformers. *arXiv preprint arXiv:2406.18400*.
- Katie Kang, Eric Wallace, Claire Tomlin, Aviral Kumar, and Sergey Levine. 2024. Unfamiliar finetuning examples control how language models hallucinate. *arXiv preprint arXiv:2403.05612*.
- Katherine Lee, Orhan Firat, Ashish Agarwal, Clara Fannjiang, and David Sussillo. 2018. Hallucinations in neural machine translation.
- Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, and Sanjiv Kumar. 2023. Large language models with controllable working memory. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1774–1793, Toronto, Canada. Association for Computational Linguistics.

- Ang Lv, Yuhan Chen, Kaiyi Zhang, Yulong Wang, Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan. 2024. Interpreting key mechanisms of factual recall in transformer-based language models. *arXiv preprint arXiv:2403.19521*.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1906–1919, Online. Association for Computational Linguistics.
- Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. 2023. Sources of hallucination by large language models on inference tasks. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In *Advances in Neural Information Processing Systems*.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2024. Language models implement simple Word2Vec-style vector arithmetic. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5030–5047, Mexico City, Mexico. Association for Computational Linguistics.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. *arXiv preprint arXiv:2410.05229*.
- Neel Nanda. 2023. Mechanistic interpretability quickstart guide.
- nostalgebraist. 2020. interpreting gpt: the logit lens.
- Chris Olah. 2022. Mechanistic interpretability, variables, and the importance of interpretable bases.
- Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. How context affects language models' factual predictions. In *Automated Knowledge Base Construction*.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can

- be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. 2023. A mechanistic interpretation of arithmetic reasoning in language models using causal mediation analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7035–7052, Singapore. Association for Computational Linguistics.

805

806

811

812 813

814

815 816

817

818

819

820

824 825

827

829

830

832

834

836

842

846

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
  - Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai Zhang, and Yanghua Xiao. 2024a. How easily do irrelevant inputs skew the responses of large language models? In *First Conference on Language Modeling*.
    - Wenhao Wu, Yizhong Wang, Guangxuan Xiao, Hao Peng, and Yao Fu. 2024b. Retrieval head mechanistically explains long-context factuality. *arXiv preprint arXiv:2404.15574*.
  - Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. 2024c. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 1819–1862, Mexico City, Mexico. Association for Computational Linguistics.
  - Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. 2024. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*.
  - Lei Yu, Meng Cao, Jackie CK Cheung, and Yue Dong.
     2024. Mechanistic understanding and mitigation of language model non-factual hallucinations. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7943–7956, Miami, Florida, USA. Association for Computational Linguistics.
- Yuhang Zhou, Paiheng Xu, Xiaoyu Liu, Bang An, Wei Ai, and Furong Huang. 2024. Explore spurious correlations at the concept level in language models for text classification. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 478–492, Bangkok, Thailand. Association for Computational Linguistics.

# A Dataset

852

853

855

857

870

871

872

877

882

888

896

The detailed breakdown of ParaRel dataset (Elazar et al., 2021) based on relation type is presented in Table 12. We categorize the sub-datasets into 5 knowledge types based on the expected class or type of the answer (column 'Ctx Type'): 'language', 'place', 'company', 'job', and if a subdataset doesn't fit into the above types then it is categorized as 'others'. This is the dataset that we use for Q-only experiments, and we construct the dataset for C + Q experiments by generating 3900 context variations spanning all knowledge types per query, resulting in a dataset of 106.2M data points. For each generation, we restrict the vocabulary to the set of tokens that begins with a capitalized English letter (Yu et al., 2024). When evaluating, we lowercase generated and gold answers and perform string matching: if the top-1 generated answer is a substring of the gold answer, then this is correct.

# **B** Class-based Generalization

We further categorize class-based generalization into two distinct cases:

• **Copying:** When a token belonging to the expected class appears in the context, the model is more likely to directly copy it as the answer. From a dataset statistics perspective, we observe a high copy rate when the context contains tokens belonging to the same class as the query.

**Example:** The mother tongue of Dominique Sanda is French. The original language of Puss in Boots was  $\rightarrow$  French.

• Non-copying: When tokens of the expected query class are not explicitly present in the input, the model combines the expected class with relevant features inferred from context or query to generate an answer.

**Example:** Honda Civic (fifth generation), produced by Honda. The original language of Tow Truck Pluck was  $\rightarrow$  Japanese.

# C Behavioral Changes Induced by Irrelevant Context

# C.1 Irrelevant Context Hallucination Evaluation

In Table 6, we provide detailed statistics of the accuracy/wrong rate for each model under each case for all three models.

Model	Q-0	ONLY	C+Q		
	CASE	Prop.	CASE	Prop.	$\Delta$ Rate
Llama	Т	47.2%	$\begin{array}{c} T \rightarrow T \\ T \rightarrow F \end{array}$	35.7% 11.5%	0% 100%
Liumu	F	52.8%	$\begin{array}{c} F \rightarrow T \\ F \rightarrow F \end{array}$	7.4% 45.4%	100% 42.7%
	Total	47.2%	Total	43.1%	38.3%
Mistral	Т	38.2%	$\begin{array}{c} T \rightarrow T \\ T \rightarrow F \end{array}$	29.4% 8.8%	0% 100%
	F	61.8%	$\begin{array}{c} F \rightarrow T \\ F \rightarrow F \end{array}$	5.9% 55.9%	100% 59.5%
	Total	38.2%	Total	35.3%	48.0%
Pythia	Т	30.9%	$\begin{array}{c} T \rightarrow T \\ T \rightarrow F \end{array}$	22.4% 8.4%	0% 100%
	F	69.1%	$\begin{array}{c} F \rightarrow T \\ F \rightarrow F \end{array}$	5.6% 63.6%	100% 67.8%
	Total	30.9%	Total	28.0%	57.1%

Table 6: Comparison of proportions (Prop.) of correct and incorrect answers in Q-only and C+Q cases, along with answer change rates ( $\Delta$  Rate) for different models. Average across 39 datasets are reported. In the 'Total' row, under 'Prop.' column, it indicates the global accuracy across different cases, while under ' $\Delta$  Rate' column, it underlies the global answer change rate.

Table 6 shows that models are not robust against irrelevant context. Even when a single irrelevant demonstration is prepended, models exhibit notable shifts in performance. For instance, in Llama, 11.5% of previously correct answers become incorrect, while 7.4% of incorrect answers are corrected after adding context. However, accuracy alone does not capture all behavioral shifts — predictions can still change even if they remain incorrect. 897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

# **C.2** Composition of $A_{C+Q}^{\text{top-3}}$

Table 7 provides counts and proportions of the breakdown of samples according to the composition of  $A_{C+Q}^{\text{top-3}}$  for three models based on 106M datapoints.

# **D** Annotation

#### **D.1** Annotation Procedure

To systematically evaluate the impact of irrelevant context on model predictions, we perform an annotation procedure for context-based candidates — those predictions that were influenced by the inclusion of extraneous context. The aim was to rigorously assess whether (i) these predictions incorporated identifiable features from the context, **and** (ii) appropriately combined them with the ex-

Case	Top-3 Candidates	Llama	Mistral	Pythia
No influence	1. All query-based	50,874,341 (47.9 %)	51,013,564 (48.0%)	41,833,760 (39.3%)
Query- dominant	2. Mix: Query + Context, top-1 is query-based	27,940,495 (27.9 %)	27,342,287 (25.7%)	28,885,252 (27.2%)
Context-	3. Mix: Query + Context, top-1 is context-based	16,069,253 (15.1 %)	17,013,397 (16.0%)	20,412,292 (19.2%)
dominant⁻	4. All context-based	11,353,892 (10.1 %)	10,963,675 (10.3%)	15,250,026 (14.3%)

Table 7: Breakdown of samples according to the composition of  $A_{C+Q}^{\text{top-3}}$ , based on 106M datapoints.

pected class as indicated by the query. Upon acceptance, we will release the annotation.

Step 1: Candidate Selection We first randomly sample a set of 500 context-based candidates from different sub-datasets, ensuring a diverse set of instances. Context-based candidates were selected for both context- and query-dominant cases.

Step 2: Context Feature Identification For each context-based candidate, we analyzed the context —specifically the subject and object — to identify any features that could have been leveraged by the model in generating the response. ('contextinfluenced?' row in Table 13).

Each feature is categorized as identifiable if it can be explicitly extracted from the context. For example, the country of origin of a figure (e.g., candidates 'South' 'Korea' for context subject 'Lee Jonghyun' in Example 5 in Table 13), country/continent of a district ('India', 'Asia' for context object 'Bihar' in Example 4 in Table 13) are classified as identifiable. In contrast, context-based candidates 'Bee', 'Beach' are categorized as non-context influenced for context subject 'Grant Green' and object 'jazz' as shown in Example 6 in Table 13.

We ensure transparency by documenting the rationale. For example, in Example 2 of Table 13, we provide the justification that 'Svend Asmussen' is a Danish violinist and jazz musician, which supports that 'Danmark' is a context-influenced candidate.

Step 3: Class Verification Next, each contextbased candidate is classified according to the abstract class suggested by the query. The candidate is compared to the expected class, and we verify whether the response falls within the correct category. For example, the context-based candidates 'Vietnamese' and 'Thai' for Example 1 in Table 13 have the correct class 'language', but 'South', 'Korea' in Example 5 in Table 13 do not have the correct class because the query is asking about continent, not country.

Step 4: Hypothesis Verification Finally, a context-based candidate is considered to satisfy the hypothesis if it meets the criteria from both Step 2 (context feature identification) and Step 3 (class verification). Only candidates that successfully integrate context features and align with the expected class are retained as valid instances.

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

#### **D.2** Annotation Examples

Details of examples and non-examples are shown in Table 13.

#### Statistical Validation of Contextual E Influence

Mean PMI values for each model are presented in Table 8. A mean PMI of around 70 across all models and expected classes confirms strong statistical dependence (Table 8).

Full results on three models in Table 8.

Value	Llama-3	Mistral	Pythia
Mean PMI	3.9	3.7	3.8
T-statistic	8.1	7.3	6.6
<i>p</i> -value	0.0006	0.0009	0.001

Table 8: Mean PMI values and T-test results for all three models.

#### **Logit Attribution** F

#### Implementation Details **F.1**

When the target candidates or class have multiple tokens, we take the maximum logit, and average this maximum logit across all data points in the dataset.

To obtain the class logits from the model, we predefine a list of tokens according to the relation type.

921

922

923

925

926

927

928

929

930

931

932

933

934

936

937

938

943 944

945

947

949

951

953

955

956

957

• Languages: languages, language, tongue, tongues, lingua, dialect, dialects

989

991

993

994

995

997

1000

1001

1002

1003

1004

1005

1006

1007

1008

1010

1011

1012

1013

1014

1016 1017

1018

1019

1020

1022

1023

1025

1026

1027

1028

1029

1030

1031

1033

- Places: country, countries, place, places, location, locations, territory, city, cities, town, towns, village, villages, state, states, province, provinces, district, districts, continent, continents
  - Companies: company, companies, manufacturer, manufacturers, make, firm, firms, business, corporation, corporations, enterprise, enterprises, organization, organizations, channel, channels, broadcaster, broadcasters, industry, industries
    - Jobs: position, positions, job, jobs, career, careers, profession, professions, occupation, occupations, role, roles, assignment, assignments, employment, employments
    - Others: expertise, area, areas, field, fields, subject, subjects, instrument, instruments, genre, music, religion, religions, concept, concepts, framework, frameworks, artifact, artifacts, type, types, part, parts, class, classes, eponym, eponyms, entity, entities, person, persons, place, places

# F.2 Logit Lens Example

We provide an example of how the model's top-1 predictions shift along the residual stream from abstract concepts to concrete instances across layers in Table 4 and Figure 5. The prompt used is *Honda Civic (fifth generation), produced by Honda. The original language of Tow Truck Pluck was.* Red indicates probability around 80%. We show predictions above layer 15 because lower than this, the predictions are not interpretable.

# F.3 Additional Logit Attribution Results

Additional results for Llama 8B are presented in Figure 6. Importantly, we point out that the classbased generalization might have existed already for the Q-only case. In Figure 6a, we observe a similar pattern as the C+Q case presented in Figure 3a – models build abstract class representation in the lower layers, before refining their answers to concrete ones. In fact, when we plot the logit difference of the abstract class tokens under C+Q and Qonly case in Figure 6b, as shown as orange and yellow lines for context-dominant and query-dominant case, the lines center around 0 – suggesting that the computation of abstract class representations exists for zero-shot case, and is not influenced by the added irrelevant context.

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1048

1049

1050

1051

1052

1053

1055

1056

1057

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1073

1074

1075

1076

1077

1078

1079

1080

1081

Logit attribution results for Mistral 7B are presented in Figure 7, and for Pythia 6.9B in Figure 8. We remark that these plots follow a similar pattern.

# **G** Activation Patching

#### G.1 Implementation Details

In the corrupted run, we corrupt the embeddings of all tokens for context subject and object in context patching, and all tokens for query subject in query patching by adding a Gaussian noise where  $\sigma$  is 3 times of the empirical standard deviation of the input embeddings over a body of text (*sigma*  $\approx$ 0.3) (Meng et al., 2022).

# G.2 Additional Activation Patching Results

Activation patching results under the C+Q condition for Mistral and Pythia are in Figure 10 and 13, respectively.

Additionally, we also visualize the query circuit under the Q-only condition in Figure 9, 11, and 12, for Llama, Mistral, and Pythia, respectively. We remark on two important observations: (i) The query circuit is the same for context-dominant and query-dominant data, without irrelevant context. (ii) The query circuit remains as is after adding the irrelevant context, as compared to Figures 4b and 4d.

# H Attention Knockout

# H.1 Implementation Details

In the attention knockout experiments, our goal is to see if we can intervene in the internal computation to change the output behavior. Specifically, in context-dominant case, we would like to flip the prediction  $A_{C+Q}$  from  $C_{\text{cand.}}$  (e.g., 'Japanese' to  $Q_{\text{cand.}}$  'French'; And in query-dominant case, we would like to flip the prediction  $A_{C+Q}$  from  $Q_{\text{cand.}}$ 'Malaysia' to  $C_{\text{cand.}}$  'Australia'.

To do this, we intervene in two layers: the first attention layer where the context information is transferred to the last token residual stream, and the attention layer where the most context information is written into the last token residual stream. These two layers correspond to the first blue spike and the highest blue spike in Figures 6d, 7d and 8d. For Llama-3, it is layers 17 and 24, respectively. For Mistral, it is layers 18 and 24, respectively. For Pythia, it is layers 19 and 24, respectively.



Figure 5: Logit lens on Llama-3 shows how model's top-1 predictions shift along the residual stream from abstract concepts (e.g., 'languages') to concrete instances (e.g., 'English' or 'Japanese') across layers. Red indicates high probability.



(a) **Q-only:** Token logit from accumulated residual stream. (b) Token logit difference (Logit in **C+Q** - Logit in **Q-only**)  $(R_{T,l}^1, R_{T,l}^2)$  are visualized per layer. from accumulated residual stream  $(R_{T,l}^1, R_{T,l}^2)$ .



(c) **C+Q**: Token logit from accumulated residual stream. (d) **C+Q**: candidate logit difference from attention and MLP  $(R_{T,l}^1, R_{T,l}^2)$  are visualized per layer. output.  $A_{T,l}$  and  $M_{T,l}$  are visualized per layer.

Figure 6: Additional logit attribution results for Llama-3 8B.

Specifically, in the context-dominant case, at the last token position, we set the attention scores corresponding to all tokens in the context to be  $-\infty$ , therefore, attention weight (which sums up to 1) is only a distribution over the query tokens. We perform this intervention to block information flow from the context to the last token position, and we only allow models to attend to the query part. Similarly, in the query-dominant case, we set the attention scores corresponding to all tokens in the query to be  $-\infty$ , allowing the models to only retrieve information from the context.

1082

1084

1085

1086

1088

1091

1092

1093

1094

To compare the knockout effect of the two criti-

cal layers with other layers, we select two random1095lower layers and two random higher layers. We re-1096port the average intervention results for three runs.1097

#### H.2 Additional Results

Results for Llama and Mistral are presented in Ta-	10
ble 10 and Table 9, respectively.	11
With the targeted two-critical-layer intervention:	11

1098

 Llama: 465/1000 context-dominant datapoints flip to query-based candidates, while
 407/1000 remain context-based. Conversely,
 225/1000 query-dominant datapoints shift to
 1102



(a) **Q-only:** Token logit from accumulated residual stream. (b) Token logit difference (Logit in **C+Q** - Logit in **Q-only**)  $(R_{T,l}^1, R_{T,l}^2)$  are visualized per layer. from accumulated residual stream  $(R_{T,l}^1, R_{T,l}^2)$ .



(c) C+Q: Token logit from accumulated residual stream. (d) C+Q: candidate logit difference from attention and MLP  $(R_{T,l}^1, R_{T,l}^2)$  are visualized per layer. output.  $A_{T,l}$  and  $M_{T,l}$  are visualized per layer.

Figure 7: Logit Attribution Results For Mistral 7B.

1106 1107	context-based candidates, while 704/1000 re- main query-based.	
1108	• Mistral: 437/1000 context-dominant data-	
1109	points flip to query-based candidates, while	Qi
1110	514/1000 remain context-based. Similarly,	
1111	232/1000 query-dominant datapoints shift to	
1112	context-based candidates, while 713/1000 re-	Čt
1113	main query-based.	Tabl

 Pythia: 470/1000 context-dominant datapoints flip to query-based candidates, while 486/1000 remain context-based. Conversely, 294/1000 query-dominant datapoints shift to context-based candidates, while 648/1000 remain query-based.

1120Across all models, approximately 950 datapoints1121remain context- or query-based candidates, instead1122of random non-identifiable answers, indicating that1123our intervention preserves model capabilities.

1114

1115

1116

1117

1118

1119

	Orig.	L17+	L17+L24		ow	2 High			
	Prob.	Prob.	$\Delta$	Prob.	$\Delta$	Prob.	$\Delta$		
	Context-Dominant								
Ctx Query	22.6 8.2	14.6 12.4	-8.0 +4.2	19.9 8.2	-2.7 +0.0	19.0 9.2	-3.6 +1.0		
Query-Dominant									
Query Ctx	33.0 6.5	25.0 10.7	-8.0 +4.2	25.5 7.7	-7.5 +1.2	31.7 6.4	-1.3 -0.1		

Table 9: Effect of attention knockout on context- (Ctx) and query-based (Query) candidate probabilities on Llama-3. "Orig." = No intervention, "2 Low" = Two lower layers (<17), "2 High" = Two higher layers (>24). "Diff." represents the probability difference, and  $\Delta$  denotes the change from the original setting. (Mistral 7B)



(a) Q-only: Token logit from accumulated residual stream. (b) Token logit difference (Logit in C+Q - Logit in Q-only)  $(R_{T,l}^1, R_{T,l}^2)$  are visualized per layer. from accumulated residual stream  $(R_{T,l}^1, R_{T,l}^2)$ .



(c) **C+Q**: Token logit from accumulated residual stream. (d) **C+Q**: candidate logit difference from attention and MLP  $(R_{T,l}^1, R_{T,l}^2)$  are visualized per layer. output.  $A_{T,l}$  and  $M_{T,l}$  are visualized per layer.

Figure 8: Logit Attribution Results For Pythia 6.9B.

	Orig.	L17-	+L24	2 L	ow	2 High				
	Prob.	Prob.	$\Delta$	Prob.	Δ	Prob.	$\Delta$			
	Context-Dominant									
Ctx Query	25.5 8.6	13.1 14.8	-12.4 +6.2	20.9 8.9	-4.6 +0.3	21.1 12.6	-4.4 +4.0			
	Query-Dominant									
Query Ctx	35.2 6.6	26.8 11.3	-8.4 +4.7	25.7 7.7	-9.5 +1.1	33.4 7.1	-1.8 +0.5			

	•
Table 10: Effect of attention knockout on context- (Ctx)	Т
and query-based (Query) candidate probabilities on	а
Llama-3. "Orig." = No intervention, "2 Low" = Two	L
lower layers (<17), "2 High" = Two higher layers (>24).	lo
"Diff." represents the probability difference, and $\Delta$ de-	"]
notes the change from the original setting. (Llama-3)	n

	Orig.	L17+	-L24	2 L	ow	2 High				
	Prob.	Prob.	$\Delta$	Prob.	$\Delta$	Prob.	$\Delta$			
		Cor	ntext-D	omina	nt					
Ctx Query	22.3 7.4	13.6 11.5	-8.7 +4.1	18.3 8.4	-4.0 +1.0	20.5 7.8	-1.8 +0.4			
		Qu	ery-D	ominan	t					
Query Ctx	26.6 6.3	20.8 9.6	-5.8 +3.3	20.5 6.9	-6.1 +0.6	25.6 6.2	-1.0 -0.1			

able 11: Effect of attention knockout on context- (Ctx) nd query-based (Query) candidate probabilities on lama-3. "Orig." = No intervention, "2 Low" = Two ower layers (<17), "2 High" = Two higher layers (>24). Diff." represents the probability difference, and  $\Delta$  denotes the change from the original setting. (Pythia)



(b) Q-only: Query circuit in query-dominant case.

Figure 9: Activation patching under Q-Only condition reveals that query circuit is the same before and after adding the irrelevant context for Llama-3 8B.



(a) C+Q: Context circuit in context-dominant case.

C_REL 0	- 0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
*C_SUBJ 1	- 0.48	8 0.60	0.61	0.61	0.63	0.62	0.62	0.60	0.58	0.60	0.58	0.56	0.56	0.53	0.55	0.48	0.46	0.47	0.37	0.32	0.30	0.28	0.28	0.27	0.15	0.15	0.13	0.13	0.11	0.07	0.04	0.00	0.00	0.00
C_REL 2	- 0.01	L 0.03	0.03	0.06	0.07	0.10	0.10	0.11	0.11	0.11	0.11	0.11	0.10	0.10	0.10	0.13	0.12	0.12	0.15	0.12	0.11	0.11	0.11	0.11	0.07	0.07	0.07	0.07	0.07	0.02	0.01	0.00	0.00	0.00
*C_OBJ 3	- 1.43	3 1.49	1.49	1.45	1.46	1.43	1.44	1.43	1.39	1.37	1.36	1.30	1.24	1.13	1.08	0.98	0.92	0.89	0.72	0.62	0.57	0.55	0.53	0.53	0.28	0.27	0.24	0.23	0.19	0.12	0.07	0.00	0.00	0.00
Q_REL 4	- 0.04	0.08	0.09	0.17	0.19	0.22	0.22	0.25	0.28	0.27	0.26	0.27	0.23	0.24	0.23	0.23	0.19	0.19	0.22	0.22	0.20	0.19	0.19	0.19	0.16	0.16	0.16	0.16	0.16	0.06	0.04	0.00	0.00	0.00
Q_SUBJ 5	-0.0	0-0.01	-0.00	-0.00	0.02	0.04	0.04	0.06	0.08	0.10	0.11	0.11	0.12	0.12	0.12	0.11	0.10	0.10	0.09	0.06	0.06	0.06	0.06	0.05	0.03	0.03	0.03	0.03	0.03	0.01	0.01	0.00	0.00	0.00
Q_REL 6	-0.0	0.01	0.01	0.02	0.05	0.07	0.09	0.12	0.15	0.17	0.20	0.23	0.25	0.30	0.36	0.44	0.46	0.47	0.60	0.69	0.72	0.74	0.75	0.77	0.93	0.94	0.96	0.98	1.01	1.09	1.12	1.19	1.19	1.19
	_	~	<u>.</u>	~	8	5	5	~	6	6	-	~	Å.	~	2	5	5	~	6	6	-	~	0.	6	2	5	5	~	0	6	-	~	~	×
	7	2	2	$\checkmark$	Å.	1	7	~	2	V	Ì	I	3	J	Ì	Ì	Ĭ	J	Ì	Ì	Ň	Y	Y	Y	Ŷ	Ŷ	Ÿ	Ŷ	Ŷ	Ŷ,	5	3	2	100

(b) C+Q: Context circuit in query-dominant case.



(c) C+Q: Query circuit in context-dominant case.



(d) C+Q: Query circuit in query-dominant case.

Figure 10: Activation patching under C+Q condition for Mistral 7B.



(b) Q-only: Query circuit in query-dominant case.

Figure 11: Activation patching under Q-Only condition reveals that query circuit is the same before and after adding the irrelevant context for **Mistral 7B**.



(b) Q-only: Query Circuit in query-dominant case.

Figure 12: Activation patching under Q-Only condition reveals that query circuit is the same before and after adding the irrelevant context for **Pythia 6.9B**.





(c) Query circuit in context-dominant case.



(d) Query circuit in query-dominant case.

Figure 13: Activation patching under C+Q condition for Pythia 6.9B.

Relation	Template	Ctx Type	<b>Total Rows</b>		
P1001	[X] is a legal term in [Y]	Place	664		
P101	The expertise of [X] is [Y].	Others	571		
P103	The mother tongue of [X] is [Y].	Language	919		
P106	[X] works as [Y].	Job	821		
P108	[X], who is employed by [Y].	Company	378		
P127	[X] owner [Y].	Company	616		
P1303	[X] plays the [Y].	Others	513		
P131	[X] is in [Y].	Place	775		
P136	[X] plays [Y].	Others	859		
P1376	[X], the capital city of [Y].	Place	179		
P138	[X], which is named after [Y].	Others	461		
P140	[X] is follower of [Y].	Others	432		
P1412	[X] communicated in [Y].		924		
P159	P159 [X] is headquartered in [Y].		801		
P17	[X], located in [Y].	Place	912		
P176	[X], produced by [Y].	Company	925		
P178	[X], a product developed by [Y].	Company	588		
P19	[X] is native to [Y].	Place	779		
P190	[X] is a twin city of [Y].	Place	671		
P20	[X] passed away at [Y].	Place	817		
P264	[X]'s label is [Y].	Company	53		
P27	[X], a citizen of [Y].	Place	958		

Relation	Template	Туре	<b>Total Rows</b>
P276	[X] is located in [Y].	Place	764
P279	[X], a type of [Y].	Others	900
P30	[X] is a part of the continent of [Y].	Place	959
P36	The capital city of [X] is [Y].	Place	471
P361	[X] is a part of [Y].	Others	746
P364	The original language of [X] was [Y].	Language	756
P37	The official language of [X] is [Y].	Language	900
P39	[X], who holds the position of [Y].	Job	485
P407	[X] was written in [Y].	Language	857
P413	[X] plays in the position of [Y].	Job	952
P449	[X] premiered on [Y].	Company	801
P463	[X] belongs to the organization of [Y].	Company	203
P47	[X] borders with [Y].	Place	649
P495	[X] was formed in [Y].	Place	905
P530	[X] ties diplomatic relations with [Y].	Place	950
P740	[X], founded in [Y].	Place	843
P937	[X] found employment in [Y].	Place	853

Table 12: Overview of Relations, Templates, Types, and Total Rows in the original Pararel Dataset. We take this dataset and construct the C + Q dataset, which has around 106.6M rows.

Category	Details
Example 1	
Context	Hanoi is a twin city of Bangkok.
Query	The mother tongue of Louis Legendre is
Class	Languages
Context Subject Possible Answers	Vietnamese, Tay, Hmong, Khmer, English, French Chinese
Context Object Possible Answers	Thai, Lao, Chinese, Malay, Khmer
Context-Based Candidates	Vietnamese, Thai
Context-Influenced?	True
Correct Class?	True
Exists Answer that Satisfies Both?	True
Example 2	
Context	Svend Asmussen plays the violin.
Query	Social-Economic Council is a legal term in
Class	Places (Countries, Cities, States, etc.)/Languages
Context Subject Possible Answers	Danmark, Danish (Svend Asmussen is a Violinis and jazz musician)
Context Object Possible Answers	Italy, Italian (Violin was originated in Italy)
Context-Based Candidates	Denmark
Context-Influenced?	True
Correct Class?	True
Exists Answer that Satisfies Both?	True
Example 3	
Context	Manchester Business School is headquartered in Manchester.
Query	Antipope Paschal III, who holds the <b>position</b> of
Class	Jobs/Positions/Roles
Context Subject Possible Answers	Professor, Lecturer, Instructor, Researcher, Depar ment Chair, Provost, Dean, Academic Adviso Teaching Assistant, Student, etc.
Context Object Possible Answers	N/A
Context-Based Candidates	Dean, Professor
Context-Influenced?	True
Correct Class?	True
Exists Answer that Satisfies Both?	True

Category	Details
Example 4	
Context	Saharsa district is in Bihar.
Query	Colbert Mountains is a part of the <b>continent</b> of
Class	Continents/Places
Context Subject Possible Answers	Asia
Context Object Possible Answers	Asia
Context-Based Candidates	Asia, India
Context-Influenced?	True
Correct Class?	True
Exists Answer that Satisfies Both?	True
Example 5	
Context	Lee Jong-hyun plays the guitar.
Query	Northern Foothills is a part of the continent of
Class	Continents
Context Subject Possible Answers	Asia
Context Object Possible Answers	Europe (Guitar originated in Spain)
Context-Based Candidates	South, Korea
Context-Influenced?	True
Correct Class?	False
Exists Answer that Satisfies Both?	False
Example 6	
Context	Grant Green plays jazz.
Query	David Gates plays the
Class	Role/Genre/Style/Position/Musical Instrument
Context Subject Possible Answers	guitarist, composer, musician, songwriter etc. (role of Grant Green), guitar (Musical Instrument that Grant Green plays), jazz, R&B, etc. (music genre of Grant Green)
Context Object Possible Answers	Jazz.
Context-Based Candidates	Bee, Beach
Context-Influenced?	False
Correct Class?	False
Exists Answer that Satisfies Both?	False
Example 7	

Category	Details
Context	Samuil Marshak passed away at Moscow.
Query	Jean Metcalfe, who is employed by
Class	Company/Person
Context Subject Possible Answers	Russia-1, Channel One Russia, RT, TV Rain, etc.
Context Object Possible Answers	Russia-1, Channel One Russia, RT, TV Rain, etc.
Context-Based Candidates	BBC, Radio
Context-Influenced?	False
Correct Class?	True
Exists Answer that Satisfies Both?	False