# Enhancing Instruction-Following Capabilities in Seq2Seq Models: DoLa Adaptations for T5

**Huey Sun**[*]       **Anabel Yong**[*]       **Lorenzo Gilly**[*]       **Felipe Jin**

**University College London**

## Abstract

Encoder–decoder models such as FLAN-T5 are finetuned to follow instructions, but often fail when the instructions conflict with memorized continuations ingrained during training. To understand this behavior, we adapt DoLa to FLAN-T5 and examine how representations evolve in the decoder. Our findings show that T5's intermediate layers undergo rapid shifts driven by cross-attention to the encoder. When projected through the language modeling head, each depth presents highly volatile token preferences, leading to unreliable behavior with contrastive decoding. Motivated by this, we introduce a gradient-based activation-steering method that injects an "instruction-compliance" direction into mid-decoder layers, where the representation is both meaningful and still malleable. This intervention dramatically improves MemoTrap performance (52% → 99.7%), demonstrating that mechanistic steering can succeed where contrastive decoding fails in Seq2Seq architectures.

## 1 Introduction

Large language models (LLMs) are increasingly expected to serve as general-purpose instruction followers and perform diverse tasks. While recent instruction-tuned models have demonstrated strong zero-shot performance across a wide range of benchmarks [1, 2], even highly tuned models often struggle to be faithful to complex prompts. These failures are most visible when an instruction directly conflicts with their language priors, as models are prone to recite memorized training data over following explicit user instructions [3].

Recent work has explored inference-time decoding strategies to improve factuality and reduce hallucinations. One such example is **D**ecoding by **C**ontrastive **La**yers (**DoLa**) [4], a decoding strategy that contrasts a model's final layer logits with those from intermediate layers to amplify "mature" predictions and suppress unstable or spurious ones. As each

layer of decoder-only models incrementally refine an autoregressive hidden state, intermediate logits serve as meaningful indicators of the model's evolving beliefs. By exploiting the modular nature of knowledge encoding [5, 6], DoLa improves factuality in the LLaMA [7] family of models, and suggests that the "knowledge neurons" in the upper layers are better expressed when lower layers serve as a contrasting baseline.

However, it remains unclear whether such methods extend to encoder–decoder architectures, where each decoder layer continuously integrates information from a fully processed encoder representation. In models such as T5 [2] and FLAN-T5 [8], intermediate decoder layers are not trained to make next-token predictions, and their hidden states may not correspond to coherent partial distributions. This raises open questions about how instruction-related signals propagate through Seq2Seq models. Although T5-style architectures are widely used for instruction following, little is known about their internal dynamics during constrained generation. Prior mechanistic analyses have largely focused on decoder-only models, characterizing layer-wise emergence of factuality, syntactic structure, or reasoning behaviors [9, 10, 11]. Comparable studies in Seq2Seq models are scarce, leaving open why faithfulness failures persist and what forms of intervention are most effective.

In this work, we investigate instruction-following behavior through the lens of layer-wise representational dynamics in the FLAN-T5 family. We adapt DoLa to the encoder–decoder setting, use it as a diagnostic tool to probe how decoder representations evolve with depth, and develop a targeted activation-steering method informed by this analysis.

## 2 Related Work

### 2.1 Instruction Following in Language Models

Instruction following enables LLMs to behave as controllable systems that can perform user-specified tasks, adapt to new constraints, and generalize beyond their training distribution. This can be thought of as a form of zero-shot generalization, where the model
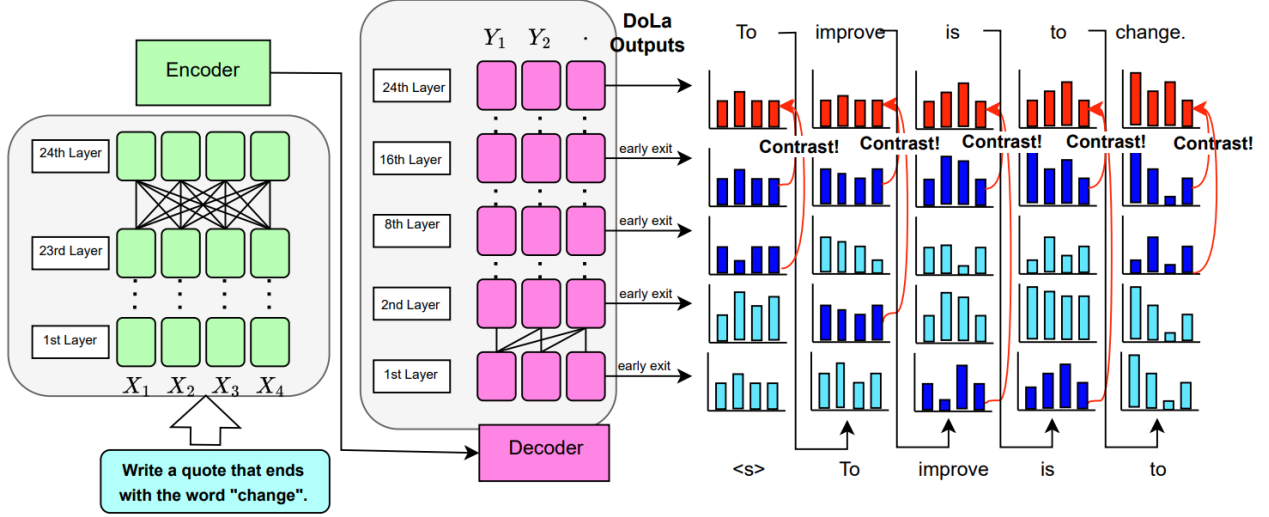
---

[*] Equal contribution.

Figure 1: Schematic representation of how dynamic premature layer selection (DoLa) works with T5 architectures. This illustration is specifically for T5-Large, which has 24 encoder and decoder layers. The different model sizes have different early-exit layers, as explained in Appendix 7.1.

must use latent knowledge to satisfy constraints it has never been explicitly trained on. While early work showed that pretrained Seq2Seq models could follow templated instructions for tasks such as classification and translation [12], the FLAN framework [13] expanded this by finetuning models on diverse instruction formats directly, substantially improving their ability to generalize to unseen tasks and phrasing.

Despite these advances, recent evaluations reveal that instruction adherence remains fragile under conflicting semantic priors. Even large models will imitate undesirable patterns in training data, get distracted, or be easily misled at easy tasks [3], indicating that instruction signals may not uniformly dominate the model's internal activations.

## 2.2 Contrastive Decoding Methods

Contrastive decoding aims to improve generative faithfulness by comparing two predictive distributions and amplifying their differences. Early approaches contrasted two separate models, treating one as an "expert" and the other as an "amateur" [14], while later methods contrasted two internal states of a single model, such as predictions with and without context history [15]. Variants of this idea introduce plausibility constraints [16] or other heuristics to guide the contrast.

DoLa [4] follows this line of work by contrasting the final-layer logits with those from an earlier "premature" layer. This layer is dynamically selected using Jensen—Shannon divergence, and the con-

trastive adjustment is applied to emphasize tokens whose probability increases from the early to the late distribution. This approach has proven effective in decoder-only models and provides a natural starting point for analyzing layer-wise signal development.

## 2.3 Mechanistic Interpretability

Mechanistic interpretability seeks to characterize how specific behaviors arise from internal model components. Recent work shows that models often encode behaviors along identifiable directions in activation space [17], and that manipulating these directions can reliably shift model outputs. Activation-steering methods, including representation engineering [18] and steering vectors [19], demonstrate that targeted interventions can modulate reasoning patterns, stylistic preferences, or safety-related behaviors without changing model weights.

Our approach follows this paradigm but is motivated by a different goal: understanding and improving instruction faithfulness. Instead of using generic steering directions, we mine a contrastive direction that separates instruction-following behavior from memorization-driven behavior. This enables a targeted intervention makes instruction-relevant representations easier for the model to express.

## 2.4 Architectural Differences: Decoder-Only vs. Seq2Seq

Decoder-only models such as LLaMA and GPT operate entirely on left-context information; once decoding begins, no new information enters the network,

and intermediate activations remain directly predictive of the next token [7]. As a result, layer-wise probing and early exiting yield coherent partial predictions.

In Seq2Seq architectures like T5, the decoder attends to a fully contextualized encoder representation at every layer [2]. This repeated incorporation of encoder features can reshape decoder activations in non-monotonic ways, making intermediate states less directly predictive of next tokens compared to decoder-only models. As a result, the notion of 'early' or 'immature' predictions is less clear in T5-style decoders, raising questions about whether contrastive early–late decoding behaves differently in these architectures and motivating a closer examination of how instruction-relevant signals evolve across depth.

## 3 Methodology

Our goal is to analyze and improve instruction-following behavior in encoder–decoder architectures. First, we adapt DoLa to T5 to study how instruction signals evolve across decoder layers. Then, we develop a gradient-based activation steering method that actively suppresses memorization circuits.

### 3.1 Adapting DoLa to T5

In DoLa, we need to identify the premature layer with the predictive distribution that differs the most from the final next-token logits. Let $q_N(\cdot)$ denote the final-layer distribution and $q_j(\cdot)$ the distribution obtained by passing the hidden state at decoder layer $j$ through the language modeling head. We select the premature layer $M$:

$$M = \arg\max_{j \in \mathcal{J}} \text{JSD}(q_N, q_j), \quad (1)$$

where JSD is the Jensen–Shannon divergence. We then modify the logits $\ell$ by amplifying the mature–premature difference:

$$\tilde{\ell} = \ell_N + \lambda(\ell_N - \ell_M), \quad (2)$$

where $\lambda$ controls the contrastive strength. As in the original implementation by Chuang et al. [4], we apply a repetition penalty of 1.2 to mitigate degenerate looping [20].

**Application to T5.** Unlike decoder-only models, T5's decoder layers are not trained to make next-token predictions, and their hidden states may not be in a logit-aligned subspace. To adapt DoLa, we follow Chuang et al. [4] and apply the shared LM head to each decoder layer's residual stream, treating these projections as approximate next-token distributions.
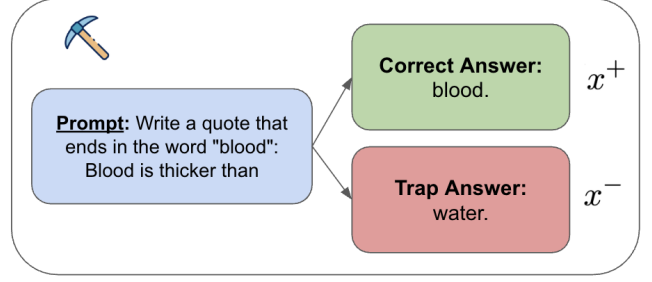


Figure 2: **Contrastive Setup** A prompt from the MemoTrap dataset [21] that sets up a contrast between an instruction and a common memorized trap.

Because the encoder does not participate in autoregressive generation, we restrict the candidate set $\mathcal{J}$ to decoder layers (Figure 1).

While this projection-based approach does not guarantee that intermediate T5 activations form coherent partial distributions, it enables a consistent layer-wise comparison within the contrastive decoding framework and serves as the basis for our diagnostic analysis in Section 4.1.

### 3.2 Gradient-Based Activation Steering

While DoLa contrasts predictions across layers, it cannot directly reshape the model's internal representations when the layers express the same bias. To enable targeted control, we develop a gradient-based activation steering method that extracts a direction separating instruction-compliant and memorized behaviors and injects it into the decoder at inference time.

#### 3.2.1 Contrastive Setup

We consider cases where generation involves a binary conflict between an instruction-specified outcome and a competing continuation (Figure 2). For each example, we define (i) a **target** token $x^+$ that satisfies the instruction, and (ii) a **competing** token $x^-$ that reflects the model's default or memorized preference.

To extract a direction, we define a contrastive loss that increases the model's preference for $x^-$ relative to $x^+$:

$$\mathcal{L} = \log P(x^-) - \log P(x^+). \quad (3)$$

For a chosen decoder layer $l$, we compute the gradient of this loss with respect to its hidden state $h_l$:

$$g_l = \nabla_{h_l} \mathcal{L}. \quad (4)$$

This gradient $g_l$ points in the direction that makes the model more likely to produce the competing token $x^-$. Its negation thus provides a steering direction that shifts the layer's representation toward the instruction-following outcome.
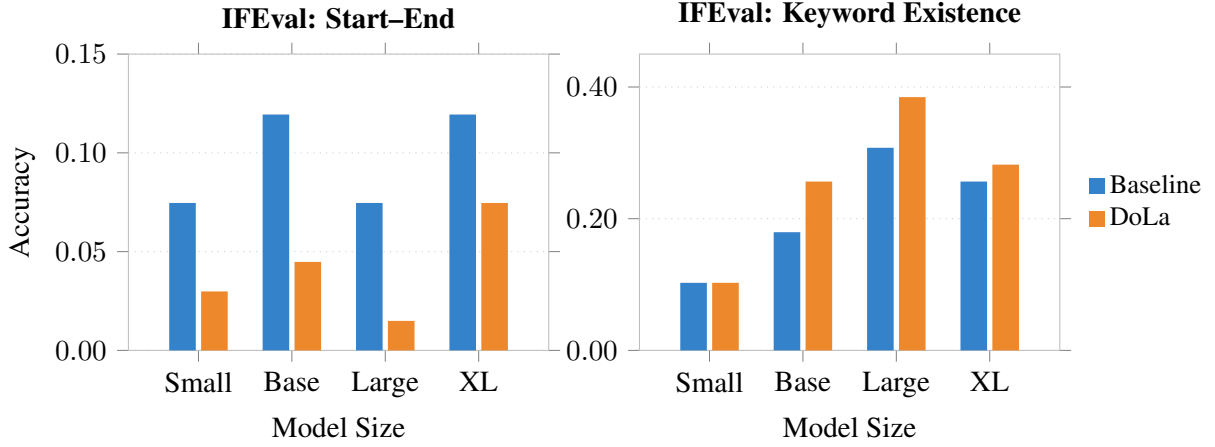
Figure 3: Baseline vs. DoLa accuracy across FLAN-T5 model sizes on two representative IFEval categories.

### 3.2.2 Mining the Steering Vector

We compute gradients across 100 examples and average them to obtain a task-specific steering direction:

$$\vec{v}_{\text{task}} = -\frac{1}{N} \sum_{i=1}^{N} \nabla_{h_l} \mathcal{L}_i. \tag{5}$$

Because each example uses a different instruction and a different target token, token-specific semantics cancel out, leaving a vector that captures the shared mechanism distinguishing trap-driven from instruction-driven behavior.

### 3.2.3 Layer Injection

We mine and inject this vector to the same layer $l$, which we sweep over. During inference, we modify the hidden state $h_l$:

$$h'_l = h_l + \alpha \cdot \frac{\vec{v}_{\text{task}}}{\|\vec{v}_{\text{task}}\|}, \tag{6}$$

where $\alpha$ controls steering strength. We experiment with a variety of values for $\alpha$, and find that $\alpha = 1000$ provides stable improvements without distorting syntax or fluency.

### 3.3 Evaluation Setup

We evaluate FLAN-T5 models (Small, Base, Large, XL) on two benchmarks:

**IFEval** A framework for evaluating instruction-following capability with verifiable criteria [22]. We report the loose prompt and instruction-level accuracy and qualitative evaluation with GPT-4 as a judge [23].

**MemoTrap** A dataset [21] to evaluate whether a model can adhere to explicit instructions or succumb to memorized completions. We use the Proverb Ending subset, reserving 100 instances for vector mining, and evaluating the free generation performance on another subset of 300 examples.

## 4 Results

### 4.1 DoLa and Layer Dynamics

| Model | Baseline (%) | DoLa (%) |
|-------|--------------|----------|
| FLAN-T5-Small | 20.38 | **22.90** |
| FLAN-T5-Base | **26.38** | 23.02 |
| FLAN-T5-Large | **26.86** | 26.74 |
| FLAN-T5-XL | **28.78** | 25.30 |

Table 1: IFEval instruction-level accuracy.

On IFEval, applying DoLa to the FLAN-T5 family leads to modest overall gains for the Small model and neutral or slightly negative effects for larger variants (Table 1). However, a deeper category-level analysis reveals that DoLa benefits specific instruction types across the models and harms others (Figure 3). This heterogeneous behavior provides critical clues about the evolution of T5's intermediate representations.

**Start–End Constraints** Figure 4 shows a representative failure case. When instructed to output text with a specific suffix, the baseline FLAN-T5-XL only outputs the mandated ending, which scores under IFEval's verifiable metric. By contrast, DoLa produces a more natural-sounding blog post, but fails to end with the required phrase.

In this case, while both the intermediate and final layers encode the required ending, the final layers introduce a more developed natural-language response. When DoLa contrasts these layers, it amplifies the components of the final distribution that differ most from the intermediate one, while leaving the shared instruction signal largely unchanged. As a result, the contrastive update strengthens a path that produces fluent but instruction-violating text, illustrating how later layers can drift toward natural-language continuation in ways that run orthogonal to rigid constraint-following.

**Prompt**: My name is Naomi. Write a blog post in my name for the canucks hockey team about why they need to be more mindful about their environments. End the blog post with "Naomi thanks you for reading."

**Base**: Naomi thanks you for reading.

**DoLa**: Naomi, a canucks hockey team member, is very interested in the environment. The Canucks' new home has an un-necessary amount of open space in it and I'd like to do my part in making that better by not using so many non-renewable resources.

Figure 4: Example IFEval Start–End prompt and FLAN-T5-XL outputs with and without DoLa.

**Keyword Constraints** In contrast, keyword-insertion tasks often show the opposite pattern. As illustrated in Figure 5, the required token ("lacking") receives extremely low probability in the early and mid decoder layers, fluctuating widely in rank (e.g., $180 \to 1156 \to 3992$) before finally becoming competitive in the final layer. Because DoLa amplifies the late-emerging differences between the final and premature distributions, it tends to promote the instruction-relevant token once it appears. In this example, DoLa raises "lacking" from the fifth to the top-ranked token in the output distribution, fulfilling the instruction.

**Seq2Seq Layer Dynamics** Across both cases, we find that instruction-following cues do not develop along a single, monotonic trajectory through the decoder. Depending on the task, the instruction-relevant token may persist across multiple depths, emerge only at the final layers, or fluctuate sharply as the decoder integrates cross-attention information. This is reflected in the Jensen–Shannon divergences, which rise dramatically in the later layers as representations consolidate into the final output. DoLa succeeds when its selected premature layer falls at a point where the instruction signal is strengthening, and fails when it contrasts against a layer where the instruction cue is weak or overshadowed by the model's natural-language priors. This task-dependent variability explains DoLa's varying performance in Seq2Seq models and motivates interventions that do not rely on stable intermediate logits.

## 4.2 A Representation-Level Intervention: Gradient-Based Steering

With activation steering, we can directly perturb the decoder's hidden state along a direction associated with instruction compliance. This removes the dependence on layer-wise token distributions and instead operates on the model's internal representations.

**MemoTrap** In MemoTrap, each example contains a paired contrastive target: a correct instruction-

**Model Response (Ongoing):** GANs can be applied to architecture and performance quality metrics such as fault tolerance or throughput, which are

**Next Token with DoLa:** lacking

| Jensen-Shannon Divergences | Layer | Top 3 Tokens with DoLa (Likelihood Rankings in Vocabulary) | | |
|---|---|---|---|---|
| | | *lacking* | *generally* | *both* |
| | DoLa | 1 | 2 | 3 |
| Final (N/A) | Layer 24 | 5 | 12 | 8 |
| **0.0266** | Layer 22 | 180 | 106 | 66 |
| 0.0189 | Layer 20 | 108 | 97 | 43 |
| 0.0162 | Layer 18 | 164 | 111 | 88 |
| 0.0160 | Layer 16 | 1156 | 122 | 108 |
| 0.0174 | Layer 14 | 3350 | 285 | 120 |
| 0.0152 | Layer 12 | 3992 | 210 | 205 |
| 0.0139 | Layer 10 | 4832 | 386 | 168 |
| 0.0112 | Layer 8 | 5585 | 494 | 414 |
| 0.0060 | Layer 6 | 4134 | 485 | 768 |
| 0.0025 | Layer 4 | 852 | 185 | 661 |
| 0.0016 | Layer 2 | 192 | 405 | 1315 |

Figure 5: Layerwise ranking of the instruction-relevant token ("lacking") for FLAN-T5-Large on IFEval Question 154

following output and a memorized trap drawn from the model's pretraining distribution. This enables us to compute a clean contrastive gradient that separates the behaviors, which even large models can struggle to do without intervention [3].

**Steering Results** After mining and averaging gradients from a held-out dataset, we inject a scaled steering vector at each decoder layer in our test split. Table 2 shows that activation steering dramatically improves instruction-following performance across all model sizes, nearly solving the benchmark for FLAN-T5-Large.

| Variant | Baseline | Best Steered | $\Delta$ |
|---|---|---|---|
| Small | 60.3% | 92.3% | +32.0% |
| Base | 65.7% | 98.7% | +33.0% |
| Large | 52.0% | 99.7% | +47.7% |

Table 2: Best MemoTrap activation-steering performance across FLAN-T5 model scales.

**Where Should We Steer?** This intervention requires a pre-specified target layer and application strength. In our sweep, we find that performance varies sharply depending on where the steering vector is injected. Figure 6 shows that steering is highly effective in a narrow mid-depth region, and substantially weaker when applied too early or too late, and depends on the strength of the intervention.
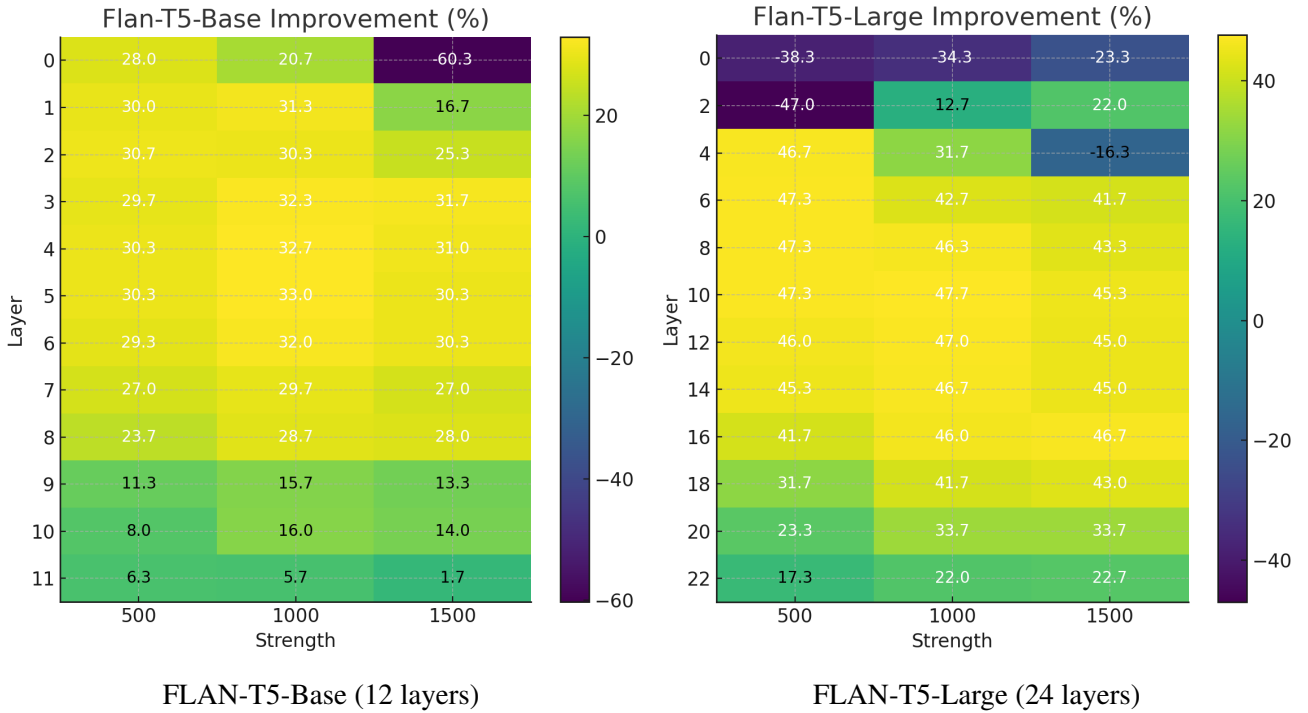
Figure 6: Tracking activation-steering improvements across FLAN-T5 model scales (more in Appendix 7.2).

| Model | Best Layer | # Layers |
|---|---|---|
| FLAN-T5-Small | 2 | 8 |
| FLAN-T5-Base | 5 | 12 |
| FLAN-T5-Large | 10 | 24 |

Table 3: Layer position of peak steering efficacy across different model depths at any strength.

**Interpretation** These results show that injecting activations in early layers has little or degenerating impact on the eventual output, perhaps due to the their role in syntax processing or cross-attention mixing. More surprisingly, steering also proves ineffective in the final layers, even though our token-tracking results show that the output logits remain turbulent. This reflects a difference between *volatile predictions* and *flexible representations*: the model may still shuffle token rankings late in the stack, but the underlying representations have already settled into a narrow path toward a particular response. By contrast, the mid layers contain a more informative but not yet committed representation, allowing the steering direction to meaningfully redirect the model away from memorized completions and toward instruction-following behavior.

## 5 Conclusion

In this work, we investigated how instruction-following signals evolve inside encoder–decoder models and examined why contrastive decoding methods, such as DoLa, exhibit mixed behavior in T5 architectures.

Our analysis revealed that intermediate decoder layers in FLAN-T5 undergo substantial representational shifts driven by repeated cross-attention, making their projected token distributions unstable and difficult to leverage for early–late contrastive decoding. This instability helps explain DoLa's task-dependent effectiveness: depending on where the instruction cue becomes identifiable, contrasting against a premature layer may either strengthen or obscure instruction-relevant behavior.

Motivated by these observations, we introduced a gradient-based activation steering method that acts directly on hidden representations rather than relying on intermediate logits. Steering in mid-decoder layers, where representations are informative but not yet committed, produces dramatic improvements on MemoTrap, raising FLAN-T5-Large accuracy from 52% to 99.7%. These findings demonstrate that mechanistic, representation-level interventions offer a promising alternative to layerwise contrast in Seq2Seq models, and highlight mid-depth activations as an important locus of controllability.

## 6 Limitations and Future Work

Our study focuses on the FLAN-T5 family and on tasks with clear contrastive supervision at the token level. While MemoTrap is useful for isolating memorization–instruction conflicts, it covers only a nar-

row subset of instruction-following behavior and does not reflect more complex settings such as multi-sentence reasoning or dialog. The steering directions we extract are also specific to this task, with unknown generalization performance across other instruction types. Furthermore, our experiments center on single-token interventions, and steering multi-token or sequence-level constraints may require different formulations. Finally, our diagnostics rely on LM-head projections, which provide a useful but incomplete view of the hidden-state geometry in Seq2Seq models.

These limitations suggest several directions for future work. One key direction is developing steering vectors that encode more general behaviors such as following structural constraints or suppressing memorized continuations rather than being tied to specific token contrasts. Investigating encoder-side interventions may clarify how instruction information is formed and propagated before decoding begins. Another promising area is designing automated, metric-driven procedures for selecting steering layers or identifying coherent gradient directions, reducing the need for manual search.

Beyond MemoTrap, evaluating steering across a wider set of instruction types, such as multi-step reasoning, rewriting, and transformation tasks, would help characterize when and how representation-level interventions improve faithfulness. Finally, combining activation steering with decoding-time algorithms may yield hybrid methods that leverage the strengths of both representation-level control and output-level contrastive adjustments.

## References

[1] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=TG8KACxEON.

[2] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2019. URL https://arxiv.org/abs/1910.10683.

[3] Ian R. McKenzie, Alexander Lyzhov, Michael Martin Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Xudong Shen, Joe Cavanagh, Andrew George Gritsevskiy, Derik Kauffman, Aaron T. Kirtland, Zhengping Zhou, Yuhui Zhang, Sicong Huang, Daniel Wurgaft, Max Weiss, Alexis Ross, Gabriel Recchia, Alisa Liu, Jiacheng Liu, Tom Tseng, Tomasz Korbak, Najoung Kim, Samuel R. Bowman, and Ethan Perez. Inverse scaling: When bigger isn't better. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=DwgRm72GQF. Featured Certification.

[4] Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. Dola: Decoding by contrastive layers improves factuality in large language models, 03 2024. URL https://arxiv.org/pdf/2309.03883.pdf.

[5] Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline, 2019. URL https://aclanthology.org/P19-1452.pdf.

[6] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, Furu Wei, and Moe Key. Knowledge neurons in pretrained transformers. 1:8493–8502, 2022. URL https://aclanthology.org/2022.acl-long.581.pdf.

[7] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *arXiv (Cornell University)*, 02 2023. doi: 10.48550/arxiv.2302.13971.

[8] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Pet-

rov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022. URL https://arxiv.org/abs/2210.11416.

[9] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main. 446. URL https://aclanthology.org/2021.emnlp-main.446/.

[10] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems*, 36, 2022. arXiv:2202.05262.

[11] Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. Function vectors in large language models. In *Proceedings of the 2024 International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=AwyxtyMwaG. arXiv:2310.15213.

[12] Zheng Zhang Dan Klein Ruiqi Zhong, Kristy Lee. Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections, 04 2021. URL https://arxiv.org/pdf/2104.04670.pdf.

[13] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *arXiv:2210.11416 [cs]*, 10 2022. URL https://arxiv.org/abs/2210.11416.

[14] Ariel Gera, Roni Friedman, Ofir Arviv, Chulaka Gunasekara, Benjamin Sznajder, Noam Slonim, and Eyal Shnarch. The benefits of bad advice: Autocontrastive decoding across model layers. 1:10406, 1042. URL https://aclanthology.org/2023.acl-long.580.pdf.

[15] Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Yih. Trusting your evidence: Hallucinate less with context-aware decoding. URL https://arxiv.org/pdf/2305.14739.pdf.

[16] Lisa Xiang, Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. URL https://arxiv.org/pdf/2210.15097.pdf.

[17] Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. Function vectors in large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=AwyxtyMwaG.

[18] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down approach to ai transparency, 2025. URL https://arxiv.org/abs/2310.01405.

[19] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering language models with activation engineering, 2024. URL https://arxiv.org/abs/2308.10248.

[20] Jin Xu, Xiaojiang Liu, Jianhao Yan, Deng Cai, Huayang Li, and Jian Li. Learning to break the loop: Analyzing and mitigating repetitions for neural text generation, 10 2022. URL https://arxiv.org/abs/2206.02369.

[21] Alisa Liu and Jiacheng Liu. The memo-trap dataset. https://github.com/liujch1998/memo-trap, 2023. GitHub repository.

[22] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan,

Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL https://arxiv.org/pdf/2311.07911.pdf.

[23] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment. 03 2023. doi: 10.48550/arxiv.2303.16634.

# 7 Appendix

## 7.1 DoLa Candidate Layer Configurations

To utilize DoLa, we must specify the set of "premature" layers used for dynamic contrastive selection. Because the T5 and FLAN-T5 model families vary significantly in depth, the candidate layers differ for each model size[cite: 173].

- **Small Models:** T5-Small has 6 decoder layers, while FLAN-T5-Small has 8 decoder layers[cite: 174, 175].

- **Base Models:** Both T5-Base and FLAN-T5-Base have 12 decoder layers[cite: 174].

- **Large & XL Models:** Both T5 and FLAN-T5 in Large and XL sizes have 24 decoder layers.

To maintain consistency across architectures without introducing complex partitioning schemes, we uniformly selected all **even-indexed layers** as candidates for the early-exit contrast. The specific configurations are listed below:

| Model Size | Candidate Early-Exit Layers |
|---|---|
| T5-Small | 0, 2, 4 |
| FLAN-T5-Small | 0, 2, 4, 6 |
| Base | 0, 2, 4, 6, 8, 10 |
| Large / XL | 0, 2, 4, ..., 20, 22 |

Table 4: Candidate premature layers used for DoLa decoding across different T5 model sizes.

For example, when running DoLa on FLAN-T5-Base, the argument provided was `-early_exit_layers 0,2,4,6,8,10,12`.
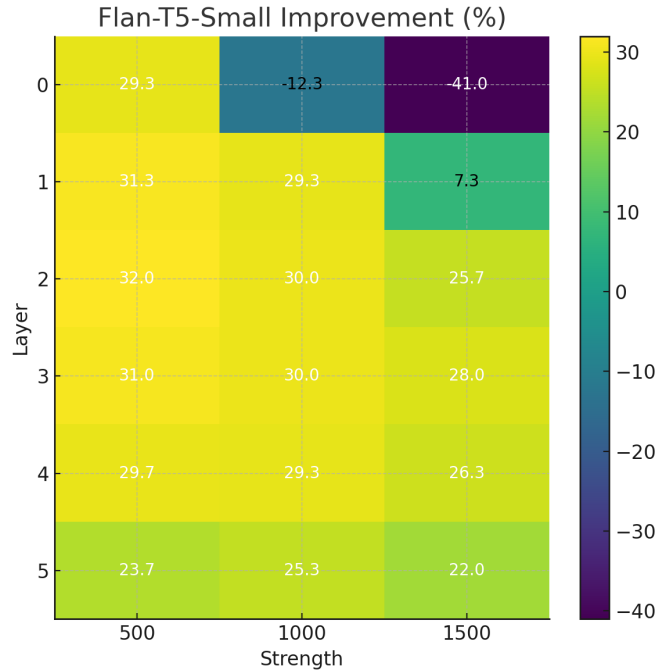
## 7.2 Activation Steering Efficacy in FLAN-T5-Small



Figure 7: Tracking activation-steering efficacy in FLAN-T5-Small.

## 7.3 IFEval Verifiable Instruction Accuracy of FLAN-T5-Small-DoLa

| Instruction Categories | FLAN-T5-Small-DoLa Accuracy (%) |
|---|---|
| change_case | 17.98 |
| combination | 4.615 |
| detectable_content | 20.75 |
| detectable_format | 4.459 |
| keywords | 29.45 |
| language | 29.03 |
| length_constraints | 34.27 |
| punctuation | 69.70 |
| startend | 2.985 |
| change_case:capital_word_frequency | 32.00 |
| change_case:english_capital | 0 |
| change_case:english_lowercase | 20.51 |
| combination:repeat_prompt | 0 |
| combination:two_responses | 12.50 |
| detectable_content:number_placeholders | 14.81 |
| detectable_content:postscript | 26.92 |
| detectable_format:constrained_response | 20.00 |
| detectable_format:json_format | 5.882 |
| detectable_format:multiple_sections | 0 |
| detectable_format:number_bullet_lists | 0 |
| detectable_format:number_highlighted_sections | 8.333 |
| detectable_format:title | 0 |
| keywords:existence | 10.26 |
| keywords:forbidden_words | 46.94 |
| keywords:frequency | 21.43 |
| keywords:letter_frequency | 36.36 |
| language:response_language | 29.03 |
| length_constraints:nth_paragraph_first_word | 0 |
| length_constraints:number_paragraphs | 0 |
| length_constraints:number_sentences | 55.77 |
| length_constraints:number_words | 38.46 |
| punctuation:no_comma | 69.70 |
| startend:end_checker | 7.692 |
| startend:quotation | 0 |

## 7.4 Qualitative Study of DoLa's impact on text generation.

To supplement our quantitative accuracy metrics, we utilized GPT-4 to provide qualitative feedback on response quality, following recent work validating LLMs as evaluators. We used the following prompt to ensure the model focused on instruction adherence and mechanics rather than factual correctness (which is intentionally violated in MemoTrap tasks):

*"Please rate the grammaticality, cohesiveness, and instruction-following performance of these responses, but not factuality. You are not required to verify the factual accuracy of the answers. Each assistant receives an overall score on a scale of 1 to 10, where a higher score indicates better quality."*
[cite: 191]

Table 6: Qualitative study of FLAN-T5-XL with and without DoLa using GPT-4 judgement

| Prompt | **My name is Naomi. Write a blog post in my name for the canucks hockey team about why they need to be more mindful about their environments. End the blog post with "Naomi thanks you for reading." No other words should follow this phrase. This phrase should be the very end of your entire response.** | |
|---|---|---|
| **Method** | **FLAN-T5-XL** | **FLAN-T5-XL + DoLa** |
| | Naomi thanks you for reading. | Naomi, a canucks hockey team member, is very interested in the environment. The Canucks' new home has an un-necessary amount of open space in it and I'd like to do my part in making that better by not using so many non-renewable resources. |
| **Score** | 5.5 | 7.5 |
| **Review** | Assistant 1 delivered a grammatically correct sentence. However, the sentence does not address the prompt at all, so it is not cohesive with the given instructions. | Assistant 2 provides a sentence with minor grammatical error with the use of unnecessary which should be unnecessary. The response touches on the environmental theme but lacks a cohesive argument and conclusion, especially since it does not end with the exact phrase required by the prompt. |

Table 7: Qualitative study of FLAN-T5-Large with and without DoLa using GPT-4 judgement

| Prompt | Write a detailed review of the movie The Social Network. Your entire response should be in English and all lower case (no capital letters whatsoever). | |
|---|---|---|
| **Method** | **FLAN-T5-Large** | **FLAN-T5-Large + DoLa** |
| | The Social Network is a smart, funny, and surprisingly affecting look at the power of social networking. | Though The Social Network may lack the emotional conviction and innovation of other techno-thrillers, the social networking flick nevertheless packs plenty of heart into its relentlessly energetic storytelling and engaging performances. |
| **Score** | 5 | 6 |
| **Review** | Assistant 1 delivers a concise response that effectively touches on the movie's key attributes—its wit, humor, and impact. Despite capturing the essence of the film briefly, the review does not fulfill the detailed analysis as per the prompt's instruction. Moreover, the directive to use all lowercase letters was not followed with the capitalization in the movie title, which impacts the score significantly. | Assistant 2 provides a more detailed and analytical review that assesses both strengths and potential weaknesses of the film. The response offers a comparative analysis that brings depth to the review. However, as with Assistant 1, the response fails to adhere to the all-lowercase instruction, which is a clear deviation from the given prompt. |

Table 8: Qualitative study of FLAN-T5-Large with and without DoLa using GPT-4 judgement

| Prompt | Write a funny haiku about moms, containing keywords "mom" and "mother" in your response. Repeat the request word for word without change, then give your answer (1. do not say any words or characters before repeating the request; 2. the request you need to repeat does not include this sentence) | |
|---|---|---|
| **Method** | **FLAN-T5-Large** | **FLAN-T5-Large + DoLa** |
| | moms are the best mothers i have ever had i love them so much i can't even imagine a life without them | mommy you always look up at me mother you are such an angel mom |
| **Score** | 2 | 3 |
| **Review** | The response provided by Assistant 1 does not conform to the traditional structure of a haiku, which should follow a 5-7-5 syllable pattern. Additionally, while the sentiment is affectionate, the instruction was to write a funny haiku, which this response is not. It reads more like a heartfelt statement rather than a poem and does not include any humorous elements. | Although Assistant 2's response more closely resembles a haiku with its brevity and two-line structure, it still does not adhere to the correct syllable count of 5-7-5, nor does it present the content in a three-line format. Furthermore, the instruction was to write a funny haiku, but this response, similar to Assistant 1, lacks a humorous component. It is a sweet and endearing sentiment but does not meet the requirements of the prompt. |

# 8 Supplementary Material

## 8.1 Documentation of Main Files

Code and documentation can be found here: `https://github.com/itshuey/DoLa/`

## 8.2 Results

- The IFEval FLAN-T5 model outputs and evaluation results can be found here: `https://github.com/itshuey/DoLa/tree/main/results/ifeval-all-results`.

- The MemoTrap FLAN-T5 model outputs and evaluation results can be found here: `https://github.com/itshuey/DoLa/tree/main/results/memo-trap-all-results`

- The logit analysis for Prompt 154 of IFEval with FLAN-T5-Large can be found here: `https://github.com/itshuey/DoLa/tree/main/results/misc`