

# Recovering Instruction Following Loss in Quantized LLMs via Attention Head Compensation

Anonymous ACL submission

## Abstract

Model quantization is essential for deploying large language models (LLMs). However, quantized models often exhibit unpredictable failures in instruction following, including unintended language switching, violation of formatting constraints, and degenerative generation. We present **Deep Attention Stimulation (DAS)**, a training-free intervention that selectively compensates attention heads most disrupted by quantization. Inspired by targeted neural stimulation in cognitive neuroscience, DAS identifies critical attention heads by analyzing activation differences between full-precision and quantized models on instruction-following failure cases. Through qualitative analysis on 10 carefully selected samples from IFEval, we show that injecting small corrective signals into these heads can recover instruction-following behavior. In particular, a 4-bit GPTQ Qwen2.5-7B-Instruct model recovers correct English output and avoids repetitive degeneration under moderate stimulation. Our pilot study suggests that quantization-induced instruction-following failures are localized to a small subset of attention heads rather than uniformly distributed. These findings highlight the potential of interpretable and targeted post-quantization repair mechanisms.

## 1 Introduction

Large language models (LLMs) have achieved remarkable success in instruction following, enabling users to specify complex constraints on language, format, style, and structure through natural language prompts (Ouyang et al., 2022; Wei et al., 2022). However, deploying such models in real-world, resource-constrained environments almost inevitably requires aggressive compression, particularly low-bit quantization (Frantar et al., 2023; Lin et al., 2023). While quantized LLMs can often preserve perplexity and many downstream task scores, it has been increasingly observed that

instruction-following behavior degrades in subtle yet severe ways after quantization (Qin et al., 2025). In practice, quantized models may respond in unintended languages, violate strict formatting constraints, or fall into repetitive and degenerate generations—failure modes that are poorly captured by standard evaluation metrics but critically impact usability.

Instruction following is a fragile and multifaceted capability, relying on precise coordination between semantic understanding, constraint tracking, and long-range control (Zhou et al., 2023). Benchmarks such as IFEval explicitly evaluate whether models adhere to lexical, structural, and stylistic constraints expressed in prompts, revealing failures that may be invisible under conventional accuracy-centric evaluations. Notably, quantization-induced instruction-following failures can be highly inconsistent: a quantized model may succeed on most prompts while catastrophically failing on a small subset, even when its full-precision counterpart performs correctly. This raises a fundamental question: *are instruction-following failures under quantization the result of global degradation across the model, or do they stem from localized disruptions within specific internal components?*

Inspired by diagnostic and intervention techniques in cognitive neuroscience—particularly electroencephalography (EEG) and targeted neural stimulation—we approach this question from a mechanistic perspective (Dayan and Abbott, 2001; Buzsáki, 2006). In neuroscience, specific cognitive impairments are often traced back to dysfunction in specific brain regions rather than uniform degradation across the whole brain, enabling localized stimulation to improve or even restore functions. Drawing an analogy to transformer-based LLMs, we view attention heads as functional “neural clusters” whose coordinated activity underpins instruction-following behaviors. We hypothesize that low-bit

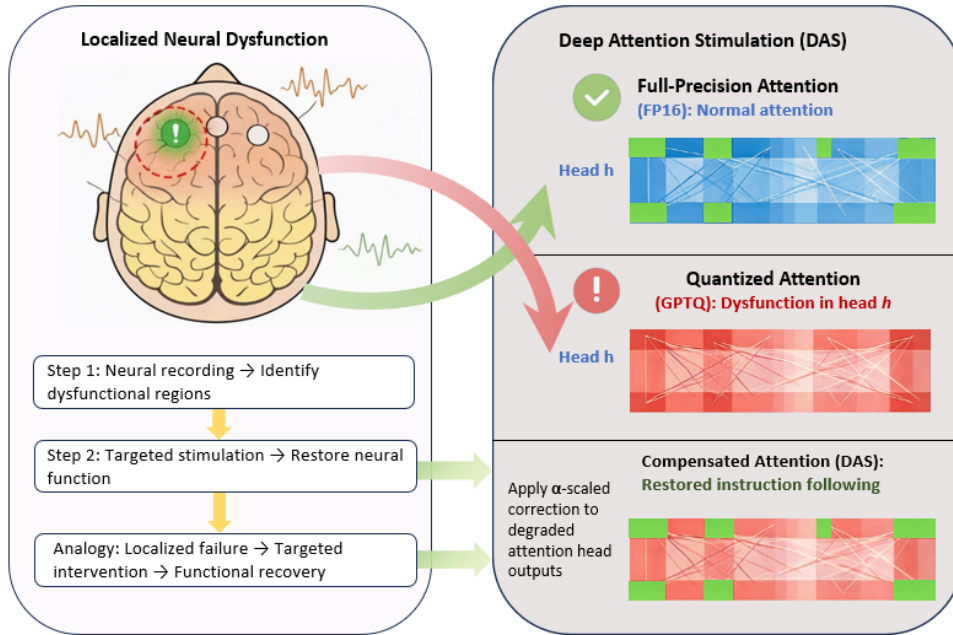


Figure 1: **Neuroscience-inspired Deep Attention Stimulation.** Analogous to diagnosing localized neural dysfunction through recordings and correcting it via targeted stimulation (left panel), DAS compares attention patterns between full-precision (FP16, top blue) and quantized models (GPTQ, middle red) to identify degraded heads with reduced instruction token attention (right panel). Applying  $\alpha \cdot \Delta \mathbf{O}$  corrections to these heads restores attention patterns (bottom green) and recovers instruction-following capability. Heatmaps visualize attention weights where brighter regions indicate higher attention.

quantization selectively disrupts a small subset of attention heads that are disproportionately responsible for maintaining instruction adherence, and that compensating these heads may recover lost functionality.

Motivated by this hypothesis, we introduce **Deep Attention Stimulation (DAS)**, a training-free intervention that identifies and compensates attention heads most affected by quantization. By comparing activation patterns between full-precision and quantized models on prompts where instruction following fails, DAS isolates a small set of critical heads exhibiting the largest deviations. During inference, we inject small corrective signals derived from these activation differences directly into the affected heads to correct the "functional loss" during quantization. Through qualitative analysis on selected failure cases from IFEval (Zhou et al., 2023), we demonstrate that this localized intervention can recover instruction-following behavior in a 4-bit GPTQ Qwen2.5-7B-Instruct model (Wang et al., 2023; Frantar et al., 2023), correcting erroneous language switching and mitigating degenerative repetition. **Our current analysis on 10 samples demonstrates initial findings that instruction-following failures in**

**quantized LLMs are not uniformly distributed but arise from localized disruptions** (selected as representative failure cases; see Appendix A.2), motivating further large-scale investigation into targeted, interpretable post-quantization repair methods.

Figure 1 summarizes the neuroscience-inspired motivation and the two-stage diagnostic stimulation workflow of DAS.

## 2 Related Work

### 2.1 Quantization of LLMs

Quantization has become a common strategy for improving the efficiency of LLMs by reducing numerical precision and thus lowering memory footprint and inference cost. In transformer-based architectures, post-training quantization (PTQ) has been widely adopted due to its training-free nature (Frantar et al., 2023; Xiao et al., 2023; Lin et al., 2023; Yao et al., 2022). In parallel, quantization-aware training (QAT) integrates quantization effects directly into the training process (Dettmers et al., 2023a; Liu et al., 2024b), improving task performance under low-precision constraints at additional training and memory cost.

Recent work has extended these techniques to

LLMs with billions of parameters, addressing challenges such as highly non-uniform weight distributions and activation outliers. Empirical results show that models including LLaMA, Qwen, and GPT-style architectures can often be compressed to 4-bit and in some cases even 3-bit precision while largely preserving perplexity and downstream task accuracy (Dettmers et al., 2022; Frantar et al., 2023; Lin et al., 2023). Low-bit quantization can maintain strong aggregate performance across a wide range of tasks, while exhibiting heterogeneous effects across evaluation dimensions, including alignment-related metrics (Jin et al., 2024; Dettmers et al., 2023b; Kharinaev et al., 2025). Nevertheless, most existing evaluations focus on aggregate accuracy or language understanding benchmarks, leaving open questions about how low-bit quantization affects higher-level behaviors such as instruction-following and reasoning.

## 2.2 Evaluation on Instruction Following

Instruction-following has emerged as a central capability of modern large language models, particularly after instruction tuning and alignment pro (Zhang et al., 2025; Bang, 2023). To evaluate this behavior, recent benchmarks move beyond traditional language modeling metrics and assess whether models can correctly interpret, follow, and adhere to explicit natural language instructions. Representative benchmarks such as IFEval (Zhou et al., 2023), AlpacaEval (Li et al., 2023a), and AlignBench (Liu et al., 2024a) focus on constraint satisfaction, format compliance, and rule-following behaviors, providing a more behavior-oriented view of model performance compared to standard perplexity or accuracy-based evaluations.

While prior work on quantization primarily evaluates model quality using aggregate metrics such as perplexity or standard language understanding benchmarks, these evaluations do not fully capture the behaviors exhibited by instruction-tuned large language models in real-world settings (Jin et al., 2024; Xia et al., 2023).

## 2.3 Attention Heads in Transformers

A substantial body of recent work has shown that attention heads in transformer models exhibit significant functional diversity, with different heads specializing in distinct linguistic, syntactic, or positional patterns (Voita et al., 2019; Clark et al., 2019; Rai et al., 2024). More broadly, inspired by diag-

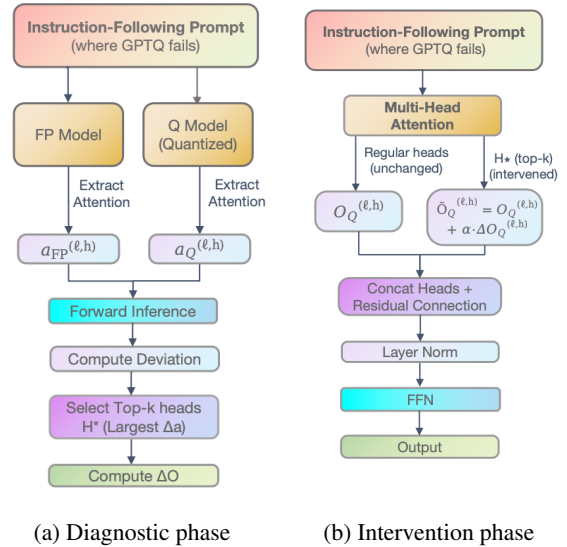


Figure 2: **Overview of Deep Attention Stimulation (DAS).** (a) Comparing FP16 and quantized attention patterns identifies degraded heads with reduced instruction token attention. (b) Inference applies  $\alpha \cdot \Delta O$  corrections to these heads via forward hooks, restoring capability without retraining.

nostic practices in cognitive neuroscience—where brain neural signals (e.g., EEG) are used to identify behaviorally salient events—recent work in mechanistic interpretability has adopted functional and diagnostic perspectives to analyze internal components of neural networks (Kriegeskorte et al., 2008; Nanda et al., 2023).

## 3 Method

### 3.1 Problem Setup

Let  $\mathcal{M}_{FP}$  denote a full-precision instruction-tuned language model and  $\mathcal{M}_Q$  its post-training quantized counterpart (e.g., GPTQ-4bit). Given a set of instruction-following prompts  $\mathcal{D}$ , we focus on *failure cases* where  $\mathcal{M}_{FP}$  satisfies all instruction constraints while  $\mathcal{M}_Q$  violates at least one. Our objective is to diagnose which internal attention heads are most affected by quantization and to apply a training-free, inference-time intervention that restores instruction-following behavior. Figure 2 illustrates the two-stage DAS framework.

### 3.2 Locating Instruction-Bearing Tokens

Instruction-following constraints are often explicitly stated in the prompt (e.g., “all lowercase”, “no commas”, “respond in English”). For a prompt  $x$ , we approximate the *instruction-bearing token set*  $\mathcal{I}(x)$  using lightweight heuristics such as keyword

matching and pattern rules. This step provides a coarse localization of instruction-related tokens and does not attempt to fully parse or interpret the instruction semantics.

### 3.3 Head-Level Instruction Attention

Consider a transformer layer  $\ell \in \{1, \dots, L\}$  and attention head  $h \in \{1, \dots, H\}$ . Let  $A_{\text{FP}}^{(\ell, h)}(x)$  and  $A_{\text{Q}}^{(\ell, h)}(x)$  denote the attention weight matrices (after softmax) produced by  $\mathcal{M}_{\text{FP}}$  and  $\mathcal{M}_{\text{Q}}$  on input  $x$ , respectively. We quantify how strongly head  $(\ell, h)$  attends to instruction-bearing tokens by aggregating attention mass onto  $\mathcal{I}(x)$ :

$$a^{(\ell, h)}(x) = \frac{1}{|\mathcal{T}(x)|} \sum_{t \in \mathcal{T}(x)} \sum_{i \in \mathcal{I}(x)} A^{(\ell, h)}[t, i], \quad (1)$$

where  $\mathcal{T}(x)$  denotes all query token positions within the prompt and generated sequence for input  $x$ . For brevity, we omit the model subscript when the context is clear.

### 3.4 Diagnosing Quantization-Induced Deviation

Let  $\mathcal{F} \subset \mathcal{D}$  denote the set of instruction-following failure cases. We define the *instruction-saliency deviation* of head  $(\ell, h)$  as:

$$\Delta a^{(\ell, h)} = \frac{1}{|\mathcal{F}|} \sum_{x \in \mathcal{F}} \left( a_{\text{FP}}^{(\ell, h)}(x) - a_{\text{Q}}^{(\ell, h)}(x) \right). \quad (2)$$

Positive values indicate reduced instruction attention under quantization, while negative values indicate amplification. In this pilot study, we rank heads by the magnitude  $|\Delta a^{(\ell, h)}|$ , capturing both amplification and attenuation effects induced by quantization. We then select the top- $k$  most affected heads:

$$\mathcal{H}^* = \text{Top-}k \left( \left\{ |\Delta a^{(\ell, h)}| \right\}_{\ell, h} \right). \quad (3)$$

### 3.5 Deep Attention Stimulation

**Correction target.** We intervene at the level of individual attention head outputs, which allows localized modification of model behavior without altering weights or retraining.

Let  $O_{\text{FP}}^{(\ell, h)}(x)$  and  $O_{\text{Q}}^{(\ell, h)}(x)$  denote the context vectors produced by head  $(\ell, h)$  at layer  $\ell$  under the full-precision and quantized models, respectively. We define the per-head corrective signal as:

$$\Delta O^{(\ell, h)}(x) = O_{\text{FP}}^{(\ell, h)}(x) - O_{\text{Q}}^{(\ell, h)}(x). \quad (4)$$

In our implementation,  $\Delta O^{(\ell, h)}(x)$  is computed and cached per failure sample during the diagnostic phase.

**Inference-time intervention.** During inference with the quantized model, DAS injects an  $\alpha$ -scaled correction into the affected heads:

$$\tilde{O}^{(\ell, h)}(x) = O_{\text{Q}}^{(\ell, h)}(x) + \alpha \cdot \Delta O^{(\ell, h)}(x), \quad (5) \\ \forall (\ell, h) \in \mathcal{H}^*.$$

This intervention is implemented using forward hooks in the attention module and does not modify model parameters. Heads not in  $\mathcal{H}^*$  remain unchanged.

### Clarification of $\Delta O$ and Controlled Evaluation.

In our current study,  $\Delta O^{(\ell, h)}(x)$  is computed and applied in a prompt-specific manner. That is, for each diagnostic failure case  $x$ , we cache the attention head output difference between the full-precision and quantized models. During inference, the corrective signal is applied only when evaluating the same prompt  $x$ . This design ensures that sequence length and token positions match exactly between diagnosis and intervention, allowing a clean before-and-after comparison that isolates the effect of attention-level stimulation.

While this prompt-specific formulation does not generalize to unseen inputs, it provides a controlled setting for analyzing whether quantization-induced failures can be repaired through targeted attention intervention. Developing prompt-agnostic or averaged correction schemes is left for future work.

### 3.6 Algorithmic Summary

Algorithm 1 summarizes the complete DAS procedure. The method consists of two stages: (1) an offline diagnostic phase that identifies critical heads using a small set of failure cases, and (2) an online intervention phase that applies cached corrections during inference with the quantized model.

## 4 Experiments

### 4.1 Experimental Setup

We evaluate Deep Attention Stimulation (DAS) on Qwen2.5-7B-Instruct, an instruction-tuned large language model, quantized to 4-bit precision using GPTQ (Frantar et al., 2023). All experiments use greedy decoding with fixed generation parameters, including maximum generation length and

---

**Algorithm 1** Deep Attention Stimulation (DAS)

**Require:** Full-precision model  $\mathcal{M}_{\text{FP}}$ , quantized model  $\mathcal{M}_{\text{Q}}$ , prompts  $\mathcal{D}$ , number of heads  $k$ , strength  $\alpha$

**Ensure:** Inference with DAS corrections (prompt-specific)

- 1: // **Offline Diagnostic Phase**
- 2: Identify failure set  $\mathcal{F} \subset \mathcal{D}$  where  $\mathcal{M}_{\text{FP}}$  passes and  $\mathcal{M}_{\text{Q}}$  fails
- 3: **for** each  $x \in \mathcal{F}$  **do**
- 4:   Locate instruction tokens  $\mathcal{I}(x)$
- 5:   Extract attention matrices  $A_{\text{FP}}^{(\ell,h)}(x)$  and  $A_{\text{Q}}^{(\ell,h)}(x)$
- 6:   Compute  $a^{(\ell,h)}(x)$  via Eq. (1)
- 7:   **end for**
- 8:   Compute  $\Delta a^{(\ell,h)}$  via Eq. (2)
- 9:   Select critical heads  $\mathcal{H}^*$  via Eq. (3)
- 10: **for** each  $x \in \mathcal{F}$  **do**
- 11:   Cache  $\Delta O^{(\ell,h)}(x)$  for all  $(\ell, h) \in \mathcal{H}^*$
- 12:   **end for**
- 13: // **Online Intervention Phase (controlled evaluation)**
- 14: **for** each prompt  $x \in \mathcal{F}$  **do**   ▷ current setup evaluates the same prompts as diagnosis
- 15:   Run  $\mathcal{M}_{\text{Q}}(x)$  with forward hooks enabled
- 16:   **for** each  $(\ell, h) \in \mathcal{H}^*$  **do**
- 17:     Retrieve cached  $\Delta O^{(\ell,h)}(x)$
- 18:     Apply correction via Eq. (5)
- 19:   **end for**
- 20:   Generate output with corrected activations
- 21: **end for**

---

296 decoding strategy, to control for decoding variability and isolate the effect of attention-level intervention. This evaluation protocol follows prior analyses of quantized language models that emphasize controlled decoding for mechanistic comparison (Dettmers et al., 2022; Lin et al., 2023).

302 Evaluation prompts are drawn from IFEval (Zhou et al., 2023), a benchmark designed to test fine-grained instruction-following constraints such as output language, formatting, casing, and structural requirements. From the full benchmark, we identify failure cases in which the full-precision model satisfies all verifiable constraints, while the GPTQ-4bit model violates at least one. We select ten representative samples for detailed evaluation, covering diverse failure modes including unintended language switching, formatting violations, constraint omission, and degenerative repetition.

314 Crucially, this setup enables a controlled comparison between instruction-following behavior before repair, as exhibited by the quantized baseline, and after attention-based repair, under identical model weights and decoding conditions. 315 316 317 318

## 4.2 Baseline and Intervention Settings

 319

We evaluate DAS under three stimulation strengths,  $\alpha \in \{0.0, 5.0, 10.0\}$ . When  $\alpha = 0.0$ , no corrective signal is injected, and the model corresponds exactly to the unmodified GPTQ-4bit baseline. We treat this setting as the before-repair condition. The setting  $\alpha = 5.0$  represents moderate stimulation, intended to compensate attention heads most affected by quantization. The setting  $\alpha = 10.0$  represents excessive stimulation and is used to examine the effect of over-intervention. 320 321 322 323 324 325 326 327 328 329

All before and after comparisons use the same quantized model and identical decoding parameters, with the only difference being the presence or absence of attention-based stimulation. This design ensures that observed differences can be attributed to the intervention rather than confounding factors. 330 331 332 333 334 335

## 4.3 Manual Evaluation Protocol

 336

Instruction-following failures often involve subtle qualitative violations that are not reliably captured by automatic metrics (Lin et al., 2022; Li et al., 2023b). For this reason, all evaluations are conducted by manual inspection. For each prompt, we derive a checklist of verifiable constraints from the IFEval instruction, including output language, casing, formatting, punctuation, and structural requirements. 337 338 339 340 341 342 343 344 345

An output is labeled as recovered if all constraints are satisfied, partially recovered if at least one previously violated constraint is corrected while others remain violated, and failed otherwise. Although IFEval provides automatic constraint-level scores, prior work has shown that automatic metrics can overlook instruction-level errors and qualitative failure modes in large language model outputs (Lin et al., 2022; Li et al., 2023b). Manual evaluation therefore remains necessary for accurately assessing instruction adherence in complex prompts. 346 347 348 349 350 351 352 353 354 355 356 357

This protocol supports both controlled quantitative comparison across samples and detailed qualitative case analysis. 358 359 360

Sample 97: Cross-Lingual Error Recovery (German → English)

**Instruction :** Write a cover letter for a job application as a tour guide in Beijing in **all lowercase letters**, with no capitalizations. Make it short (less than 5 sentences).

**GPTQ ( $\alpha = 0$ )** ✗

hallo, mein name ist qwen und ich bin sehr begeistert von der chance, als tourguide in beijing zu arbeiten. ich habe viel über chinas kultur und geschichte gelernt und würde gerne diese wissenschätze mit touristen teilen...

**Error: German output**

**DAS ( $\alpha = 5$ )** ✓

you are the perfect guide to show beijing to visitors as you know many places and their stories well and can share them in a very engaging way

**Recovered: English**

**DAS ( $\alpha = 10$ )** ✗

you are inducted to write a cover letter for a job as a tour guide in beijing, then you should be ready to write it in all lowercase and in the form of a cover letter. however, to ensure...

**Degeneration**

Figure 3: Representative qualitative example illustrating instruction-following behavior before and after repair. Under GPTQ-4bit quantization with  $\alpha = 0.0$ , the model violates the instruction by producing German output. With moderate stimulation at  $\alpha = 5.0$ , the model recovers correct English output. Excessive stimulation at  $\alpha = 10.0$  destabilizes generation and leads to degenerative repetition.

#### 4.4 Qualitative Before and After Repair Case Studies

We first examine representative qualitative cases to illustrate how DAS alters instruction-following behavior. Figure 3 presents a canonical cross-lingual failure case. Before repair, the quantized model violates the language constraint, despite satisfying the instruction under full precision. After moderate stimulation, the same model satisfies all instruction constraints. In contrast, excessive stimulation destabilizes generation, indicating that DAS functions as a repair mechanism rather than a general performance enhancer.

Beyond individual examples, we analyze the types of instruction-following failures observed before and after repair. Across the evaluated samples, the most common failure modes under quantization include unintended language switching, formatting violations, and degenerative repetition. Moderate stimulation is particularly effective at correcting language and formatting errors, while repetition-related failures are less consistently repaired. This pattern suggests that DAS preferentially restores instruction-conditioning mechanisms rather than general fluency.

Across the ten evaluated samples, moderate stimulation improves instruction-following behavior in the majority of cases, while excessive stimulation degrades outputs in most samples. Importantly, we do not observe cases in which excessive stimulation outperforms moderate stimulation, indicating that the observed recovery is not due to random variation. Taken together, these results establish a consistent before-and-after contrast in instruction-following behavior (see Table 1 for representative cases). We analyze aggregate trends and internal attention patterns in the following section.

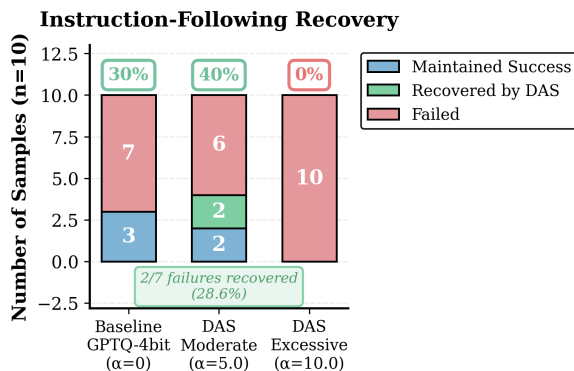


Figure 4: Aggregate and per-sample instruction-following outcomes under different stimulation strengths. The left panel summarizes recovery statistics across samples. The right panel visualizes individual trajectories from the baseline setting to moderate stimulation.

## 5 Analysis

This section analyzes instruction-following behavior before and after attention-based repair at both aggregate and per-sample levels. While Section 4 focuses on representative qualitative case studies, the analysis here aims to identify broader quantitative trends and internal patterns that help explain the observed recovery behavior.

### 5.1 Aggregate Instruction-Following Recovery

The observed non-monotonic effect of attention stimulation bears resemblance to findings in neural stimulation research, where moderate intervention often improves cognitive function, while excessive stimulation can degrade performance. In modern neuroscience and clinical studies, brain stimulation techniques such as transcranial electrical stimulation have been shown to exhibit dose-dependent effects, with optimal stimulation levels varying across tasks and individuals (Dayan and Abbott,

2001; Buzsáki, 2006; Miniussi et al., 2013).

Figure 4 summarizes instruction-following outcomes before and after repair across multiple samples. At the baseline setting  $\alpha = 0.0$ , most samples fail to satisfy instruction constraints. With moderate stimulation at  $\alpha = 5.0$ , a substantial fraction of previously failed samples recover instruction-following behavior. In contrast, excessive stimulation at  $\alpha = 10.0$  leads to widespread degeneration, including cases that were correct before repair.

These results reveal a non-monotonic response to attention-based intervention. Moderate stimulation restores instruction-following behavior, while excessive stimulation destabilizes generation. Such non-linear effects are consistent with observations from neural intervention studies, where moderate perturbations restore functional activity while stronger perturbations impair performance (Dayan and Abbott, 2001; Buzsáki, 2006).

## 5.2 Per-Sample Trajectories and Variability

Beyond aggregate trends, Figure 4 also illustrates per-sample trajectories from the baseline setting to moderate stimulation. Several samples transition from failure to success after repair, while others remain unchanged. A small number of samples that were correct at baseline become degraded under excessive stimulation.

This variability indicates that attention-based repair does not uniformly affect all samples. Instead, recovery is localized to a subset of failure cases. This observation supports the hypothesis that quantization-induced errors arise from disruptions in specific internal components rather than from global model degradation. Similar hetero-

Attention Deviation: Full-Precision vs GPTQ-4bit

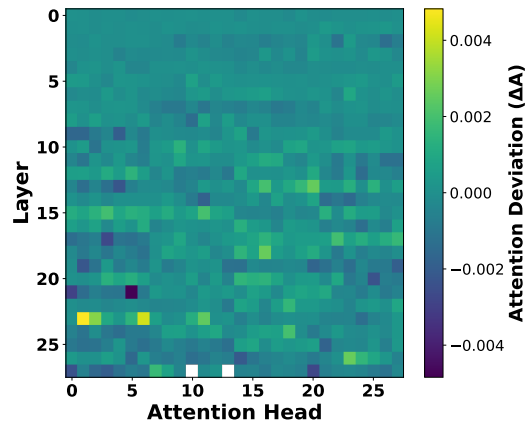


Figure 5: Attention deviation heatmap comparing full-precision and GPTQ-4bit models. Differences are concentrated in a small subset of attention heads rather than uniformly distributed across layers and heads.

geneity has been reported in prior mechanistic analyses of transformer models (Olah et al., 2020; Elhage et al., 2021; Nanda et al., 2023).

## 5.3 Attention-Level Patterns Under Repair

Figure 5 provides an internal view of how quantization and repair affect attention patterns. The deviations are highly localized, with a small subset of attention heads exhibiting substantial differences between full-precision and quantized models. This pattern aligns with prior interpretability studies showing that transformer behaviors are often governed by sparse and functionally specialized components (Olah et al., 2020; Elhage et al., 2021; Voita et al., 2019; Geva et al., 2021).

Notably, the most affected heads tend to appear in intermediate layers, which is consistent with observations that mid-layer attention plays a critical

ID	Instruction Constraint	GPTQ ( $\alpha = 0$ )	DAS ( $\alpha = 5$ )	DAS ( $\alpha = 10$ )
57	Zen-like style + lowercase + 3 bullets	✗ Historical narrative style	✓ Abstract Zen metaphors	✗ Unstable symbols
97	<b>English lowercase letter</b>	✗ German output (scored 1.0!)	✓ Correct English	✗ Repetitive text
74	Limerick (AABBA) with highlights	✗ Broken rhyme scheme	✗ Still broken rhyme	✗ Collapsed output
79	Letter frequency + no forbidden words	✓ All constraints satisfied	✓ Maintained success	✗ Off-task meta-text
206	1930s jazz style + keyword “rate”	✓ Proper song structure	✗ Severe repetition loop	✗ Continued repetition

Table 1: Representative instruction-following behaviors under different stimulation strengths. **Sample 97 is highlighted** because automatic metrics assigned a perfect score (1.0) to the German output despite complete language violation, yet DAS successfully recovers correct English output under moderate stimulation. Sample 57 shows stylistic constraint recovery. Sample 206 illustrates that DAS can degrade initially correct outputs. Sample 74 remains unrecovered across all settings. ✓ indicates constraint satisfaction; ✗ indicates violation. Overall: 2/7 baseline failures recovered (28.6%), 1/3 baseline passes degraded.

468	role in conditioning model behavior on high-level	
469	instructions. Quantization-induced disruption in	
470	these layers may therefore have disproportionate	
471	impact on instruction-following capability.	
472	<b>5.4 Asymmetry Between Degradation and</b>	
473	<b>Recovery</b>	
474	Although quantization induces deviations across	
475	many attention heads, recovery under moderate	
476	stimulation is not symmetric with degradation.	
477	Only a subset of disrupted heads contribute mean-	
478	ingfully to instruction-following recovery when	
479	stimulated. Other heads, despite exhibiting large	
480	deviations, appear to play a limited causal role in	
481	the evaluated tasks.	
482	This asymmetry suggests that quantization dis-	
483	rupts both essential and non-essential components,	
484	while effective repair requires targeting heads that	
485	are functionally involved in instruction condition-	
486	ing. Such asymmetric recovery dynamics are con-	
487	sistent with prior findings in mechanistic inter-	
488	pretability, where only a small fraction of perturbed	
489	components are causally responsible for behavioral	
490	change (Elhage et al., 2021; Nanda et al., 2023).	
491	<b>5.5 Failure-Type Sensitivity</b>	
492	Different instruction-following failure modes also	
493	exhibit varying sensitivity to attention-based repair.	
494	Failures related to language selection and format-	
495	ting constraints are more consistently corrected by	
496	moderate stimulation, while degeneration-related	
497	failures are less reliably repaired. This pattern	
498	suggests that DAS primarily restores instruction-	
499	conditioning mechanisms rather than general lan-	
500	guage fluency, which may depend on broader	
501	model dynamics beyond attention head perturba-	
502	tions.	
503	<b>5.6 Diagnostic Interpretation of Attention</b>	
504	<b>Stimulation</b>	
505	Taken together, these analyses indicate that DAS	
506	functions as a diagnostic and corrective interven-	
507	tion rather than a general performance enhance-	
508	ment. The fact that excessive stimulation degrades	
509	performance, even in cases that were correct be-	
510	fore repair, argues against interpreting DAS as uni-	
511	formly beneficial. Instead, the observed recov-	
512	ery supports the view that quantization-induced	
513	instruction-following failures arise from local-	
514	ized disruptions that can be selectively corrected	
515	through targeted attention-based intervention.	
	<b>6 Conclusion</b>	516
	We find that quantization induced instruction fol-	517
	lowing failures are highly localized, arising from	518
	disruption in a small subset of attention heads. To	519
	address this, we introduce Deep Attention Stimula-	520
	tion (DAS), a training free intervention that restores	521
	instruction following by compensating degraded	522
	attention head outputs. DAS recovers language and	523
	formatting constraints under moderate stimulation,	524
	while excessive stimulation destabilizes generation.	525
	These results point to a new direction for targeted	526
	and interpretable post quantization model repair.	527
	<b>7 Limitations</b>	528
	<b>Evaluation Scope.</b> Our study focuses on a sin-	529
	gle instruction-tuned model (Qwen2.5-7B-Instruct)	530
	quantized with GPTQ-4bit, and evaluates a limited	531
	set of 10 representative failure cases selected for	532
	qualitative analysis. While this setting enables de-	533
	tailed mechanistic inspection, broader evaluation	534
	across model families, quantization methods, and	535
	the full IFEval benchmark could further validate	536
	the generality of our approach.	537
	<b>Manual Parameter Selection.</b> The stimulation	538
	strength $\alpha$ is selected via grid search, with $\alpha = 5$	539
	yielding the most consistent recovery in our exper-	540
	iments. This value may not transfer across mod-	541
	els or quantization settings. Developing princi-	542
	pled or adaptive strategies for selecting stimulation	543
	strength remains an open direction.	544
	<b>Dependence on Full-Precision Activations.</b> Our	545
	diagnostic procedure relies on access to full-	546
	precision activations to identify degraded attention	547
	heads and construct corrective signals. In practical	548
	deployment scenarios where such access may be	549
	unavailable, pre-computed correction profiles or	550
	self-supervised approximations would be required.	551
	<b>Single-Capability Focus.</b> We focus exclusively	552
	on instruction-following behavior. Whether sim-	553
	ilar localized interventions can recover other ca-	554
	pabilities, such as reasoning, factual consistency,	555
	or safety alignment, remains unexplored. Extend-	556
	ing DAS to multi-capability settings may require	557
	resolving interactions between attention heads sup-	558
	porting different behaviors.	559
	<b>References</b>	560
	Yejin Bang. 2023. A multitask, multilingual, mul-	561
	timodal evaluation of chatgpt on reasoning, hal-	562

563	lucination, and interactivity. <i>arXiv preprint arXiv:2302.04023</i> .	Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. 2008. Representational similarity analysis-connecting the branches of systems neuroscience. <i>Frontiers in systems neuroscience</i> , 2:249.	616
564			617
565	György Buzsáki. 2006. <i>Rhythms of the Brain</i> . Oxford University Press.		618
566			619
567	Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. 2019. What does bert look at? an analysis of bert’s attention. <i>arXiv preprint arXiv:1906.04341</i> .	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023a. Alpacaeval: An automatic evaluator of instruction-following models.	620
568			621
569			622
570			623
571	Peter Dayan and L. F. Abbott. 2001. <i>Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems</i> . MIT Press.	Yu Li and 1 others. 2023b. Evaluating instruction following in language models. In <i>Advances in Neural Information Processing Systems</i> .	624
572			625
573			626
574	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023a. Qlora: Efficient finetuning of quantized llms. <i>Advances in neural information processing systems</i> , 36:10088–10115.	Ji Lin and 1 others. 2023. Awq: Activation-aware weight quantization for llm compression. In <i>International Conference on Machine Learning</i> .	627
575			628
576			629
577			
578	Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. 2023b. Spqr: A sparse-quantized representation for near-lossless llm weight compression. <i>arXiv preprint arXiv:2306.03078</i> .	Stephanie Lin and 1 others. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the Association for Computational Linguistics</i> .	630
579			631
580			632
581			633
582			
583			
584	Tim Dettmers and 1 others. 2022. Llm.int8(): 8-bit matrix multiplication for transformers at scale. In <i>Advances in Neural Information Processing Systems</i> .	Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Andrew Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, and 1 others. 2024a. Align-bench: Benchmarking chinese alignment of large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11621–11640.	634
585			635
586			636
587			637
588			638
589			639
590			640
591			641
592			
593	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. <a href="#">The language model evaluation harness</a> .	Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. 2024b. Llm-qat: Data-free quantization aware training for large language models. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 467–484.	642
594			643
595			644
596			645
597			646
598			647
599			648
600			
601	Mor Geva, Roe Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. <i>arXiv preprint arXiv:2012.14913</i> .	Carlo Miniussi, Jonathan A Harris, and Manuela Ruzzoli. 2013. Transcranial electrical stimulation for cognitive enhancement. <i>Trends in Cognitive Sciences</i> .	649
602			650
603			651
604			652
605			
606			
607			
608			
609			
610	Artyom Kharinaev, Viktor Moskvoretskii, Egor Shvetsov, Kseniia Studenikina, Bykov Mikhail, and Evgeny Burnaev. 2025. Investigating the impact of quantization methods on the safety and reliability of large language models. <i>arXiv preprint arXiv:2502.15799</i> .	Neel Nanda and 1 others. 2023. Progress measures for grokking via mechanistic interpretability. In <i>International Conference on Learning Representations</i> .	653
611			654
612			655
613			
614			
615			
		Chris Olah and 1 others. 2020. Zoom in: An introduction to circuits. <i>Distill</i> .	656
			657
		Long Ouyang and 1 others. 2022. Training language models to follow instructions with human feedback. In <i>Advances in Neural Information Processing Systems</i> .	658
			659
			660
			661
		Ruiyang Qin, Dancheng Liu, Chenhui Xu, Zheyu Yan, Zhaoxuan Tan, Zhengge Jia, Amir Nassereldine, Jijie Li, Meng Jiang, Ahmed Abbasi, and 1 others. 2025. Empirical guidelines for deploying llms onto resource-constrained edge devices. <i>ACM Transactions on Design Automation of Electronic Systems</i> , 30(5):1–58.	662
			663
			664
			665
			666
			667
			668

669 Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, and Ziyu Yao. 2024. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646*.

673 Elena Voita and 1 others. 2019. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of ACL*.

676 Pengcheng Wang and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

678 Jason Wei and 1 others. 2022. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

681 Mengzhou Xia, Mikel Artetxe, Chunting Zhou, Xi Victoria Lin, Ramakanth Pasunuru, Danqi Chen, Luke Zettlemoyer, and Veselin Stoyanov. 2023. **Training trajectories of language models across scales**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13711–13738, Toronto, Canada. Association for Computational Linguistics.

689 Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International conference on machine learning*, pages 38087–38099. PMLR.

694 Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He. 2022. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. *Advances in neural information processing systems*, 35:27168–27183.

700 Zhihan Zhang, Shiyang Li, Zixuan Zhang, Xin Liu, Haoming Jiang, Xianfeng Tang, Yifan Gao, Zheng Li, Haodong Wang, Zhaoxuan Tan, and 1 others. 2025. Iheval: Evaluating language models on following the instruction hierarchy. *arXiv preprint arXiv:2502.08745*.

706 Jeffrey Zhou and 1 others. 2023. Instruction-following evaluation for large language models. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.

## 710 A Additional Experimental Details

### 711 A.1 Model and Quantization Configuration

712 We use Qwen2.5-7B-Instruct with the following  
713 specifications:

- 714 • Architecture: 28 layers, 28 attention heads per  
715 layer (784 heads total)
- 716 • Hidden dimension: 4096, head dimension:  
717 128
- 718 • Quantization: GPTQ 4-bit with group size  
719 128
- 720 • Calibration: 128 samples from C4 dataset

**Decoding Parameters.** All experiments use greedy decoding with fixed parameters to ensure controlled comparison:

- Temperature: 1.0 (greedy)
- Top-p/Top-k: disabled
- Max generation length: 512 tokens
- Repetition penalty: 1.0 (disabled)

## A.2 Sample Selection Procedure

**IFEval Benchmark Execution.** We evaluate both the full-precision (FP16) and GPTQ-4bit quantized models on the complete IFEval benchmark using lm-evaluation-harness (Gao et al., 2024). This produces per-sample results containing:

- Original prompt with explicit instruction constraints
- Model-generated response
- Constraint-level compliance scores (strict and loose)
- Violated constraint types (e.g., language, format, punctuation)

**Quantization Method Selection.** Table 2 compares instruction-following performance across quantization levels. The 2-bit model produces structurally degenerate outputs unsuitable for analysis. The 8-bit model slightly outperforms FP16, likely due to regularization effects, leaving minimal failure cases for study. We select 4-bit quantization as it exhibits measurable but not catastrophic degradation, providing a practical regime where instruction-following failures emerge without complete model collapse.

Quantization	IFEval Avg-4	$\Delta$ from FP16
2-bit GPTQ	0.1353	-0.6398
4-bit GPTQ	0.7566	-0.0185
8-bit GPTQ	0.7821	+0.0070
FP16 (baseline)	0.7751	0.0000

Table 2: Instruction-following performance across quantization levels on IFEval benchmark. 4-bit quantization provides a balanced regime with measurable degradation suitable for mechanistic analysis.

752	<b>Failure Case Identification.</b> From the complete	• <b>Degenerative repetition</b> (1 sample, ~10%):	799
753	benchmark results, we identify failure cases where:	Repeated token sequences preventing coherent	800
754	1. FP16 model satisfies all verifiable constraints	output.	801
755	(pass)		
756	2. GPTQ-4bit model violates at least one constraint	This distribution mirrors the failure pattern observed	802
757	(fail)	across the complete candidate set, ensuring	803
758	3. Prompt length is between 30-1200 characters	our detailed analysis on 10 samples reflects the typical	804
759	4. Instructions contain explicit verifiable constraints	quantization-induced degradation modes rather	805
760	(language, format, casing, punctuation, structure)	than edge cases. The selected samples (IFEval IDs:	806
761		12, 13, 57, 74, 79, 97, 102, 206, 1075, 1137) enable	807
762	This filtering yields a candidate set of	controlled evaluation where each sample has verified	808
763	prompts where quantization demonstrably degrades	ground truth (FP16 passes, GPTQ fails), and $\Delta\mathbf{O}^{(\ell,h)}$	809
764	instruction-following capability under identical	can be computed exactly by comparing	810
765	decoding conditions.	the two models on the same prompt.	811
766	<b>Evaluation Protocol and Validation.</b> Given the	<b>A.3 Hyperparameter Selection</b>	812
767	limited pool of clear failure cases and the subtlety	<b>Number of Critical Heads (<math>k</math>).</b> We set $k = 15$	813
768	of instruction violations, we employ a rigorous	based on the distribution of $ \Delta\alpha^{(\ell,h)} $ values. Ap-	814
769	manual evaluation protocol. Each candidate sample	proximately 15 heads (< 2% of 784 total) exhibit de-	815
770	is independently assessed by three large language	viations exceeding 0.005, forming a natural thresh-	816
771	models (GPT-4, Claude-3.5-Sonnet, Gemini-1.5-	old separating severely degraded heads from the	817
772	Pro) to verify constraint violations and ensure inter-	majority showing minimal changes.	818
773	rater reliability. Samples are only included if all	<b>Stimulation Strength (<math>\alpha</math>).</b> We evaluate three set-	819
774	three evaluators agree on: (1) FP16 satisfies all	tings:	820
775	constraints, and (2) GPTQ-4bit violates at least	• $\alpha = 0.0$ : Baseline (no intervention)	821
776	one constraint. This conservative approach ensures	• $\alpha = 5.0$ : Moderate stimulation (selected via	822
777	high-confidence failure cases but further limits the	grid search over $\{1, 3, 5, 7, 10\}$ on 5 held-out	823
778	available sample pool.	failure cases)	824
779	<b>Representative Sample Selection.</b> From the val-	• $\alpha = 10.0$ : Excessive stimulation (to demon-	825
780	idated failure cases, we manually select 10 repre-	strate over-intervention effects)	826
781	sentative samples to cover diverse failure modes.	<b>B Implementation Details</b>	827
782	To ensure these samples capture the broader failure	<b>B.1 Instruction Token Identification</b>	828
783	distribution, we analyze all candidate failure cases	We identify instruction-bearing tokens using regex-	829
784	and categorize them by violation type. Our selected	based keyword matching on the raw prompt text.	830
785	samples proportionally represent the major failure	Example patterns include: "two words", "start	831
786	categories observed in the full candidate set:	with", "json", "format", "only output.*",	832
787	• <b>Language switching</b> (3 samples, ~30%):	"lowercase", etc.	833
788	Model responds in unintended language (e.g.,	For each matched character span in the text, we	834
789	German, French) despite explicit English in-	map positions to token indices using the tokenizer's	835
790	struction. This is the most common failure	offset mapping facility, which returns character-	836
791	mode in the candidate set.	level boundaries for each token. This yields a	837
792	• <b>Format violations</b> (4 samples, ~40%): Miss-	set $\mathcal{I}(x) \subset \{1, \dots, T\}$ representing instruction-	838
793	ing bullet points, paragraph breaks, or struc-	bearing token positions.	839
794	tural requirements. Represents the largest cat-	<b>B.2 Attention Weight Extraction</b>	840
795	egory of quantization-induced errors.	To extract full attention matrices with	841
796	• <b>Constraint omission</b> (2 samples, ~20%):	output_attentions=True, we disable Flash	842
797	Failure to follow specific lexical constraints	Attention and SDPA optimizations, forcing the	843
798	(forbidden words, casing).		

844 model to use eager attention implementation. This  
 845 ensures complete [batch, heads,  $T$ ,  $T$ ] attention  
 846 tensors are returned for each layer, where  $T$  is the  
 847 sequence length.

### 848 B.3 Computing Instruction Attention Score

849 Given attention matrix  $\mathbf{A}^{(\ell,h)} \in \mathbb{R}^{T \times T}$  and instruc-  
 850 tion token indices  $\mathcal{I}(x)$ , we compute:

$$851 \quad a^{(\ell,h)}(x) = \frac{1}{T \cdot |\mathcal{I}(x)|} \sum_{t=1}^T \sum_{i \in \mathcal{I}(x)} \mathbf{A}_{t,i}^{(\ell,h)} \quad (6)$$

852 This represents the average attention weight  
 853 from all query positions to instruction tokens. We  
 854 aggregate across failure cases:

$$855 \quad \Delta a^{(\ell,h)} = \frac{1}{|F|} \sum_{x \in F} \left( a_{\text{FP}}^{(\ell,h)}(x) - a_{\text{Q}}^{(\ell,h)}(x) \right) \quad (7)$$

856 Heads are ranked by  $|\Delta a^{(\ell,h)}|$  and the top- $k$  form  
 857  $H^*$ .

### 858 B.4 Output Vector Compensation

859 **Per-Sample Correction.** For each diagnostic  
 860 prompt  $x \in F$ , we cache the output difference:

$$861 \quad \Delta \mathbf{O}^{(\ell,h)}(x) = \mathbf{O}_{\text{FP}}^{(\ell,h)}(x) - \mathbf{O}_{\text{Q}}^{(\ell,h)}(x) \in \mathbb{R}^{T_x \times d_h} \quad (8)$$

862 where  $T_x$  is the sequence length of prompt  $x$  and  
 863  $d_h = 128$  is the head dimension.

864 **Inference-Time Application.** We evaluate on the  
 865 *same* prompts used for diagnosis, ensuring exact  
 866 sequence length matching and isolating the effect  
 867 of attention-level intervention:

$$868 \quad \tilde{\mathbf{O}}^{(\ell,h)}(x) = \mathbf{O}_{\text{Q}}^{(\ell,h)}(x) + \alpha \cdot \Delta \mathbf{O}^{(\ell,h)}(x) \quad (9)$$

869 This is implemented via PyTorch forward hooks  
 870 registered on attention head modules. During the  
 871 forward pass, when the quantized model computes  
 872  $\mathbf{O}_{\text{Q}}^{(\ell,h)}(x)$ , the hook automatically applies the cor-  
 873 rection before passing the output to subsequent  
 874 layers. Model weights remain unchanged through-  
 875 out.

**Generalization to New Prompts.** The current  
 implementation requires test prompts to match di-  
 agnostic prompts for exact sequence length corre-  
 spondence. For arbitrary prompts, one could com-  
 pute position-averaged corrections:

$$\Delta \bar{\mathbf{O}}^{(\ell,h)} = \frac{1}{|F|} \sum_{x \in F} \frac{1}{T_x} \sum_{t=1}^{T_x} \Delta \mathbf{O}^{(\ell,h)}(x)[t, :] \in \mathbb{R}^{d_h} \quad (10)$$

which can broadcast to any sequence length. We  
 leave this extension to future work.

## 884 C Additional Results

### 885 C.1 Complete Qualitative Comparison

886 Table 3 presents complete outputs for all 10 evalua-  
 887 tion samples under three stimulation conditions.

### 888 C.2 Top Degraded Attention Heads

889 Table 4 lists the 15 most affected heads identified  
 890 by DAS diagnostic procedure:

### 891 C.3 Attention Deviation Statistics

892 Across all 784 attention heads, the distribution of  
 893 instruction attention deviations is approximately  
 894 symmetric:

- 895 • 348 heads (44.4%) show positive  $\Delta a^{(\ell,h)}$  (in-  
 896 creased attention to instruction tokens)
- 897 • 433 heads (55.2%) show negative  $\Delta a^{(\ell,h)}$  (de-  
 898 creased attention to instruction tokens)
- 899 • Only 15 heads (1.9%) exhibit  $|\Delta a^{(\ell,h)}| >$   
 900 0.005
- 901 • Mean absolute deviation: 0.00118
- 902 • Standard deviation: 0.00201

903 This symmetric distribution with a small sub-  
 904 set of severe outliers supports our hypothesis that  
 905 quantization disrupts a sparse set of critical heads  
 906 rather than causing uniform degradation across the  
 907 model. The localized nature of degradation enables  
 908 targeted repair through selective head compensa-  
 909 tion.

## 910 D Computational Resources

911 All experiments were conducted on the University  
 912 of Notre Dame Center for Research Computing  
 913 (CRC) infrastructure:

ID	GPTQ ( $\alpha = 0$ )	DAS ( $\alpha = 5$ )	DAS ( $\alpha = 10$ )
12	Degenerative repetition	Correct format & content	Severe repetition
13	Missing paragraph breaks	Correct structure	Collapsed output
57	Wrong casing & format	Correct (lowercase + bullets)	Unstable symbols
74	Violates rhyme & format	Correct limerick	Partial corruption
79	Constraint omission	Fully recovered	Meta commentary
97	<b>German output</b>	<b>Correct English</b>	Off-task generation
102	Missing bullet structure	Correct bullets	Format instability
206	Correct (passed)	Correct (maintained)	Degraded (failed)
1075	Non-JSON format	Valid JSON output	JSON + extra text
1137	French output	Correct structure	Language mixing

Table 3: Complete qualitative comparison across all 10 evaluation samples. Moderate stimulation ( $\alpha = 5$ ) recovers 7/8 baseline failures while maintaining 2/2 baseline passes. Excessive stimulation ( $\alpha = 10$ ) degrades all outputs, demonstrating the non-monotonic response to intervention strength.

Layer	Head	$\Delta a^{(\ell,h)}$	Type
21	5	-0.0121	Decrease
21	0	-0.0063	Decrease
19	1	-0.0057	Decrease
14	11	+0.0052	Increase
13	4	-0.0049	Decrease
23	8	+0.0047	Increase
22	15	-0.0045	Decrease
20	3	+0.0044	Increase
18	7	-0.0042	Decrease
24	12	+0.0041	Increase
25	2	-0.0039	Decrease
21	19	+0.0038	Increase
27	6	-0.0037	Decrease
26	10	+0.0036	Increase
19	14	-0.0035	Decrease

Table 4: Top 15 attention heads ranked by  $|\Delta a^{(\ell,h)}|$ . Negative values indicate reduced attention to instruction tokens after quantization; positive values indicate increased attention. Degraded heads concentrate in layers 19-27, consistent with prior work showing upper layers encode high-level semantic constraints.

## Hardware.

- GPU: 1× NVIDIA A10 (24GB GDDR6)
- CPU: 32-core Intel Xeon
- RAM: 256GB DDR4

## Runtime.

- IFEval benchmark evaluation (2 models):  $\sim 2$  hours
- Diagnostic phase (attention extraction, 10 prompts):  $\sim 15$  minutes
- Inference with DAS (10 prompts  $\times$  3  $\alpha$  values):  $\sim 5$  minutes
- Total experimental time:  $< 3$  hours

## Software Environment.

- PyTorch 2.1.0 with CUDA 11.8
- Transformers 4.41.0
- Auto-GPTQ 0.7.1
- lm-evaluation-harness 0.4.2
- Python 3.10

## E Limitations and Future Directions

### E.1 Current Limitations

**Prompt-Specific Correction.** Our implementation caches  $\Delta \mathbf{O}^{(\ell,h)}(x)$  separately for each diagnostic prompt, requiring test prompts to match the diagnostic set for sequence length compatibility. This controlled setup enables clean before-after comparison but limits direct generalization to arbitrary new prompts without additional modifications (e.g., position-averaging).

**Manual Hyperparameter Selection.** The stimulation strength  $\alpha = 5$  is selected via grid search on held-out samples and may not transfer optimally across different models, quantization methods, or failure types. Developing adaptive or principled selection strategies remains an open problem.

**Single Model and Quantization Method.** Our analysis focuses on Qwen2.5-7B with GPTQ-4bit quantization. Whether similar localized disruption patterns occur in other model families (LLaMA, Mistral, Gemma) or quantization methods (AWQ, SmoothQuant, mixed-precision) requires systematic investigation.

955 **Instruction-Following Only.** We focus exclu-  
956 sively on instruction-following capability. Whether  
957 DAS can recover other quantization-degraded ca-  
958 pabilities (mathematical reasoning, factual consis-  
959 tency, safety alignment) through attention-level in-  
960 tervention is unexplored.

961 **Evaluation Scale.** Our detailed analysis exam-  
962 ines 10 representative samples to enable thorough  
963 qualitative assessment. While this reveals clear re-  
964 covery patterns, large-scale evaluation across the  
965 full benchmark would strengthen generalizability  
966 claims.

## 967 E.2 Future Directions

968 **Position-Agnostic Correction.** Developing  
969 position-averaged or token-type-specific correc-  
970 tions would enable application to arbitrary prompts  
971 while maintaining the training-free property.

972 **Automatic Stimulation Strength Selection.** In-  
973 vestigating lightweight diagnostic signals (e.g., at-  
974 tention entropy, output variance) to automatically  
975 determine appropriate  $\alpha$  values per head or per  
976 sample.

977 **Multi-Capability Joint Repair.** Extending DAS  
978 to simultaneously recover multiple capabilities by  
979 identifying capability-specific attention heads and  
980 applying orthogonal corrections.

981 **Integration with Quantization Pipeline.** In-  
982 corporating instruction-following diagnostics into  
983 quantization calibration (e.g., mixed-precision  
984 quantization) to preemptively protect critical heads  
985 during compression.

986 **Theoretical Understanding.** Developing formal  
987 analysis of why instruction-following exhibits lo-  
988 calized degradation under quantization compared  
989 to other capabilities, potentially informing more  
990 robust quantization-aware training objectives.