

Route-Induced Density and Stability (RIDE): Controlled Intervention and Mechanism Analysis of Routing-Style Meta Prompts on LLM Internal States

Anonymous ACL submission

Abstract

Routing is widely used to scale large language models, from Mixture-of-Experts gating to multi-model/tool selection. A common belief is that routing to a task “expert” activates sparser internal computation and thus yields more certain and stable outputs (the Sparsity–Certainty Hypothesis). We test this belief by injecting routing-style meta prompts as a textual proxy for routing signals in front of frozen instruction-tuned LLMs. We quantify (C1) internal density via activation sparsity, (C2) domain-keyword attention, and (C3) output stability via predictive entropy and semantic variation. On a RouterEval subset with three instruction-tuned models (Qwen3-8B, Llama-3.1-8B-Instruct, and Mistral-7B-Instruct-v0.2), meta prompts consistently densify early/middle-layer representations rather than increasing sparsity; natural-language expert instructions are often stronger than structured tags. Attention responses are heterogeneous: Qwen/Llama reduce keyword attention, while Mistral reinforces it. Finally, the densification–stability link is weak and appears only in Qwen, with near-zero correlations in Llama and Mistral. We present RIDE as a diagnostic probe for calibrating routing design and uncertainty estimation.

1 Introduction

As large language models (LLMs) are increasingly deployed in search assistants, coding assistants, and multi-tool systems, routing has become a foundational mechanism. In Mixture-of-Experts (MoE) architectures, a lightweight gating network activates only a small subset of expert subnetworks conditioned on the input, achieving parameter-level sparsification and computational savings. In multi-model routing systems, a policy model selects among multiple LLMs, toolchains, or

APIs to balance performance, cost, and safety. As routing structures grow more complex, a new challenge emerges: we must understand not only whether routing decisions are “correct,” but also how routing signals shape the backbone model’s internal representations and reasoning patterns. A compelling intuition in current practice and informal discussions is that routing effectively selects a narrower, more specialized “expert pathway”, making internal computation sparser and more focused, ultimately producing more certain and stable outputs. We formalize this intuition as the Sparsity–Certainty Hypothesis: when a model receives a clear routing signal (e.g., “this is a math task”), its internal representations propagate along more selectively gated sparse pathways, suppressing task-irrelevant activations and thereby reducing output entropy and increasing cross-sample consistency. This hypothesis has been implicitly adopted in various settings—for instance, using internal signals such as activation sparsity or attention concentration as proxies for uncertainty estimation, early exit, or routing decisions. However, modern instruction-tuned LLMs have highly complex internal structures that strongly depend on training details, and whether this hypothesis holds has not been systematically tested. Directly “opening up” real MoE gating mechanisms or complex policy routers is often impractical. Engineering systems are typically tightly encapsulated, making expert activation dynamics difficult to access; moreover, modifying routers to insert diagnostic signals may compromise stability and safety in production. We therefore take a pragmatic alternative: instead of intervening on the router itself, we inject routing-style meta prompts at the input level as a textual proxy for routing decisions. Concretely, we prepend

different types of prefixes to frozen instruction-tuned LLMs, including structured route tags (e.g., [RouteTag=math]), natural-language expert instructions (e.g., "You are a Math Expert."), as well as control conditions such as incorrect tags and placebo tags. For each instance, we keep the backbone parameters and random seeds identical across prefix conditions and vary only the prefix, thereby constructing paired, controlled intervention experiments. By comparing internal activation density, domain-keyword attention, and output stability across conditions, we can analyze how "routing-style signals" modulate internal states. Building on this idea, we propose RIDE (Route-Induced Density and Stability), centered on three research questions:

- RQ1 (Density effects): Under a frozen backbone, how do different routing-style meta prompts (route tags vs. expert prompts) change the density/sparsity of internal activations? Which layer segments are most affected?
- RQ2 (Attention strategies): Do routing-style meta prompts change how the model leverages domain information? After receiving explicit labels, do models "offload" reliance on domain keywords, or instead increase their focus on keywords?
- RQ3 (Densification–stability relation): At the instance level, is there a stable association between changes in internal densification (C1) and changes in output stability (C3)? Is this association consistent across models, and can it serve as a general proxy for routing decisions?

Importantly, the "routing signals" analyzed in this work are not the internal states of real MoE gating modules or multi-model routers; they are routing-style meta instructions injected as textual prefixes. Our conclusions thus constitute causal-style interventional evidence from controlled experiments, rather than full-fledged structural causal identification.

1.1 Contributions

Within this framework, our main contributions are:

- **RIDE metrics and a unified controlled-intervention pipeline.** We design a suite of routing-style meta prompts and control conditions (control / correct tag / incorrect tag / placebo tag / expert instruction), construct paired interventions on three open-source instruction-tuned LLMs, and define three metric families—C1 (density), C2 (domain-keyword attention), and C3 (output stability)—as a systematic toolkit for analyzing routing-signal effects.
- **Model heterogeneity in densification and attention responses to meta prompts.** Experiments show that all models exhibit significant densification responses to task-oriented meta prompts in the Early/Middle layers, while natural-language expert instructions induce stronger densification than structured route tags in most settings. At the attention level, different models can even exhibit opposite strategies: Llama/Qwen align more with cognitive offloading, whereas Mistral displays attention reinforcement.
- **A systematic test of the "densification \Rightarrow stability" chain with informative negative results.** We observe a small-to-moderate positive correlation consistent with densification–stability coupling in Qwen3-8B, but the link is weak or statistically insignificant at the instance level for Llama-3.1-8B-Instruct and Mistral-7B-Instruct-v0.2. These results suggest limited empirical support for the assumption that internal density serves as a cross-model, general-purpose uncertainty measure or routing signal. Consequently, RIDE is better positioned as a diagnostic probe for revealing heterogeneous model responses to routing-style meta prompts.

2 Background and Related Work

We briefly review three lines of work—routing and routing benchmarks, sparsity–certainty assumptions, and interpretability/meta-prompt analysis—and position RIDE among them.

181	Routing mechanisms and benchmarks.	232
182	In Mixture-of-Experts (MoE) models, a gating network activates a small subset of experts for conditional computation, with extensive work on scaling and training stability (e.g., load balancing and regularization). (Shazeer et al., 2017; Fedus et al., 2022; Du et al., 2022; Lepikhin et al., 2021) However, these advances are largely evaluated by downstream performance and FLOPs, and rarely analyze how gating signals shape intermediate representations. (Zoph et al., 2022; Rajbhandari et al., 2022) In multi-model/tool routing, a policy selects among models or tools under capability–cost trade-offs. (Chen et al., 2024; Schick et al., 2023; Yao et al., 2023) Recent benchmarks and learned routing strategies enable systematic evaluation (e.g., across tasks, difficulty, and cost), but mainly assess input–output selection quality rather than internal-state effects. (Hu et al., 2024; Huang et al., 2025; Song et al., 2025) In contrast, we do not build or optimize routers; we treat routing-style meta prompts as controllable textual proxies of routing signals and analyze their internal effects. (Hendel et al., 2023; Liu et al., 2024; Stolfo et al., 2025)	233 234 235 236 237 238 239 240 241 242 243 244 245
208	Sparsity and certainty. Across compression, MoE, and adaptive inference, activation sparsity/attention concentration and entropy are often used as proxies for specialization or uncertainty, and as signals for early exiting or model switching. (Frankle and Carbin, 2019; Sanh et al., 2020; Fedus et al., 2022; Du et al., 2022; Zhou et al., 2020; Xin et al., 2020; Chen et al., 2024; Vazhentsev et al., 2022) These practices implicitly suggest that “sparser/more concentrated internal states \Rightarrow more stable outputs.” (Zhai et al., 2023) We formalize this as the Sparsity–Certainty Hypothesis and test it via controlled meta-prompt interventions on frozen instruction-tuned LLMs, separating mechanism probing from performance optimization. (Chung et al., 2024; Stolfo et al., 2025)	
226	Interpretability and meta-prompt analysis. Mechanistic interpretability probes internal computation via activations, attention, and circuit-level discovery/attribution. (Meng et al., 2022; Conmy et al., 2023; Ameisen et al., 2025; Zhang et al., 2025; Jain and Wal-	
	lace, 2019; Wiegrefe and Pinter, 2019) Recent work also shows that prompts/instructions can steer internal representations and induce measurable distraction/attraction phenomena that can be localized to specific heads or circuit components. (Hendel et al., 2023; Liu et al., 2024; Niu et al., 2025) We follow this direction but focus on routing-style meta prompts and, unlike single-model qualitative studies, provide a unified multi-model intervention pipeline to compare heterogeneity along density–attention–stability pathways. (Huang et al., 2025; Hu et al., 2024; Rimsky et al., 2024; Yu et al., 2025; Meng et al., 2022)	232 233 234 235 236 237 238 239 240 241 242 243 244 245
	3 RIDE: Routing-Style Meta Prompts and Experimental Design	246 247 248
	3.1 Routing-Style Meta Prompts	249
	We inject routing-style meta prompts by prepending short prefixes to the user instruction. For each instance, we construct five prefix conditions: 1.control : No routing-related prefix is added; the task instruction is fed directly. 2.tag_correct : We add a structured tag that matches the instance domain, e.g., [RouteTag=math] or [RouteTag=code]. The tag format is fixed, with the internal field indicating the domain. 3.tag_wrong : We add a tag that does not match the instance domain (e.g., [RouteTag=code] before a math problem) to probe intervention effects of “incorrect routing signals.” 4.tag_placebo : The tag has the same surface format as tag_correct , but the internal string carries no recognizable semantics (e.g., random tokens), controlling for prefix length and formatting effects. 5.instr_expert : We use natural-language expert instructions such as “You are a Math Expert.” or “You are a coding assistant.”, representing widely used role-setting prompts.	250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271
	All conditions share the same core instruction text, frozen backbone parameters, and decoding configuration; only the prefix differs. Thus, each instance forms a strictly paired intervention experiment across the five conditions, and differences can be interpreted as controlled intervention effects of the “routing-style signal.”	272 273 274 275 276 277 278 279

280	3.2 Dataset and Domain Partitioning	4 Metrics and Analysis Methods	330
281	We construct our experimental samples from	We define three metric families: C1 (activa-	331
282	the RouterEval dataset and focus on three	tion sparsity/density), C2 (domain-keyword	332
283	sub-domains: Math : multi-step mathemat-	attention), and C3 (output stability). Unless	333
284	ical reasoning problems; Format / IFE-	otherwise noted, all metrics are computed per	334
285	val : tasks with strict output-format con-	instance and per prefix condition, then com-	335
286	straints; Commonsense : primarily multiple-	pared via paired differences.	336
287	choice commonsense reasoning tasks. Follow-		
288	ing RouterEval’s easy/hard difficulty split, we	4.1 C1: Activation Sparsity / Density	337
289	sample instances from the training split within	We quantify the sparsity of a hidden vector	338
290	each sub-domain, ensuring that every model	$\mathbf{h} \in \mathbb{R}^d$ using Hoyer sparsity: $\text{Hoyer}(\mathbf{h}) =$	339
291	sees exactly the same inputs under all prefix	$(\sqrt{d} - \ \mathbf{h}\ _1 / \ \mathbf{h}\ _2) / (\sqrt{d} - 1) \in [0, 1]$, where	340
292	conditions. To avoid data leakage, we use only	larger values indicate higher sparsity (thus	341
293	the official training partition and keep the per-	lower values indicate <i>densification</i>). We aggreg-	342
294	domain sample size on the order of hundreds	ate token-level Hoyer scores over layer seg-	343
295	to a thousand. Appendix A reports detailed	ments (Early/Middle/Late). In addition, we	344
296	statistics of sample distributions by domain	compute the Top- k energy ratio (TopK), i.e.,	345
297	and difficulty.	the fraction of ℓ_2 energy captured by the top-	346
		k dimensions after sorting by $ h_i $. We define	347
298	3.3 Models and Decoding	C1 as a segment-level <i>combined metric family</i>	348
299	Configuration	based on Hoyer and TopK; due to space, the	349
300	We evaluate three open-source instruction-	main paper primarily reports Hoyer changes,	350
301	tuned models: Llama-3.1-8B-Instruct ,	while full definitions and additional TopK re-	351
302	Mistral-7B-Instruct-v0.2 , Qwen3-8B .	sults are provided in Appendix B.	352
303	The three models have similar parameter		
304	scales, are instruction-tuned, and perform	4.2 C2: Domain-Keyword Attention	353
305	well across multiple tasks. We use a unified	To measure reliance on domain signals, we	354
306	decoding configuration (temperature, top-p,	build a small keyword lexicon for each sub-	355
307	maximum length, etc.). For each instance	domain (e.g., math, code) and compute the	356
308	and each fixed random seed, we generate K	<i>attention share</i> paid to matched keyword po-	357
309	candidate outputs to estimate output entropy	sitions. Concretely, from the last-layer atten-	358
310	and semantic variation. All prefix conditions	tion matrix, for a query position t we sum the	359
311	share the same set of random seeds, minimiz-	attention mass assigned to keyword token po-	360
312	ing the influence of sampling randomness on	sitions; we then average over (relevant) po-	361
313	paired differences.	sitions and over instances to obtain the domain-	362
		keyword attention share. We report two com-	363
314	3.4 Representation Sampling and	plementary viewpoints: Prompt-last (query	364
315	Layer-Segment Partitioning	at the last input token) and First-gen (query	365
316	We sample hidden states from every layer and	at the first generated token), corresponding	366
317	uniformly partition layers by depth into three	to domain-information usage when “finishing	367
318	segments: - Early: the first third of layers clos-	reading” vs. “starting to answer.” Formal def-	368
319	est to the input embeddings; - Middle: the	initions are in Appendix B.	369
320	middle third; - Late: the final third closest		
321	to the output head. Within each segment,	4.3 C3: Output Stability	370
322	we aggregate over the token dimension (e.g.,	We characterize output stability with two com-	371
323	by mean or max pooling) to obtain a repre-	plementary proxies.	372
324	sentative hidden vector for computing metrics		
325	such as Hoyer sparsity and Top-k energy. For	Predictive entropy (Entropy) . During	373
326	attention-based metrics, we primarily focus on	decoding, we compute token-level softmax en-	374
327	the last layer (or several late layers) as an ap-	tropy and average over the generated sequence;	375
328	proximation of how much the decoder attends	for stochastic decoding, we further average	376
329	to domain keywords during decision making.	this quantity over the K generations to obtain	377

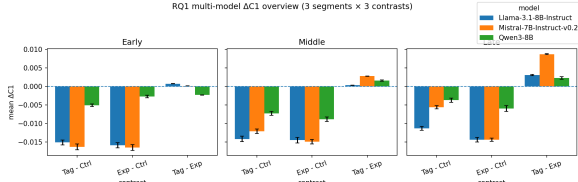


Figure 1: Multi-model overview of ΔHoyer across Early/Middle/Late segments.

a per-instance mean entropy. Lower values indicate a more concentrated predictive distribution.

Semantic variation (SemVar). For K stochastic generations from the same prompt, we compute pairwise embedding similarities and use $\text{Var} = 1 - \text{mF1}$ as a variation score; smaller values indicate higher semantic consistency across generations. Implementation details (encoder choice, K , sampling settings) are deferred to Appendix C. In the main text, we analyze Entropy and SemVar trends separately.

4.4 Paired Differences and Correlation Estimation

Our primary analyses use instance-level paired differences. For a metric M and instance i , we define $\Delta M_i = M_i(m) - M_i(\text{control})$ (e.g., $\Delta\text{Hoyer}_i = \text{Hoyer}_i(\text{tag_correct}) - \text{Hoyer}_i(\text{control})$). Paired differences reduce variance from instance-specific difficulty. We test Δ effects using paired t -tests or Wilcoxon signed-rank tests (with Benjamini–Hochberg FDR correction for multiple comparisons). Unless otherwise stated, correlations between Δ metrics are Pearson correlations; Spearman results are reported in Appendix D and are consistent with the main conclusions.

5 Experimental Results

This section presents the experimental results and analyses for RQ1–RQ3 in turn.

5.1 RQ1: How Do Routing meta prompts Change Internal Density?

Figure 1 reports instance-level paired differences in Hoyer sparsity (ΔHoyer) for `tag_correct` and `instr_expert` relative to `control`, as well as `tag_correct` relative to `instr_expert`.

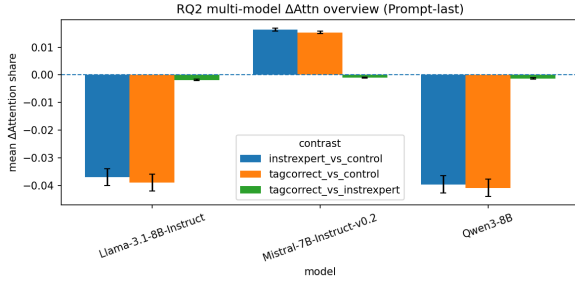
Densification vs. control. Across all three models and most sub-domains, both `instr_expert` and `tag_correct` yield negative ΔHoyer in the Early/Middle segments (typically 0.005–0.015 in magnitude), indicating a shift toward denser and more evenly distributed activations. In contrast, effects in the Late segment are markedly smaller or near zero, suggesting that routing-style meta prompts primarily reshape earlier semantic representations rather than layers close to the output head. As a sanity check, `tag_placebo` shows ΔHoyer close to zero, implying that prefix length or formatting alone cannot explain the densification effect (see Appendix C for full statistics).

Expert instructions are stronger than tags. Directly comparing `tag_correct` with `instr_expert`, we observe a consistently positive global-average gap ($\text{Hoyer}(\text{tag_correct}) - \text{Hoyer}(\text{instr_expert}) > 0$) across all three models. Since larger Hoyer values indicate higher sparsity, this means that natural-language expert instructions induce stronger densification than structured route tags. This reverses the naive expectation that short, formatted tags would behave as a closer proxy to “true” routing signals, and instead highlights the sensitivity of instruction-tuned LLMs to semantically rich natural-language task framing (Appendix B).

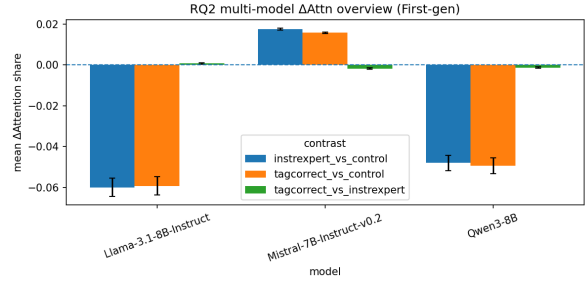
5.2 RQ2: How Do Routing meta prompts Modulate Domain-Keyword Attention?

We measure C2 as the last-layer attention share assigned to domain-keyword tokens. We estimate this share from two query viewpoints: **Prompt-last** (query at the last input token) and **First-gen** (query at the first generated token). Figure 2 reports instance-level paired differences ΔAttn under three contrasts (`tag_correct` vs. `control`, `instr_expert` vs. `control`, and `tag_correct` vs. `instr_expert`).

Model-specific attention strategies. Llama-3.1-8B-Instruct and Qwen3-8B show consistent *decreases* in keyword attention after adding routing-style prefixes (about -0.04 in Prompt-last and $-0.05 \sim -0.06$ in First-gen for Tag/Instr vs. `control`), suggesting an



(a) Prompt-last.



(b) First-gen.

Figure 2: Multi-model overview of ΔAttn (mean \pm SEM). Negative values indicate reduced attention share to domain keywords.

466 *offloading-like* pattern: the explicit routing
 467 cue partially substitutes for lexical domain
 468 cues when the model begins responding. In
 469 contrast, Mistral-7B-Instruct-v0.2 shows a
 470 stable *increase* (about +0.016 in both views),
 471 indicating *reinforcement*: the model attends
 472 to the routing cue *and* the domain keywords
 473 more strongly rather than trading one for the
 474 other. Across models, Tag and Instr typically
 475 have comparable magnitudes within each
 476 view.

477 **ΔAttn vs. ΔC3 : weak coupling.** Across
 478 contrasts and models, changes in keyword at-
 479 tention provide limited explanatory power for
 480 output stability: correlations between ΔAttn
 481 and ΔC3 are small overall (at best weak-to-
 482 moderate for entropy in Qwen, and typically
 483 closer to zero for semantic variation; see Ap-
 484 pendix C.2). Thus, keyword attention is bet-
 485 ter interpreted as a task-identification covari-
 486 ate than as a direct driver of stability.

487 5.3 RQ3: Linking Internal 488 Density to Output Stability

489 RQ3 tests a core link in the Sparsity–Certainty
 490 Hypothesis: whether instance-level changes in
 491 internal density (ΔC1) co-vary with changes
 492 in output stability (ΔC3). We instantiate C1
 493 with prompt-segment Hoyer sparsity and mea-
 494 sure C3 by predictive entropy and semantic
 495 variation ($1 - \text{mF1}$). For each model and each
 496 prefix contrast (*instr_expert* vs. *control*,
 497 *tag_correct* vs. *control*, and *tag_correct*
 498 vs. *instr_expert*), we compute Pearson cor-
 499 relations on paired differences; full statistical
 500 details are in Appendix C.3.1. Figure 3 sum-
 501 marizes the results.

Model-specific coupling. Qwen shows a
 502 consistent positive association between densi-
 503 fication and lower entropy, with correlations
 504 typically in the 0.2–0.3 range across contrasts
 505 (Appendix C.3.1). In contrast, for Llama
 506 and Mistral the corresponding correlations are
 507 near zero across contrasts, indicating that den-
 508 sification provides little predictive signal for
 509 entropy changes in these models.
 510

Semantic variation is weaker overall.
 511 Replacing entropy with semantic variation
 512 ($1 - \text{mF1}$) further attenuates the relationship:
 513 Qwen retains only a small positive association,
 514 while Llama and Mistral remain close to zero
 515 (Appendix C.3.1).
 516

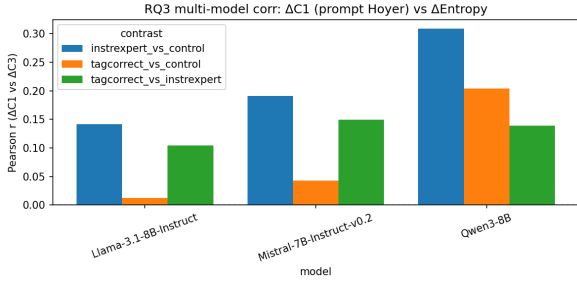
Takeaway. The densification–stability link
 517 is not a model-invariant property: it appears
 518 in Qwen but largely breaks in Llama and Mis-
 519 tral, suggesting that using internal density as
 520 a general-purpose uncertainty proxy requires
 521 model-specific calibration.
 522

523 6 Discussion: RIDE as a Diagnostic 524 Probe Rather Than a Universal 525 Routing Law

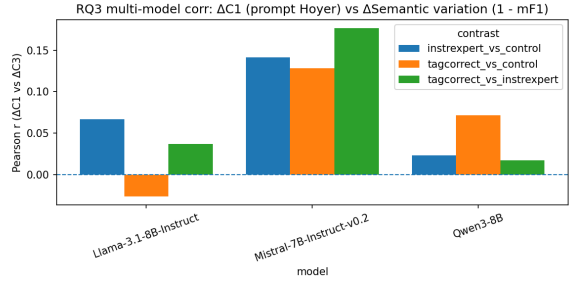
526 Taken together, RQ1–RQ3 paint a picture that
 527 is far more complex than the linear assumption
 528 “routing \Rightarrow sparsity \Rightarrow stability.” Here we dis-
 529 cuss the positioning and implications of RIDE
 530 from three perspectives.

531 6.1 From the Sparsity–Certainty 532 Hypothesis to a Model-Specific 533 Density–Stability spectra

534 Our original motivation was to test the intu-
 535 ition that “effective routing selects sparser and



(a) Entropy



(b) Semantic variation

Figure 3: Instance-level correlations between $\Delta C1$ (prompt Hoyer) and $\Delta C3$ across models and contrasts. Positive values indicate that instances with larger decreases in Hoyer tend to also exhibit larger decreases in the corresponding instability measure (entropy or semantic variation).

more certain pathways.” However, the experimental results show that:

- At the density level: task-oriented meta prompts generally induce densification rather than increased sparsity;
- At the stability level: a local “densification–stability coupling” is observed only in Qwen, while it is largely absent in Llama and Mistral.

Rather than a universal “sparsity–certainty law,” our findings suggest that different models form their own distinctive density–stability spectra under routing-style meta prompts: in some models, densification has a partially positive association with stability, whereas in others the two factors are nearly decoupled. This provides an important caution against directly transferring internal proxies across models.

6.2 Implications for Routing Design and Uncertainty Estimation

From an engineering standpoint, our results suggest that:

- Routing signals should be calibrated in a model-specific manner. A proxy that works well for one model (e.g., C1) may completely fail for another, or even yield an opposite signal.
- Natural-language instructions can themselves serve as strong routing signals. For instruction-tuned models without specialized fine-tuning, short expert instructions can substantially reshape internal density and attention structure, opening up design space for “prompt-level routing.”

- Incorrect or ambiguous routing tags may introduce additional instability. In the absence of a real router, indiscriminate use of “expert tags” may inject noise into internal states, potentially harming stability and safety.

In practical systems, RIDE can be used as a diagnostic tool: prior to deployment, one can run RIDE analyses on candidate models to determine which internal proxies carry stable semantics for that model, and only then consider using them for routing or uncertainty estimation.

6.3 Scope and Limitations of RIDE

RIDE offers a statistical diagnostic perspective that combines meta-prompt interventions with internal metrics, rather than precise identification of causal structure. We do not advocate directly using C1–C3 as decision signals in production systems. Instead, we suggest treating them as auxiliary tools for model understanding and iterative design—for example, selecting models that are better suited for “prompt-level routing,” or identifying architectures that are particularly sensitive to incorrect routing tags.

7 Future Directions

Future work can extend RIDE along several directions:

- Extend the framework to real MoE gating and multi-model routing logs, combining real routing signals with meta prompts;
- Incorporate additional internal metrics (e.g., concept activations, specific neu-

rons/subspaces) to enrich the characterization of density–stability dynamics;

- Systematically evaluate the fairness and robustness of RIDE metrics on multilingual, multi-group, and safety-critical tasks;
- Explore using RIDE as part of the routing-system design workflow, e.g., tracking how the density–stability spectra evolves before and after training.

8 Conclusion

We propose RIDE (Route-Induced Density and Stability), a framework that performs controlled interventions on frozen instruction-tuned LLMs via routing-style meta prompts, and systematically analyzes how routing signals affect internal density (C1), domain-keyword attention (C2), and output stability (C3). Based on large-scale experiments on a RouterEval subset across three open-source instruction-tuned models, we find that:

- Task-oriented meta prompts consistently induce densification of internal representations rather than increased sparsity;
- Models exhibit pronounced heterogeneity in attention redistribution (cognitive offloading vs. attention reinforcement);
- The “densification \Rightarrow stability” link receives only limited support in a subset of models and lacks cross-model robustness.

Taken together, these findings suggest that directly treating internal density or attention concentration as a universal “good routing signal” or a general-purpose uncertainty proxy can be risky. We argue that RIDE is better positioned as a diagnostic probe: through controlled, statistical analyses, it reveals model-specific density–stability spectra under routing-style meta prompts, providing fine-grained evidence for routing design, uncertainty estimation, and model selection.

9 Limitations

This work has several important limitations (which we state explicitly to clarify the scope of generalization):

1. Proxy nature of routing signals and external validity. We use routing-style meta prompts as a textual proxy for real MoE gating or multi-model router signals. This proxy cannot capture routing distributions, logging features, or training couplings present in real systems; therefore, our conclusions primarily apply to intervention analyses of “prompt-level routing signals.”
2. Limited coverage of tasks and models. Our experiments focus on a subset of RouterEval and three mid-sized open-source instruction-tuned models. Larger-scale models, different alignment strategies or specialist models, and different task distributions—especially long-context and tool-use settings—may exhibit different RIDE spectra.
3. Dependence on metric construction and decoding settings. The concrete implementations of C1–C3 (e.g., the exact Hoyer sparsity formulation, keyword-list coverage, sampling counts, and temperature) affect numerical scales. We provide several robustness checks in the appendix, but cannot exhaustively cover all configurations; using RIDE as a system signal requires model- and scenario-specific recalibration.
4. Boundaries of causal interpretation. Although we use paired interventions and control decoding and randomness, RIDE remains a “causal-style” statistical comparison rather than structured causal identification. We cannot fully rule out the influence of unobserved confounders (e.g., training history or alignment preferences) on the observed effects.

Ethical considerations

All experiments in this work are conducted using open-source models and publicly available datasets, without involving user privacy or sensitive data. The code and derived artifacts we release will not contain any personally identifiable information. RIDE is primarily intended as an analysis and diagnostic tool; any application that uses internal proxies for safety-critical decisions (e.g., filtering sensitive con-

698	tent, or routing decisions in medical or legal settings) should undergo rigorous safety evaluation and ethical review.	
699		
700		
701	References	
702	Emmanuel Ameisen, Jack Lindsey, Adam Pearce, Wes Gurnee, Nicholas L. Turner, Brian Chen, Craig Citro, and 1 others. 2025. Circuit tracing: Revealing computational graphs in language models . <i>Transformer Circuits (Anthropic)</i> .	
703		
704		
705		
706		
707	Lingjiao Chen, Matei Zaharia, and James Zou. 2024. Frugalgpt: How to use large language models while reducing cost and improving performance . <i>Transactions on Machine Learning Research</i> . Published: 20 Dec 2024.	
708		
709		
710		
711		
712	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, and 16 others. 2024. Scaling instruction-finetuned language models . <i>Journal of Machine Learning Research</i> , 25(70):1–53.	
713		
714		
715		
716		
717		
718		
719		
720		
721		
722	Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	
723		
724		
725		
726		
727		
728	Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten P. Bosma, Zongwei Zhou, Tao Wang, Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, and 8 others. 2022. GLaM: Efficient scaling of language models with mixture-of-experts . In <i>Proceedings of the 39th International Conference on Machine Learning (ICML)</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 5547–5569. PMLR.	
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740	William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity . <i>Journal of Machine Learning Research</i> , 23(120):1–39.	
741		
742		
743		
744		
745	Jonathan Frankle and Michael Carbin. 2019. The lottery ticket hypothesis: Finding sparse, trainable neural networks . In <i>International Conference on Learning Representations (ICLR)</i> .	
746		
747		
748		
749	Roe Hendel, Mor Geva, and Amir Globerson. 2023. In-context learning creates task vectors . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9318–9333, Singapore. Association for Computational Linguistics.	751
750		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807

The resulting `<items_routereval_csv>` serves as a unified “analysis corpus table” across models and prefix conditions in subsequent experiments, recording metadata such as the task text, domain labels, and source scenarios.

A.2.3 Sample Filtering and Summary Statistics

Based on the items table above, we apply lightweight filtering:

- We remove overly long instances (exceeding a predefined maximum token length) to avoid model-dependent truncation effects;
- We remove extremely short instances (fewer than a minimum number of tokens) to prevent unstable estimates for metrics such as C2/C3.

After filtering, we retain approximately 200 instances per domain, forming a fixed analysis corpus. Table A.2 reports summary statistics for each domain, including the number of instances, the average token length, and the keyword matching rate. Because our focus is on internal change patterns under prefix interventions rather than generalization performance on this dataset, we do not further split this corpus into train/dev/test. Instead, we treat it as a fixed analysis set and compare different models on the same collection.

A.3 Router-Style Meta-Prompt Templates and Tag Construction

A.3.1 Route Tags

We adopt a unified structured format for route tags, as follows: `[RouteTag=DOMAIN]` where `DOMAIN` \in `{math, code, format, commonsense}`.

- `tag_correct`: uses a tag that matches the instance domain (e.g., using `[RouteTag=math]` for a math problem);
- `tag_wrong`: uniformly samples an incorrect tag from the other domains (e.g., using `[RouteTag=code]` for a math instance).

A.3.2 Placebo Tags

Placebo tags are used to control for form factors—i.e., adding extra tokens in a similar format—without introducing meaningful

semantics. We construct tags of the form: `[RouteTag=XXXXX]` where `XXXXX` is a meaningless string that rarely appears in the RouterEval corpus and has a length comparable to real route tags. Semantically, `tag_placebo` should be approximately equivalent to control, differing only in surface form by the additional prefix tokens.

A.3.3 Expert Prompts

We design a small set of natural-language expert prompt templates for different domains, for example:

- **Math** You are a Math Expert. Please solve the problem step by step. (A Chinese variant is also used: “你是一名数学专家，请逐步推理并给出答案。”)
- **Format / IFEval** You are a formatting assistant. Please strictly follow the required output format.
- **Commonsense** You are a reasoning assistant. Please choose the most plausible option and briefly explain why.

In the main experiments, we use a fixed template per domain to reduce additional variance. In Appendix C, we conduct sensitivity checks with partially randomized template sets.

A.3.4 Prefix Injection Location and Dialogue Structure

Unless otherwise specified, all route tags and expert prompts are injected before the user prompt, e.g., `[RouteTag=math] <user prompt ...>` or `You are a Math Expert. Please solve the problem step by step. <user prompt ...>` Other parts of the system/assistant roles (e.g., the system prompt and assistant prefix) are kept unchanged under each model’s default settings. Across prefix conditions, the only differences lie in the routing-style meta-prompt segment described above.

A.4 Domain Keyword Lexicon (DOMAIN_KEYWORDS)

The domain keyword lexicon is used to construct the C2 metric (domain-keyword attention share). We build it as follows:

1. **Collecting seed keywords** For each sub-domain, we manually curate a small set of seed keywords:

1101	• Math: equation, integer, sum, proof,	• Early: Layers 1 to $\lfloor L/3 \rfloor$;	1148
1102	factor		
1103	• Code: function, variable, class, com-	• Middle: Layers $\lfloor L/3 \rfloor + 1$ to $\lfloor 2L/3 \rfloor$;	1149
1104	pile, error		
1105	• Format: json, table, markdown, field,	• Late: remaining layers.	1150
1106	column		
1107	• Commonsense: choose, option, likely,	For each segment, we aggregate across in-	1151
1108	most, best).	stances and tokens (e.g., by taking the mean)	1152
1109	These seeds aim to cover the most cen-	to obtain segment-level metrics.	1153
1110	tral semantic concepts or task-triggering		
1111	terms in the domain.	A.5.2 Caching Hidden States and	1154
1112	2. Automatic expansion and filtering:	Attention Weights	1155
1113	• We use an auxiliary script to auto-	For each instance and each prefix condition, we	1156
1114	matically expand the candidate list	cache the following during the forward pass:	1157
1115	on the RouterEval training table us-	• The hidden vector at every token for every	1158
1116	ing statistics such as TF-IDF / PMI;	layer, $\mathbf{h}_t^{(\ell)}$;	1159
1117	• We remove stopwords, punctuation,	• Attention weights used for C2:	1160
1118	and extremely frequent but domain-	• Prompt-last: the attention over all	1161
1119	agnostic words (e.g., the, is, do), and	prompt tokens at the final step before	1162
1120	normalize case and inflectional forms.	generation begins;	1163
1121	3. Manual screening and finalization:	• First-gen: the attention from the first	1164
1122	• We perform a quick manual inspec-	generated token to all prompt tokens.	1165
1123	tion of the expanded candidates and	In practice, we enable out-	1166
1124	discard clearly inappropriate entries;	put_hidden_states=True and out-	1167
1125	• We then form the final DO-	put_attentions=True in the inference API,	1168
1126	MAIN_KEYWORDS dictionary	or use framework-provided hooks to write	1169
1127	and save it as a JSON file for	intermediate results to files (e.g., Parquet /	1170
1128	downstream analysis scripts.	HDF5) for offline analysis.	1171
1129	For each instance, after tokenization, we	A.5.3 Attention Normalization and	1172
1130	match terms from DOMAIN_KEYWORDS	Keyword-Share Computation	1173
1131	and record their token positions:	When computing the keyword attention share,	1174
1132	• When computing the domain-keyword at-	we preprocess the attention vector of each	1175
1133	tention share, the target set consists of	layer as follows:	1176
1134	these token positions;	• We retain attention mass only over visible	1177
1135	• If an instance has no keyword match un-	tokens, including system/user/prompt to-	1178
1136	der its domain, we define its domain-	kens, while excluding padding and special	1179
1137	keyword attention share as 0 (since no	tokens (e.g., BOS/EOS);	1180
1138	attention mass falls on the keyword set),	• We re-normalize over the visible-token set	1181
1139	and mark it as “missing keywords” in	so that the weights sum to 1;	1182
1140	the summary table (used to analyze the	• We sum the weights over positions that	1183
1141	signal-to-noise ratio of C2).	match DOMAIN_KEYWORDS, yielding	1184
1142	A.5 Extracting Activations and	the layer-level keyword attention share;	1185
1143	Attention	• We then aggregate across layer segments	1186
1144	A.5.1 Layer-Segment Partitioning	and instances to obtain the segment-level	1187
1145	Suppose the model contains L Transformer	C2 metric.	1188
1146	blocks. We evenly partition them into three		
1147	segments:		

1189	A.6 Decoding and Output Sampling			1232
1190	(C3 Metric)			1233
1191	A.6.1 Decoding Hyperparameters			1234
1192	For experiments related to output stability, we			
1193	use stochastic sampling with temperature. Un-			
1194	less otherwise specified, we adopt the following			
1195	unified hyperparameters:			
1196		• Sampling temperature $T=0.7$;		1235
1197		• $\text{top-}p = 0.9$;		
1198		• No additional repetition penalty;		1236
1199		• A fixed generation length limit		1237
1200		<code>max_new_tokens</code> (e.g., 64).		1238
1201	These settings are applied uniformly across all			1239
1202	models and all prefix conditions.			1240
1203	A.6.2 Multi-Sample Generation			1241
1204	(Multi-sample Forward)			1242
1205	To estimate C3 (predictive entropy and seman-			1243
1206	tic variation), we generate K independent sam-			1244
1207	ples for each prompt under each prefix condi-			1245
1208	tion. In the main experiments:			1246
1209		• We use $K=5$ samples by default to es-		1247
1210		timate predictive entropy and semantic		
1211		variation;		
1212		• For the same instance, we reuse the same		
1213		sequence of random seeds across different		
1214		prefix conditions, ensuring that C3 differ-		
1215		ences are primarily attributable to prefix		
1216		differences rather than sampling random-		
1217		ness.		
1218	Predictive entropy is computed at the token			
1219	level and then averaged over the sequence. Se-			
1220	matic variation is computed based on simi-			
1221	larities among sentence embeddings obtained			
1222	from multiple samples; the formal definition is			
1223	provided in Appendix B. Appendix C reports			
1224	sensitivity analyses under different K values			
1225	and temperatures.			
1226	A.6.3 Sharding and Parallel Inference			
1227	To improve computational efficiency, we shard			
1228	<code><items_routereval_csv></code> by splitting the in-			
1229	stances into multiple shards and run inference			
1230	and metric computation in parallel on a multi-			
1231	GPU setup. The overall workflow is:			
		• Use an auxiliary script to split the items		1232
		CSV into multiple subfiles by row count		1233
		(e.g., <code>items_part0.csv</code> , <code>items_part1.csv</code> ,		1234
		<code>...</code>);		1235
		• Assign one or more shards to each GPU		1236
		and run a unified inference script (which		1237
		loads the model, applies prefixes, extracts		1238
		hidden states and attention, and gener-		1239
		ates multi-sample outputs);		1240
		• Save intermediate results for each shard		1241
		(e.g., MFV metrics and generated texts)		1242
		as structured files, and merge them during		1243
		post-processing.		1244
	All shards are executed under the same config-			1245
	uration, ensuring comparability across GPUs			1246
	and shards.			1247
	A.7 Statistical Tests and Correlation			1248
	Estimation			1249
	To obtain causal-style paired evidence and to			1250
	avoid over-interpreting any single scalar result,			1251
	we adopt the following statistical procedures:			1252
	Paired difference tests			1253
		• For paired metric differences (e.g.,		1254
		ΔHoyer , $\Delta\text{Entropy}$, $\Delta\text{KeywordShare}$),		1255
		we use a paired t-test or the Wilcoxon		1256
		signed-rank test to assess whether the		1257
		differences significantly deviate from		1258
		zero;		1259
		• When appropriate, we report effect sizes		1260
		(e.g., Cohen's d) to reflect the practical		1261
		magnitude of the differences.		1262
	Correlation analysis			1263
		• To quantify the linear relationship be-		1264
		tween ΔC1 and ΔC3 , we use the Pearson		1265
		correlation coefficient;		1266
		• We estimate the 95% confidence interval		1267
		of Pearson r via the Fisher transforma-		1268
		tion;		1269
		• In Appendix B, we additionally report		1270
		Spearman's rank correlation as a robust-		1271
		ness check.		1272
	Multiple-comparison correction			1273

- For the large number of tests across models, domains, and layer segments, we control the false discovery rate (FDR) using the Benjamini–Hochberg procedure, reporting significance at a unified q-level;
- In the main paper, we report only conclusions with clear direction and magnitude that remain stable under multiple-comparison correction, avoiding over-interpretation of marginally significant results.

B Full Mathematical Definitions of Metrics

This appendix provides the formal mathematical definitions of C1/C2/C3 and several auxiliary metrics used in Section 4 of the main paper. Notation is kept consistent with the main text.

B.1 Notation

- We denote an instance by x , and a routing prefix condition by $m \in \mathcal{M}$ (e.g., `control`, `tag_correct`, `tag_wrong`, `tag_placebo`, `instr_expert`).
- The model has L Transformer blocks. The hidden vector at layer ℓ and token position t is $\mathbf{h}_{x,t}^{(\ell)} \in \mathbb{R}^d$.
- Let the generated sequence length be T_x . The predictive distribution for the t -th generated token under condition m is $p_t^{(m)}(v)$.
- For attention, we denote the last layer by $\ell = L$, and its attention vector by $\mathbf{a}_{x,q}^{(L,m)}$, where q is the query position (e.g., `prompt-last` or `first-gen`).
- The set of token positions in the prompt segment is \mathcal{P}_x , and the set of domain-keyword positions is $\mathcal{K}_x \subseteq \mathcal{P}_x$.

Unless otherwise stated, instance-level metrics are first averaged over tokens and/or layers, and then subjected to statistical analysis over the set of instances.

B.2 C1: Activation Sparsity / Density

In Section 4.1 of the main paper, C1 is described as a combined metric of segment-level

Hoyer sparsity and Top- k energy. Here we provide precise definitions of both components and explain how C1 is constructed.

B.2.1 Hoyer Sparsity

Given a hidden vector $\mathbf{h} \in \mathbb{R}^d$, the Hoyer sparsity is defined as:

$$\text{Hoyer}(\mathbf{h}) = \frac{\sqrt{d} - \frac{\|\mathbf{h}\|_1}{\|\mathbf{h}\|_2}}{\sqrt{d} - 1}.$$

where $\mathbf{h} \in \mathbb{R}^d$ and $\text{Hoyer}(\mathbf{h}) \in [0, 1]$. Larger values indicate higher sparsity and stronger concentration on a small subset of dimensions. For an instance x , prefix condition m , layer ℓ , and token position t , the token-level Hoyer score is:

$$s_{x,t}^{(\ell,m)} = \text{Hoyer}\left(\mathbf{h}_{x,t}^{(\ell,m)}\right).$$

We average over a layer segment. For example, for the Early segment (layer set $\mathcal{L}_{\text{Early}}$), we define:

$$\text{C1_Hoyer}^{\text{Early}}(x, m) = \frac{1}{|\mathcal{L}_{\text{Early}}| \cdot |\mathcal{P}_x|} \times \sum_{\ell \in \mathcal{L}_{\text{Early}}} \sum_{t \in \mathcal{P}_x} s_{x,t}^{(\ell,m)}.$$

The Middle and Late segments are defined analogously. For the first-generated-token variant (C1-firstgen), we replace \mathcal{P}_x with the token position corresponding to the first generated token. The ‘‘Hoyer part of C1’’ in the main paper refers to this type of segment-level aggregation.

B.2.2 Top- k Energy Ratio

Let $|h_{(1)}| \geq \dots \geq |h_{(d)}|$ denote the components of $\mathbf{h} \in \mathbb{R}^d$ sorted by absolute value. The Top- k energy ratio is defined as:

$$\text{TopK}(\mathbf{h}) = \frac{\sum_{i=1}^k h_{(i)}^2}{\sum_{i=1}^d h_i^2}, \quad k = \lfloor \alpha d \rfloor,$$

with $\alpha = 0.1$ by default.

In our setting, decreases in Hoyer (densification) are typically accompanied by decreases in the Top- k energy ratio, indicating that energy is more evenly distributed across a larger number of dimensions.

Similarly, we define the token-level $\text{TopK}_{x,t}^{(\ell,m)}$ and average over layer segments

and prompt-token positions. For example, for the Early segment:

$$\begin{aligned} \text{C1_TopK}^{\text{Early}}(x, m) &= \left(|\mathcal{L}_{\text{Early}}| |\mathcal{P}_x| \right)^{-1} \\ &\times \sum_{\ell \in \mathcal{L}_{\text{Early}}} \sum_{t \in \mathcal{P}_x} \text{TopK}_{x,t}^{(\ell, m)}. \end{aligned}$$

B.2.3 Combining C1 and Paired Differences

Combining C1. We treat C1 as a combined metric consisting of two components, C1_Hoyer and C1_TopK:

- In Sections 4.1–5.1 of the main paper, almost all figures and conclusions are based on C1_Hoyer.
- C1_TopK and its relationship with Hoyer are reported in supplementary figures/tables in Appendix C, mainly to verify that “densification” is not an artifact caused by extreme amplification of only a few dimensions.

Paired differences. At the instance level, we focus on paired differences between a given prefix condition and the control. For an instance x , the Hoyer difference between `tag_correct` and `control` is defined (consistent with Section 4.4) as:

$$\Delta\text{Hoyer}_x = \text{C1_Hoyer}(x, \text{tag_correct}) - \text{C1_Hoyer}(x, \text{control}).$$

In implementation, ΔC1 corresponds to the combination of ΔHoyer and ΔTopK (we report one or both depending on the figure/table). In the main paper, “densification” typically refers to the case $\Delta\text{Hoyer}_x < 0$.

B.3 Auxiliary Metrics for Activation-Distribution Shape

To characterize activation distributions in a more fine-grained manner, we additionally compute the following metrics for each token-level vector $\mathbf{h} \in \mathbb{R}^d$:

- **Gini coefficient:** computed from the sequence $(|h_i|)_{i=1}^d$ using the standard definition, measuring inequality in energy allocation across dimensions;
- **Kurtosis:** quantifies whether the distribution exhibits sharp peaks and heavy tails;

• Positive ratio:

$$\text{PosRatio}(\mathbf{h}) = \frac{1}{d} \sum_{i=1}^d \mathbf{1}[h_i > 0],$$

capturing the sign pattern of activations.

These metrics can likewise be averaged over layer segments and used to form paired differences, but they are only used for robustness checks in Appendix C and are not treated as primary evidence for the main-paper conclusions.

B.4 Cross-Layer Energy Concentration (Layer-Energy Gini)

To examine whether meta prompts alter the allocation of “computational load” across layers, we define a cross-layer energy Gini metric as follows. At a key position in the target sequence (e.g., prompt-last or first-gen), we compute the ℓ -th layer’s ℓ_2 energy:

$$E_\ell(x, m) = \left\| \mathbf{h}_{x,t^*,m}^{(\ell)} \right\|_2^2, \quad \ell = 1, \dots, L,$$

where t^* denotes the token position of interest.

We then normalize the energies across layers:

$$\tilde{E}_\ell(x, m) = \frac{E_\ell(x, m)}{\sum_{\ell'=1}^L E_{\ell'}(x, m)}.$$

Treating $\{\tilde{E}_\ell(x, m)\}_{\ell=1}^L$ as a one-dimensional discrete distribution, we compute its Gini coefficient, denoted as $\text{Gini}_{\text{layer}}(x, m)$. Larger values indicate that energy is highly concentrated in a small subset of layers, whereas smaller values indicate a more even distribution across layers. This metric is used for supplementary analyses in Appendix C but is not directly referenced in the main paper.

B.5 C2: Domain-Keyword Attention

Section 4.2 of the main paper defines C2 as the attention share to domain keywords in the last layer, estimated from two viewpoints: PROMPT-LAST and FIRST-GEN. This section presents the formal definition.

Let \mathcal{P}_x denote the set of prompt-token positions for an instance x , and let $\mathcal{K}_x \subseteq \mathcal{P}_x$ denote the subset of positions matched by domain keywords. Under prefix condition m , we denote

the attention vector at the last layer ($\ell = L$) for a query position q by $\mathbf{a}_{x,q}^{(L,m)} \in \mathbb{R}^{|\mathcal{P}_x|}$, where the element $a_{x,q,t}^{(L,m)}$ is the attention weight assigned to position t .

We first renormalize attention over the set of visible tokens $\mathcal{V}_x \subseteq \mathcal{P}_x$ such that

$$\sum_{t \in \mathcal{V}_x} a_{x,q,t}^{(L,m)} = 1.$$

The domain-keyword attention share at query position q is then defined as

$$\text{C2}(x, m, q) = \sum_{t \in \mathcal{K}_x} a_{x,q,t}^{(L,m)}.$$

In the main paper, we consider two query viewpoints:

- **PROMPT-LAST**: set $q = q_{\text{promptlast}}$ as the last input-token position, corresponding to the state after “reading the question”;
- **FIRST-GEN**: set $q = q_{\text{firstgen}}$ as the first generated-token position, corresponding to the state when “starting to answer”.

For each viewpoint, we average over instances:

$$\text{C2}^{\text{PromptLast}}(m) = \frac{1}{|X|} \sum_{x \in X} \text{C2}(x, m, q_{\text{promptlast}}),$$

and analogously for the **FIRST-GEN** variant.

At the instance level, paired differences follow the Δ notation in Section 4.4. For example,

$$\Delta \text{C2}_x^{\text{PromptLast}} = \text{C2}(x, m, q_{\text{promptlast}}) - \text{C2}(x, \text{control}, q_{\text{promptlast}}).$$

A positive ΔC2 indicates that prefix m increases the attention share to domain keywords relative to the control.

Note. In some figures/tables in Appendix C, we further extend C2 to Early/Middle/Late segments by averaging attention across multiple layers. However, the core results shown in Sections 4.2–5.2 of the main paper are based on the last-layer definition above.

If an instance has no keyword matches in its domain ($\mathcal{K}_x = \emptyset$), we define $\text{C2}(x, m, q) = 0$ and track the proportion of such “missing-keyword” instances separately in the summary tables.

B.6 C3: Output Stability

Section 4.3 of the main paper defines C3 as a pair of measures: predictive entropy and semantic variation. This section provides the formal definitions.

B.6.1 Predictive Entropy

For an instance x under a prefix condition m , let $p_t^{(m)}(v)$ denote the predictive distribution over vocabulary items v for the t -th generated token. The token-level predictive entropy is

$$H_t(x, m) = - \sum_v p_t^{(m)}(v) \log p_t^{(m)}(v).$$

Let T_x be the generated sequence length. The sequence-level average entropy is

$$\bar{H}(x, m) = \frac{1}{T_x} \sum_{t=1}^{T_x} H_t(x, m).$$

We use $\bar{H}(x, m)$ as the first C3 measure: lower values indicate a more concentrated predictive distribution and hence higher confidence.

At the instance level, we define the paired difference

$$\Delta \text{Entropy}_x = \bar{H}(x, m) - \bar{H}(x, \text{control}),$$

where m can be a prefix condition such as `tag_correct` or `instr_expert`. The $\Delta \text{Entropy}$ reported in Section 5.3 refers to this quantity.

Note. In implementation, we compute $p_t^{(m)}(v)$ directly from a single forward-pass softmax. Alternatively, one may average distributions over multiple samples; the qualitative conclusions remain essentially the same.

B.6.2 Semantic Variation

Semantic variation is computed via repeated sampling and sentence-embedding similarity. For the same prompt-condition pair (x, m) , we generate K independent outputs $y^{(1)}(x, m), \dots, y^{(K)}(x, m)$ and encode them using a fixed sentence encoder $f(\cdot)$:

$$\mathbf{e}^{(k)}(x, m) = f(y^{(k)}(x, m)), \quad k = 1, \dots, K.$$

We then compute all pairwise cosine similarities

$$s_{k\ell}(x, m) = \frac{\mathbf{e}^{(k)}(x, m) \cdot \mathbf{e}^{(\ell)}(x, m)}{\|\mathbf{e}^{(k)}(x, m)\| \|\mathbf{e}^{(\ell)}(x, m)\|}$$

$$1 \leq k < \ell \leq K.$$

We treat the average pairwise similarity as a proxy for semantic F1 (mF1):

$$\text{mF1}(x, m) = \frac{2}{K(K-1)} \sum_{1 \leq k < \ell \leq K} s_{k\ell}(x, m).$$

Semantic variation is then defined as

$$\text{Var}(x, m) = 1 - \text{mF1}(x, m).$$

Smaller $\text{Var}(x, m)$ indicates that multiple generations are closer in semantic space and thus more stable. The paired difference is

$$\Delta\text{Var}_x = \text{Var}(x, m) - \text{Var}(x, \text{control}).$$

Following Section 4.3, concrete choices of the encoder, the sampling count K , and other hyperparameters are deferred to Appendix C. In this paper, C3 = {Entropy, Var}, and we analyze $\Delta\text{Entropy}$ and ΔVar separately.

B.7 Commonsense Confidence Margin (Optional Analysis)

For multiple-choice commonsense tasks, let \mathcal{C} denote the set of answer options, and let $c_{\text{gold}} \in \mathcal{C}$ be the correct option. We define the average probability of an option c as

$$\bar{p}(c; x, m) = \frac{1}{G_x} \sum_{t \in \mathcal{G}_x} p_t^{(m)}(c),$$

where $p_t^{(m)}(c)$ is the probability assigned to choosing option c at position t (the exact scoring procedure is described in Appendix C), and \mathcal{G}_x denotes the set of positions used for option scoring.

We define the confidence margin as

$$\Delta_{\text{conf}}(x, m) = \bar{p}(c_{\text{gold}}; x, m) - \max_{\substack{c \in \mathcal{C} \\ c \neq c_{\text{gold}}}} \bar{p}(c; x, m).$$

Larger $\Delta_{\text{conf}}(x, m)$ indicates a stronger confidence advantage for the correct option and hence lower uncertainty. In Appendix C.5, we use this quantity as an auxiliary metric and analyze it alongside the two primary C3 metrics (Entropy and Var), without separately elaborating it in the main text.

With these definitions, all terms appearing in Sections 4.1–4.4 of the main paper—such as C1/C2/C3, ΔHoyer , Prompt-last, and First-gen—have precise counterparts in this appendix, and this section can be directly included in Appendix B.

Model	Δ Hoyer Tag-Ctrl	Δ Hoyer Instr-Ctrl	Δ Hoyer Tag-Instr
Llama-3.1-8B-Instruct	-0.015 ± 0.001	-0.016 ± 0.001	$+0.001 \pm 0.000$
Mistral-7B-Instruct-v0.2	-0.016 ± 0.001	-0.017 ± 0.001	$+0.000 \pm 0.000$
Qwen3-8B	-0.005 ± 0.000	-0.003 ± 0.000	-0.002 ± 0.000

Table 1: Early segment: mean \pm SEM of Δ Hoyer under three contrasts.

Model	Δ Hoyer Tag-Ctrl	Δ Hoyer Instr-Ctrl	Δ Hoyer Tag-Instr
Llama-3.1-8B-Instruct	-0.014 ± 0.001	-0.015 ± 0.001	$+0.000 \pm 0.000$
Mistral-7B-Instruct-v0.2	-0.012 ± 0.001	-0.015 ± 0.001	$+0.003 \pm 0.000$
Qwen3-8B	-0.007 ± 0.001	-0.009 ± 0.001	$+0.002 \pm 0.000$

Table 2: Middle segment: mean \pm SEM of Δ Hoyer under three contrasts.

C Supplementary Experimental Results

This appendix summarizes detailed statistical results based on real experimental data (from `routereval_rq1_multimodel`, `rq2_multimodel`, and `rq3_multimodel`).

C.1 Densification Results by Model and Layer Segment (RQ1)

Tables 1–3 report changes in Hoyer sparsity (Δ Hoyer) for each model across layer segments. Values are reported as mean \pm standard error of the mean (SEM). Negative values indicate densification (lower sparsity).

Interpretation note. For the contrast `Tag-Instr`, a *positive* value means $\text{Hoyer}(\text{Tag}) > \text{Hoyer}(\text{Instr})$, i.e., expert instructions induce stronger densification than tags.

C.2 Domain-Keyword Attention (RQ2)

C.2.1 RQ2 Statistical Details

For RQ2, we compute instance-level paired differences in the domain-keyword attention share, $\Delta\text{Attn} = \text{Attn}(\text{condition}) - \text{Attn}(\text{control})$, under two query viewpoints (Prompt-last and First-gen). Tables 4–5 summarize these effects as mean \pm standard error of the mean (SEM). When the text describes an effect as statistically significant, it is based on paired *t*-tests of ΔAttn against zero with Benjamini–Hochberg FDR control within each viewpoint/table family (test statistics and adjusted *p*-values are omitted from the tables for space).

To relate attention changes to output stability, we compute instance-level Pearson corre-

lations between ΔAttn and ΔC3 metrics; Table 6 reports Pearson’s *r* with two-sided *p*-values.

C.2.2 ΔAttn from the Prompt-last View

Table 4 reports paired differences in the domain-keyword attention share under the Prompt-last view ($\Delta\text{Attn} = \text{condition} - \text{control}$). Negative values indicate a reduced attention share to domain keywords under the condition, whereas positive values indicate an increased share. Values are reported as mean \pm standard error of the mean (SEM).

C.2.3 ΔAttn from the First-gen View

Table 5 reports paired differences in the domain-keyword attention share under the First-gen view (i.e., at the first generated token), defined in the same way as Table 4.

We observe that:

- For Llama and Qwen, ΔAttn is significantly negative under both views (reduced reliance on keywords), roughly consistent with “cognitive offloading.”
- For Mistral, ΔAttn is significantly positive under both views (increased reliance on keywords; “attention reinforcement”), with a slightly larger magnitude under `Instr` than `Tag`.

C.2.4 Correlation Between ΔAttn and Delta Entropy

Table 6 reports the instance-level Pearson correlation coefficient between First-gen ΔAttn and $\Delta\text{Entropy}$ under the `instr_expert` vs. `control` contrast. Here, ΔAttn uses `attn_domain_share_firstgen_`

Model	Δ Hoyer Tag-Ctrl	Δ Hoyer Instr-Ctrl	Δ Hoyer Tag-Instr
Llama-3.1-8B-Instruct	-0.008 ± 0.000	-0.012 ± 0.000	$+0.004 \pm 0.000$
Mistral-7B-Instruct-v0.2	-0.002 ± 0.000	-0.015 ± 0.000	$+0.013 \pm 0.000$
Qwen3-8B	-0.001 ± 0.000	-0.007 ± 0.001	$+0.006 \pm 0.000$

Table 3: Global segment: mean \pm SEM of instance-level paired differences in Δ Hoyer under three contrasts.

Model	Δ Attn Tag-Ctrl (Mean \pm SEM)	Δ Attn Instr-Ctrl (Mean \pm SEM)
Llama-3.1-8B-Instruct	-0.0389 ± 0.0031	-0.0371 ± 0.0030
Mistral-7B-Instruct-v0.2	$+0.0154 \pm 0.0005$	$+0.0164 \pm 0.0005$
Qwen3-8B	-0.0409 ± 0.0032	-0.0396 ± 0.0031

Table 4: Mean \pm SEM of Δ Attn under Tag-Ctrl and Instr-Ctrl across three models (Prompt-last view).

1632 `delta_instrexpert_vs_control,` and
1633 `Δ Entropy` uses `gen_mean_entropy_delta_`
1634 `instrexpert_vs_control.`

1635 Brief interpretation:

- 1636 • Qwen shows the highest correlation ($r \approx$
1637 0.23), suggesting a relatively stronger link-
1638 age between keyword-attention changes
1639 and entropy changes in this model.
- 1640 • Llama exhibits a moderate-to-weak posi-
1641 tive correlation ($r \approx 0.12$).
- 1642 • Mistral shows a significant negative corre-
1643 lation, consistent with its “attention rein-
1644 forcement” pattern accompanied by lim-
1645 ited gains in stability.

1646 C.3 Correlations Between C1 and C3 1647 (RQ3)

1648 **Notation.**

- 1649 • **C1:** `c1_prompt_hoyer_mean` (paired
1650 difference in Hoyer sparsity over the
1651 prompt segment).
- 1652 • **C3-Entropy:** paired difference of
1653 `gen_mean_entropy`.
- 1654 • **C3-SemVar:** paired difference of
1655 `semantic_variation_1mF1`.
- 1656 • **Instr vs. Ctrl:**
1657 `instrexpert_vs_control`.
- 1658 • **Tag vs. Ctrl:** `tagcorrect_vs_control`.
- 1659 • **Tag vs. Instr:**
1660 `tagcorrect_vs_instrexpert`.

- **Significance markers** appended to
correlation coefficients: $*p < 0.05$,
 $**p < 0.01$, $***p < 0.001$, and ns for not
significant.

1665 C.3.1 RQ3 Statistical Details

1666 We report Pearson correlations between paired
1667 differences Δ C1 (prompt Hoyer) and Δ C3 (en-
1668 tropy or semantic variation) for each model
1669 and prefix contrast. Two-sided p -values are
1670 computed per correlation; multiple compar-
1671 isons are controlled with Benjamini–Hochberg
1672 correction within each table family. Full re-
1673 sults are in Tables 7–8.

1674 C.3.2 Correlation Between C1 and 1675 C3-Entropy

1676 Table 7 reports Pearson correlation coeffi-
1677 cients between C1 (Δ Hoyer) and C3-Entropy
1678 (Δ Entropy).

1679 *Rounding.* Values correspond to `pearson_r` in
1680 the CSV; e.g., for Qwen3-8B under Instr vs.
1681 Ctrl, $r = 0.308859$ rounds to 0.309.

1682 C.3.3 Correlation Between C1 and 1683 C3-SemVar

1684 Table 8 reports Pearson correlation coefficients
1685 between C1 (Δ Hoyer) and C3-SemVar (Δ Var,
1686 i.e., `semantic_variation_1mF1`).

1687 D Reproducibility and Resource 1688 Consumption

1689 D.1 Hardware and Runtime

- 1690 • All experiments are conducted on multi-
1691 ple high-memory GPUs ($8 \times A100$). We
1692 disable gradient computation and only
1693 perform forward inference with activa-
1694 tion/attention caching;

Model	ΔAttn Tag-Ctrl (Mean \pm SEM)	ΔAttn Instr-Ctrl (Mean \pm SEM)
Llama-3.1-8B-Instruct	-0.0592 ± 0.0045	-0.0599 ± 0.0046
Mistral-7B-Instruct-v0.2	$+0.0158 \pm 0.0004$	$+0.0175 \pm 0.0005$
Qwen3-8B	-0.0493 ± 0.0039	-0.0480 ± 0.0038

Table 5: Mean \pm SEM of ΔAttn under Tag-Ctrl and Instr-Ctrl across three models (First-gen view). Negative values indicate reduced domain-keyword attention share relative to control.

Model	$r(\Delta\text{Attn}, \Delta\text{Entropy})$	p -value
Llama-3.1-8B-Instruct	0.117	1.3×10^{-10}
Mistral-7B-Instruct-v0.2	-0.072	7.8×10^{-5}
Qwen3-8B	0.229	5.4×10^{-37}

Table 6: Pearson correlation between ΔAttn and $\Delta\text{Entropy}$ (instance-level), with corresponding p -values.

- For a single model, running MFV and extracting the metrics on the RouterEval subset used in this paper takes on the order of several hours, mainly depending on the model size and the number of stochastic generations (K);
- The analysis scripts (C1-C3 computation and statistical tests) are primarily bottlenecked by I/O and aggregation, and typically take less time than forward inference.

D.2 Randomness and Reproducibility

- All stochastic components (instance shuffling, incorrect-tag sampling, and decoding-time sampling) use fixed random seeds;
- For stochastic metrics such as semantic variation, we report the variance across repeated runs in the appendix;
- When supported by the framework, we enable deterministic options as much as possible to reduce the impact of numerical nondeterminism on the results.

D.3 Code and Data Release Plan

During the anonymous review phase, we do not disclose the repository URL. If the paper is accepted, we plan to:

- Release the RouterEval preprocessing scripts and the routing-style meta-prompt templates;
- Release the code for C1-C3 computation and statistical analysis, including MFV scripts for multi-GPU parallelization;

- Provide configuration files and command-line examples to facilitate reproducing experiments and figures;
- Provide, within licensing constraints, the indices or hashes of the RouterEval subset used, enabling other researchers to compare results on the same set of instances.

E Extended Ethics Discussion

When deploying a similar analysis framework in real-world systems, the following ethical and safety considerations should be taken into account:

1. Potential manipulation risks A finer-grained understanding of how routing signals relate to internal representations could be misused to more precisely steer model behavior or to circumvent safety mechanisms. Our work is intended for diagnosis and robustness analysis; we do not endorse any application of RIDE or related techniques for deceptive, surveillance-oriented, or otherwise harmful purposes.
2. Bias and unfair routing If “densification-stability” signals are directly used as routing cues or as proxies for uncertainty, they may exhibit systematic differences across languages, demographic groups, or topics, thereby amplifying biases present in the training data. Future work should evaluate the fairness properties of RIDE metrics on multilingual and multi-group datasets.

Model	r (Instr vs. Ctrl)	r (Tag vs. Ctrl)	r (Tag vs. Instr)
Llama-3.1-8B-Instruct	0.141 ^{***}	0.013 ^{ns}	0.105 ^{***}
Mistral-7B-Instruct-v0.2	0.191 ^{***}	0.043 [*]	0.150 ^{***}
Qwen3-8B	0.309 ^{***}	0.204 ^{***}	0.139 ^{***}

Table 7: Pearson correlations between C1 (Δ Hoyer) and C3-Entropy (Δ Entropy) under three contrasts. Superscripts indicate significance: ^{***} $p < 0.001$, ^{**} $p < 0.01$, ^{*} $p < 0.05$, and ns not significant.

Model	Instr-Ctrl	Tag-Ctrl	Tag-Instr
Llama-3.1-8B-Instruct	0.067 ^{***}	-0.027 ^{ns}	0.037 [*]
Mistral-7B-Instruct-v0.2	0.142 ^{***}	0.128 ^{***}	0.177 ^{***}
Qwen3-8B	0.023 ^{ns}	0.072 ^{***}	0.017 ^{ns}

Table 8: Pearson r between C1 (Δ Hoyer) and C3-SemVar (Δ Var). Superscripts: ^{***} $p < 0.001$, ^{**} $p < 0.01$, ^{*} $p < 0.05$, ns not significant.

- 1761
- 1762
- 1763
- 1764
- 1765
- 1766
- 1767
- 1768
- 1769
- 1770
- 1771
- 1772
- 1773
- 1774
- 1775
- 1776
- 1777
- 1778
3. Privacy and data usage This paper uses only public benchmarks and open-source models, without involving real user data. If future work extends the framework to real routing logs or production data, strict anonymization should be applied, and relevant privacy regulations and platform policies must be followed.
 4. Cautious claims about “interpretability” We emphasize that RIDE provides a statistical lens for interpretation rather than a strict causal characterization of underlying mechanisms. In external communication or system documentation, it should not be presented as a “fully interpretable” solution; instead, its role as an analysis tool and auxiliary signal should be stated explicitly.