

Discovery report for Path Engineering: Causal Manipulation of the "Road Not Taken"

Research Objective

Run this idea for me:

Path Engineering: Causal Manipulation of the "Road Not Taken"

TL;DR: We'll test whether the correlation between uncertainty and path representation is causal by artificially creating or destroying alternate paths in the model's representation space. If the core paper is right, we should see uncertainty change accordingly when we engineer the path landscape.

Research Question: Can we causally manipulate language model uncertainty by directly controlling the availability of alternate reasoning paths in the model's representation space?

Hypothesis: Artificially constraining the model's path space (by reducing the dimensionality of hidden representations along critical directions) will decrease uncertainty, while expanding the path space will increase uncertainty, even when the reasoning problem itself remains unchanged.

Experiment Plan: Identify the low-rank subspaces that encode possible reasoning paths using the methods from Dai et al. (Binding ID mechanism). Develop "path constraint" interventions that reduce activation variance along these subspace directions, effectively limiting the number of representable paths. Conversely, create "path expansion" interventions that increase variance along these directions. Apply these interventions during chain-of-thought reasoning on fixed problems and measure changes in: token-level entropy (uncertainty), ease of steering (following Zur et al.'s methodology), and final reasoning accuracy. Expected outcome: Path constraints will reduce uncertainty and make the model harder to steer, while path expansion will increase uncertainty and steerability.

References: [Zur, A., Geiger, A., Lubana, E., & Bigelow, E.J. (2025). Are language models aware of the road not taken? Token-level uncertainty and hidden state dynamics.], [Dai, Q., Heinzerling, B., & Inui, K. (2024). Representational Analysis of Binding in Language Models. Conference on Empirical Methods in Natural Language Processing.]

Summary of Discoveries

Discovery 1: Low-Rank Subspace Identification and Control for Path Engineering

This discovery assembles an end-to-end toolkit for identifying low-rank subspaces that carry high-level computations in language models and for manipulating those subspaces to constrain or expand the set of representable reasoning paths. It also systematizes dimensionality-sensitive diagnostics—stable rank and entropy-based effective rank—and fills a methodological gap by proposing Marchenko–Pastur spike detection as a principled spectral validator for transformer activations.

Discovery 2: Quantifying Causal Effects of Path Manipulation on Uncertainty

This report proposes and operationalizes a causal test of whether language model uncertainty is controlled by the availability of alternate internal reasoning paths. By applying low-rank, subspace-targeted interventions that either constrain or expand a model's "path space" while holding the problem fixed, the study quantifies resulting changes in token- and sequence-level uncertainty, steerability, and accuracy.

Discovery 3: Causal Loci and Targeting: Attention-Head Circuits and Entropy-Modulating Neurons

This report localizes causal levers for "path engineering" in language models: specialized attention-head circuits in intermediate-to-late layers drive multi-step reasoning, while a small set of final-layer MLP units ("entropy neurons") modulate token-level predictive entropy via LayerNorm scaling.

Selective, token-position-aware intervention schedules consistently yield stronger, cleaner control than continuous steering, providing the when, where, and how to manipulate the availability of alternate reasoning paths and test their causal link to uncertainty.

Discovery 4: Entropy-Regularized Control as a Baseline for Uncertainty Manipulation

This discovery establishes a principled, non-representational baseline for manipulating language model uncertainty by importing entropy-regularized control from reinforcement learning to decoding. Maximum Entropy objectives and dual temperature tuning provide explicit mechanisms to increase or target token-level entropy, but the literature lacks a documented per-step update rule for the temperature during generation, creating a methodological gap for rigorous baseline comparisons.

Low-Rank Subspace Identification and Control for Path Engineering

Summary

This discovery assembles an end-to-end toolkit for identifying low-rank subspaces that carry high-level computations in language models and for manipulating those subspaces to constrain or expand the set of representable reasoning paths. It also systematizes dimensionality-sensitive diagnostics—stable rank and entropy-based effective rank—and fills a methodological gap by proposing Marchenko–Pastur spike detection as a principled spectral validator for transformer activations.

Background

Across mechanistic interpretability and representation engineering, converging evidence suggests that many high-level behaviors in language models are mediated by low-dimensional structures embedded in the residual stream and MLP outputs. Methods that localize these structures and apply targeted low-rank interventions now enable precise control over internal computations, offering a route to experimentally test how altering the “landscape” of available representational directions changes downstream behavior. Establishing reliable subspace discovery, intervention, and dimensionality diagnostics is therefore central to causally probing and engineering reasoning pathways.

Results & Discussion

The core finding is that reproducible pipelines now exist to (i) harvest chain-specific activations, (ii) isolate low-dependency subspaces associated with branching computations, and (iii) implement low-rank interventions that either remove access to those directions (path constraint) or add new, orthogonal directions (path expansion). Three complementary pipelines illustrate this. First, unsupervised subspace discovery via an orthogonal rotation that partitions the residual stream into low mutual-information components (e.g., NDM) enables subspace patching by projecting $\hat{h} = Rh$, swapping selected coordinates with a counterfactual, and inverting with R ; layer/token sites are chosen where a variable is present but not yet consumed, and ef-

fects are quantified by logit-difference or probability shifts ($\Delta LD/\Delta P$) with concentration measured by Gini > 0.6 and mutual information via KSG estimators [r10, huang2508]. Second, PCA/probe manifold methods compute a small set of principal components (e.g., $k = 5$) at a target layer/token, fit a low-dimensional “circular probe,” and intervene with a pseudoinverse update calibrated by off-manifold sweeps ($r \in [0, 2]$, θ grid), with early layers often most causal [r10, engels2024, liao2024]. Third, low-rank subspace steering (ReFT/DII) learns a rank- r orthonormal basis R and a projected affine map $\Phi(h) = h + R(Wh + b - Rh)$, with small ranks ($r \leq 8$; rank-1 often strong) trained with $lr \approx 4 \times 10^{-3}$ over ~ 1000 epochs; directional controls also include difference-in-means vectors for ablation $x = x - \hat{r}x$ or addition $x = x \pm \alpha r$, with layer selection constrained by $\text{induce}_{\text{score}} > 0$ and $KL < 0.1$ to preserve fluency [r10, wu2024, arditi2024].

These building blocks yield an operational path-constraint protocol for chain-of-thought tasks. Residual activations are collected token-wise across layers, prioritizing bottlenecks immediately before usage of a variable; mid-upper layers often emerge as effective sites (e.g., MSRS finds layer ≈ 15 across tasks), and an “Important token” can be selected by the magnitude of subspace projection [r10, jiang2025]. Subspaces can be identified by sequential NDM partitioning/merging with Gini-guided selection and low mutual information, or by difference-in-means/SVD, optionally complemented by sparse-autoencoder feature reconstruction; selection is validated by high patching Gini and low MI [r10, huang2508, wehner2025, engels2024]. Constraint is then applied either by projecting away the discovered subspace ($x = x - P_{\text{sub}}x$), by pseudoinverse manifold edits when local geometry is known, or by a low-rank ReFT Φ at selected layers/tokens; intervention strength is tuned on held-out prompts via grids over r or α under $\Delta LD/\Delta P$ improvement and a $KL < 0.1$ guard to limit distributional drift [r10, engels2024, wu2024, arditi2024]. Evidence from GSM8K-like pipelines suggests that

branch-specific variables often concentrate into 64-D subspaces, and that rank-r = 8 ReFT at the Important token in layer ≈15 matches R-based patching on ΔLD while inducing lower KL, consistent with gated, orthonormal low-rank edits [r10, jiang2025, wu2024].

Path expansion is supported by replicable methods that derive multiple, distinct steering vectors and enforce their approximate orthogonality, enabling linear combinations that open additional reasoning routes. Contrastive activation differences (paired positive/negative sets), sparse autoencoder-based features, and PCA/SVD all yield low-dimensional bases whose columns are iteratively orthogonalized by projecting new candidates onto the complement of the current subspace; cosine similarity and inner products verify distinction [r13, yu2510, wehner2025]. In parallel, widely used protocols explicitly form contrastive subspaces by computing set-wise differences $A_+ - A_-$ and applying PCA/SVD to extract principal contrast directions for steering or suppression (e.g., RepE’s Linear Artificial Tomography, LoRRA, and low-rank affine editors), while contrastive PCA frameworks (CPCA/PCPCA) add statistical grounding and principal-angle tests for contrastive dimensionality [r33, bartoszcz2025, wehner2025, hawke2024]. Mean-difference vectors (DiM/CAA) commonly match or exceed PCA on target identification, so both DiM and PCA/CPCA bases are useful for constructing expandable, orthogonal steering sets [r33, wehner2025, bartoszcz2025].

The discovery also consolidates dimensionality-sensitive diagnostics to quantify how interventions compress or expand the accessible path space. Two complementary measures are emphasized. The stable rank is $\text{srank}(A) = \|\mathbf{A}\| \|\mathbf{F}^2\| / \|\mathbf{A}\|^2$, where the numerator is the sum of squared singular values and the denominator is the largest squared singular value; the entropy-based effective rank computes $p = \sigma^2 / \sum \sigma^2$ followed by $\exp(-\sum p \log p)$ to summarize spectral concentration [r28, balcan2019, gu2015]. Best practice for building the activation matrix \mathbf{A} is to select representative token activations from relevant layers (e.g., post-GELU MLP outputs or residual stream at mid-late depth), con-

struct tokens×dims matrices, and apply centering before SVD/PCA [r28, chang2025, simionato2023]. Empirically, intervention-driven dimensionality shifts are measurable: safety fine-tuning concentrates unsafe variance into a single dominant direction in deep MLPs ($\approx 62\%$ of nuclear norm) while safe remains diffuse ($\approx 12\%$), indicating reduced effective/stable rank; removing attention-sink “massive activations” increases spectral entropy/effective rank and reduces top-singular-value dominance; and detached-Jacobian analyses report stable-rank changes that track the strength and locus of layer-level steering [r36, jain2024, queipodellan2025, golden2025]. These metrics therefore provide quantitative readouts of path-space compression (constraint) and diversification (expansion).

Finally, the sources identify a methodological gap and a solution. Despite extensive use of random matrix theory on weight spectra and overlaps with activation eigenvectors, formal, data-driven rank-estimation procedures have not been applied directly to transformer activation covariances in these settings [r68, staats2024]. Yet explicit, step-by-step Marchenko–Pastur spike-counting algorithms exist: biwhiten the data (e.g., diagonal Sinkhorn–Knopp scaling) to address heteroskedasticity, compute the sample covariance and aspect ratio γ , estimate bulk edges $\beta \pm = \sigma^2(1 \pm \sqrt{\gamma})^2$ (or equivalently singular-value thresholds $\sqrt{n} \pm \sqrt{m}$), and count eigenvalues beyond the upper bulk edge as signal dimensions, with robustness refinements for heavy tails [r75, landa2022, martin2021, ali2025a, vallet2015]. Integrating MP-based spike detection alongside stable/entropy ranks supplies a principled spectral validator for rank selection on centered activation covariances, closing the validation gap and enabling standardized, quantitative assessment of how path-constraint and path-expansion interventions reshape the model’s representational degrees of freedom [r28, r36, r68, r75].

Trajectory Sources

Trajectory r10: Yes—the literature provides end-to-end, reproducible building blocks to design a concrete path-constraint protocol that (i) collects divergent-chain trajectories, (ii) identifies a low-dependency subspace, and (iii) constrains activations via an explicit projection/intervention operator at selected ...

Trajectory r13: The literature supports that replicable unsupervised and semi-supervised methods exist to identify multiple, distinct, and approximately orthogonal steering vectors whose linear combinations can serve as effective path expansion interventions for reasoning tasks ([yu2510](#) pages 2-...)

Trajectory r28: The literature confirms that stable rank is defined as $\text{srank}(A) = \|A\|F^2/\|A\|^2$ and that an entropy-based effective rank is obtained by forming a probability distribution from the squared singular values and then applying Shannon's entropy, while best practices for activation matrix A call for sele...

Trajectory r33: Yes—the literature contains multiple protocols that explicitly form contrastive subspaces by applying PCA/SVD (or related decompositions) to matrices of activation differences (including difference-of-means) between opposing concept conditions, and these directions are then injected at inference to ...

Trajectory r36: Yes—at least three studies quantify intervention-driven changes in transformer activation dimensionality with singular-value-based effective/stable-rank measures, and they specify which layers to probe, how to build activation matrices, and how to interpret rank shifts. ([jain2024](#) pages 6...)

Trajectory r68: Based on the provided sources, I find no study that applies a formal, data-driven rank-estimation procedure (RMT thresholding, MDL, or Bayesian rank selection) directly to transformer activation covariance matrices; RMT is used on weight spectra with overlaps to activation-covariance eigenvectors, b...

Trajectory r75: The literature confirms that explicit, step-by-step algorithms with

complete mathematical formulas for applying Marchenko–Pastur analysis to estimate the rank of an empirical covariance matrix are indeed available ([landa2022](#) pages 3-5, [landa2022](#) pages 1-3).

Quantifying Causal Effects of Path Manipulation on Uncertainty

Summary

This report proposes and operationalizes a causal test of whether language model uncertainty is controlled by the availability of alternate internal reasoning paths. By applying low-rank, subspace-targeted interventions that either constrain or expand a model’s “path space” while holding the problem fixed, the study quantifies resulting changes in token- and sequence-level uncertainty, steerability, and accuracy.

Background

Uncertainty in language models has been tied to the diversity of latent reasoning trajectories, but most support to date is correlational. Advances in representation engineering now allow precise, low-rank edits to internal activations, providing a plausible intervention to alter the “landscape” of possible reasoning paths without retraining. In parallel, uncertainty quantification for chain-of-thought has matured, with reliable single-pass metrics and semantically informed, sampling-based measures. The confluence of these tools creates a timely opportunity to move from correlation to causation by directly manipulating internal path availability and measuring the downstream effects on uncertainty and reasoning behavior.

Results & Discussion

The discovery formalizes a test of the hypothesis that token-level uncertainty is causally governed by the internal diversity of reasoning paths: “path constraints” should lower uncertainty and reduce steerability, whereas “path expansions” should elevate uncertainty and increase steerability under otherwise fixed tasks [r0]. Prior work indicates that high-level variables concentrate in low-rank subspaces and can be manipulated via steering vectors and binding-identifier edits to predictably influence outputs and uncertainty without full retraining, motivating a causal intervention on the path subspaces themselves [r0]. The gap this fills is a systematic, controlled link from engineered changes in the latent path landscape to measured changes in

token-level entropy, steering ease, and accuracy, rather than relying on correlational associations [r0].

The protocol assembles reproducible building blocks for both subspace discovery and intervention. Unsupervised subspace identification proceeds by learning an orthogonal rotation that partitions the residual stream into low-mutual-information subspaces (NDM), then patching by projecting activations, replacing selected coordinates at layers and tokens where variables are present but not yet consumed; effects are summarized by changes in logit difference and probability ($\Delta LD/\Delta P$), with concentration scored by Gini and dependence by KSG mutual information, using documented thresholds (e.g., merge 0.015–0.04; cosine preimage $0.85 \rightarrow 0.7$) and search settings (e.g., N and steps) [r10]. Alternative manifolds are accessed with PCA/probe pseudoinverse edits (e.g., k=5 PCs into a $2 \times k$ probe with off-manifold sweeps r [0, 2] and θ -grids), or with low-rank operators such as ReFT, where small rank r (often 8; rank-1 can be strong) is trained or calibrated to gate edits at chosen layers/tokens with modest learning rates (e.g., 4×10^{-3}) [r10]. Directional ablations and additions use difference-in-means vectors with layer selection guarded by $\text{induce}_{\text{score}} > 0$ and $\text{KL} < 0.1$ against neutral data, and strength is tuned by grids while enforcing KL constraints [r10]. For path expansion, unsupervised and semi-supervised procedures derive multiple, approximately orthogonal steering vectors by contrastive activation differences, sparse autoencoder features, or PCA/SVD; these directions are validated by cosine/inner-product tests and can be linearly combined to span distinct strategies, enabling controlled increases in path diversity [r13]. Together, these pieces instantiate “path constraint” (variance reduction or projection onto an orthogonal complement of path-encoding directions) and “path expansion” (addition/combination of orthogonal path vectors) interventions that can be applied at bottleneck layers or “Important tokens” identified

by subspace-projection magnitude [r10].

Uncertainty is measured at both token and sequence levels using complementary metrics with established reliability. Predictive entropy is adopted as the primary single-pass, logit-based indicator due to its consistent superiority over max-logit or margin as a first-stage filter for correctness and output variability [r30, yang2025]. To capture full-distribution effects in chain-of-thought, the analysis augments predictive entropy with logit variance, semantic entropy (i.e., entropy over meaning-level clusters), and multi-sample disagreement, which collectively provide a stronger and more robust signal than next-token entropy alone in reasoning settings [r6, abbasli2025, zhang2025c]. Token-level metrics are aggregated into sequence-level scores via mean, quantile (e.g., max or high-quantile), and cumulative cross-entropy operators; these aggregation schemes are empirically validated for predicting correctness and deferral, albeit with context-dependent trade-offs that motivate testing multiple operators in parallel [r42, gupta2024, sharma2510]. To quantify realized path diversity under fixed prompts, the protocol generates many samples (potentially with CoT branching or speculative decoding), clusters outputs into semantic equivalence classes using entailment and embeddings with practical dedup thresholds (e.g., cosine 0.8), and, when available, uses execution-based validators (e.g., unit tests) to count distinct valid solutions via pass@n and “different implementation” measures; the resulting semantic entropy and cluster counts track the diversity induced by interventions [r25, foodeei2506, farquhar2024, chen2024, havrilla2024]. Benchmarks such as GSM8K, logic grids, Sudoku, and graph puzzles contain multiple valid reasoning paths and thus provide appropriate testbeds for mapping engineered path landscapes to uncertainty and accuracy outcomes [r7].

The expected causal signature is bidirectional and testable. Under path-constraint interventions (e.g., projection that removes identified path-encoding components, or stronger low-rank gating at a critical layer), token-level predictive entropy and sequence-level aggregates should decrease, semantic cluster counts should shrink, and steering ease (quantified

by $\Delta\text{LD}/\Delta\text{P}$ under held-out steering prompts) should drop; conversely, path-expansion interventions that add or combine orthogonal steering vectors should increase these quantities [r0]. Because representation edits are sensitive to strength, layer choice, and sample selection—with documented trade-offs between accuracy and fluency—the protocol calibrates intervention strength by grid sweeps (e.g., radius and angle on probe manifolds), keeps KL divergence below 0.1 on neutral data, and compares low-rank ReFT (rank r 8 at an Important token around a mid-upper layer, often ≈ 15) against subspace patching for matched ΔLD at lower distributional shift [r10]. Demonstrating that engineered reductions (or expansions) of representable paths reliably decrease (or increase) uncertainty on fixed problems would move the field from correlational observations to a causal account of the “road not taken,” while establishing a practical knob for calibrating deferral policies and steering behavior in complex reasoning systems [r0, r42].

Trajectory Sources

Trajectory r0: Path engineering aims to causally control a language model’s uncertainty by directly manipulating the availability of alternate reasoning paths in its latent representation space ([bartoszce2025](#) pages 12-14). Prior work has shown that internal activations encode high-leve...

Trajectory r6: Our findings support the hypothesis that uncertainty metrics leveraging the full output distribution—such as semantic entropy, logit variance, and disagreement measures from multiple sampled decoding paths—yield a significantly more robust signal in chain-of-thought reasoning than next-token entropy...

Trajectory r7: Existing chain-of-thought reasoning benchmark datasets, including GSM8K, logic grid puzzles, and graph-based as well as game puzzles, demonstrably contain problems with multiple valid reasoning paths, confirming their suitability as a testbed for our experiment ([chen2025b](#) pages...)

Trajectory r10: Yes—the literature provides end-to-end, reproducible building blocks to design a concrete path-constraint protocol that (i) collects divergent-chain trajectories, (ii) identifies a low-dependency subspace, and (iii) constrains activations via an explicit projection/intervention operator at selected ...

Trajectory r13: The literature supports that replicable unsupervised and semi-supervised methods exist to identify multiple, distinct, and approximately orthogonal steering vectors whose linear combinations can serve as effective path expansion interventions for reasoning tasks ([yu2510](#) pages 2-...)

Trajectory r25: The hypothesis is supported: there are established, largely automatic pipelines that generate many samples for a single prompt and use semantic clustering/equivalence tests plus task validators to count distinct, valid solutions or reasoning paths.

Trajectory r30: Predictive entropy consistently emerges as the most reliable logit-based predictor of task correctness and semantic diversity, making it the prime candidate for a

first-stage uncertainty filter.

Trajectory r42: The literature supports aggregating token-level uncertainty into sequence-level scores using operators such as the mean, maximum (or extreme quantiles), and cumulative measures (e.g., sum of negative log probabilities), with empirical evidence showing these methods can predict correctness and variab...

Causal Loci and Targeting: Attention-Head Circuits and Entropy-Modulating Neurons

Summary

This report localizes causal levers for “path engineering” in language models: specialized attention-head circuits in intermediate-to-late layers drive multi-step reasoning, while a small set of final-layer MLP units (“entropy neurons”) modulate token-level predictive entropy via LayerNorm scaling. Selective, token-position-aware intervention schedules consistently yield stronger, cleaner control than continuous steering, providing the when, where, and how to manipulate the availability of alternate reasoning paths and test their causal link to uncertainty.

Background

Mechanistic studies increasingly suggest that language models’ uncertainty and reasoning behaviors arise from structured, editable internal circuits rather than diffuse, uniform computation. Causal tracing and activation patching have mapped hierarchical roles to attention heads and MLPs, while targeted activation-space steering has revealed the importance of intervention timing. In parallel, a growing body of work has identified dedicated neurons that regulate predictive entropy without necessarily altering the model’s top choice. Together, these advances outline a principled strategy to intervene on the “road not taken”—the set of latent alternatives the model can represent—and to test whether uncertainty is causally determined by the richness of these representational paths.

Results & Discussion

Causal analyses converge on specialized attention-head circuits in intermediate-to-late layers as the primary drivers of multi-step reasoning, with MLPs contributing secondary or task-specific computations. Across arithmetic and logical tasks, activation patching shows that perturbing specific attention heads—often in the mid-to-late stack (e.g., layer 10 in a 12-layer model)—produces the strongest causal effects on reasoning performance, whereas final MLP blocks, while involved in computations such as modular addition, are generally less critical than attention-mediated retrieval and induc-

tion operations [r41, goyal2025, zhang2510, altabaa2510, hong2024]. Although some recurrent transformer settings assign aggregation roles to late MLP layers, the bulk of controlled interventions still implicates attention circuits as the decisive loci for multi-step reasoning, suggesting that path-formation mechanisms are predominantly attention-driven [r41, hong2024].

The timing and scope of interventions prove equally important: selective, token-position-aware steering achieves stronger, more stable control than continuous every-token interventions. Comparative studies that target specific prompt or generation positions (e.g., a single appended-instruction token, all verbs in a prompt versus only the last verb) demonstrate that intervention location and duration modulate both steering strength and side effects such as topic shift or degeneration; continuous interventions maintain influence but increase degeneration risk [r46, chang2508, klerings2025]. Notably, prior work has not dynamically identified “reasoning-step” tokens, but the consistent sensitivity to timing underscores that control should be scheduled at task-relevant positions rather than applied uniformly across all tokens [r46, klerings2025].

Independently, final-layer MLP “entropy neurons” exert causal control over predictive uncertainty by regulating the final LayerNorm scale. These units often have large output weight norms yet minimal direct logit contributions; clamping or ablating their activations reliably increases or decreases token-level entropy without dramatically changing the top prediction. Moreover, interventions on certain induction heads can modulate these entropy signals indirectly by altering head output norms that feed into LayerNorm, yielding coupled attention–MLP control over uncertainty. These phenomena have been validated by large-scale activation–uncertainty correlations and cross-model replications, establishing a robust link between specific neurons/heads and predictive entropy during reasoning [r74,

[gurnee2024](#), [stolfo2024a](#), [choi2510](#)]. In this context, token-level entropy refers to the Shannon entropy of the next-token distribution, and the LayerNorm-mediated mechanism connects upstream attention dynamics to the final uncertainty readout [[r74](#), [stolfo2024a](#)].

Synthesizing these findings yields an actionable blueprint for path engineering to test the causal relationship between alternate-path availability and uncertainty. First, engineer the path space by intervening on the attention-head circuits that causally drive multi-step reasoning: constrain paths by reducing activation variance along low-rank directions that encode candidate reasoning trajectories within mid-to-late attention layers, and expand paths by injecting variance or reinstating multiple candidate directions. Second, schedule these interventions at task-relevant tokens to maximize control and minimize degeneration (e.g., before key reasoning steps or immediately prior to decisive outputs), rather than at every token [[r41](#), [r46](#), [hong2024](#), [klerings2025](#)]. Third, separately modulate final-layer entropy neurons to adjust the uncertainty readout without necessarily altering the top prediction, enabling a clean test of whether changes in path availability alone shift token-level entropy and steering properties [[r74](#), [gurnee2024](#), [stolfo2024a](#)]. The primary metrics are: token-level entropy (Shannon entropy of the next-token distribution), steering strength/efficiency (change in target behavioral or representational metrics per unit intervention magnitude), side-effect rates (topic shift/degeneration frequency), and end-task accuracy. Based on prior causal evidence, constraining attention-path subspaces should reduce entropy and lower steering efficiency, whereas expansion should increase entropy and steering efficiency; jointly tuning entropy neurons provides a sensitive, orthogonal handle on the uncertainty channel to isolate causal effects. Methodological caveats remain—causal roles can vary across tasks and training regimes, and dynamic identification of reasoning-step tokens is still an open challenge—but the attention–MLP coupling and token-aware scheduling identified here specify where, when, and how to manipulate the “road not taken” to establish a causal link to uncertainty [[r41](#), [r46](#), [r74](#), [klerings2025](#), [choi2510](#)].

Trajectory Sources

Trajectory r41: The literature shows that while intermediate to mid-late layers are critical for multi-step reasoning, causal interventions consistently implicate specialized attention head circuits—not exclusively middle-to-late MLP blocks—as the primary drivers of reasoning performance on benchmarks like GSM8K (p...)

Trajectory r46: The literature demonstrates that the timing and scope of activation-space interventions are critical—studies consistently find that targeted, selective application (rather than continuous every-token intervention) tends to yield more controlled outputs with fewer side effects. ([chang2508unveilingthe...](#)

Trajectory r74: Studies consistently show that specific final-layer MLP neurons (often termed “entropy neurons”) and distinct attention heads reliably correlate with and causally influence predictive uncertainty during reasoning tasks ([gurnee2024](#) pages 8-10).

Entropy-Regularized Control as a Baseline for Uncertainty Manipulation

Summary

This discovery establishes a principled, non-representational baseline for manipulating language model uncertainty by importing entropy-regularized control from reinforcement learning to decoding. Maximum Entropy objectives and dual temperature tuning provide explicit mechanisms to increase or target token-level entropy, but the literature lacks a documented per-step update rule for the temperature during generation, creating a methodological gap for rigorous baseline comparisons.

Background

Uncertainty in language models is often studied through the lens of hidden-state dynamics and representational capacity, yet decoding-time control of uncertainty provides an orthogonal lever that does not alter internal representations. Entropy-regularized control from reinforcement learning offers mature, mechanistic tools for steering policy stochasticity, suggesting a direct pathway to regulate token-level entropy during generation. Establishing such a baseline is essential to disentangle representational causes of uncertainty from generic entropy steering when testing whether changes to the availability of alternate reasoning paths causally modulate uncertainty.

Results & Discussion

The central finding is that well-studied entropy-aware control frameworks provide explicit, mechanistic objectives for increasing or targeting policy entropy that translate directly to token-level uncertainty control in language models. In Maximum Entropy reinforcement learning, the policy objective augments rewards with an entropy bonus $H(\pi) = E[-\log \pi]$, leading to optimization of $J(\pi) = E[\sum(r_t + \alpha H(\pi(\cdot|s_t)))]$, soft value recursions, and policy updates that explicitly include the term $-\alpha \log \pi(a|s)$ inside the expectation [r70, han2021, tao2024, liu2019]. Soft Actor-Critic implements this by alternating soft policy evaluation and improvement, maximizing $E[Q(s, a) - \alpha \log \pi(a|s)]$ while learning a temperature α as a Lagrange dual to

enforce a target entropy via minimizing $E[-\alpha \log \pi(a|s) - \alpha H_0]$ [r70, han2021, kim2023a]. This maps one-to-one to language modeling: steer the token distribution by adding a per-token entropy bonus $-\alpha \log p\theta(y_t|context)$ or by dynamically adapting α to meet a target token-level entropy during generation [r70, tao2024, kim2023a].

Beyond SAC, related formulations broaden the control toolkit and remain compatible with LM decoding. Maximum-Minimum Entropy variants modify Bellman targets to emphasize low-entropy states while keeping an actor that still ascends $E[Q_R - \alpha \log \pi]$, with practical mechanisms such as double-Q minima, EMA targets, and reward scaling [r70, han2021]. Proximal control based on f-divergence regularization optimizes $E[R] - \eta D_f(\rho_\pi \rho_{\pi^0})$, yielding closed-form advantage-weighted reweighting updates and trust-region behavior that can be used to shape exploration jointly with an explicit $-\alpha \log \pi$ penalty [r70, belousov2019]. Convex occupancy-measure programs add entropic penalties with weight η and give dual optimality conditions with a Bregman divergence structure, reinforcing the theoretical grounding for entropy-regularized control [r70, neu2017]. Importantly, the policy-level regularizer $-\alpha \log \pi$ is directly applicable to discrete-token LMs without a value critic, though stability considerations for temperature adaptation still apply [r70].

Notwithstanding this strong conceptual basis, the surveyed excerpts do not provide a fully specified, per-step update rule for α during autoregressive decoding, nor a concrete equation that couples α to logit scaling based on instantaneous entropy errors. While SAC-style dual tuning suggests minimizing $J(\alpha) = E[-\alpha \log \pi - \alpha H_{target}]$ and standard temperature sampling rescales logits as $\exp(z/T)$, the exact gradient update for α and its integration with decoding remain unspecified in the reviewed material [r77, huang2025a, zhang2403, kang2510]. This omission marks a baseline-method gap: to serve as

a rigorous comparator to representational interventions, the entropy-control baseline requires an explicit, stepwise α update and a precise logit modification rule tied to token-level entropy. Two concrete hypotheses emerge: implement a dual update defined by $L(\alpha) = \alpha \cdot (H_0 - \text{current entropy})$ with gradient descent, and couple α to logits via $z' = z/\alpha$ during generation [r77].

As a baseline for testing causal claims about path-space manipulation, entropy-regularized control offers a powerful, non-representational lever to adjust uncertainty while holding the hidden-state geometry fixed. Prior evidence suggests that automatic dual temperature tuning to a target token entropy increases diversity at fixed perplexity relative to fixed-temperature sampling, and that f-divergence proximal regularization can deliver controllable exploration with lower repetition rates than KL-only proximal baselines, providing concrete readouts for evaluation such as diversity–perplexity trade-offs and repetition metrics [r70, kim2023a, belousov2019]. In the context of Path Engineering, matching the baseline’s target entropy and perplexity while observing differences in token-level entropy trajectories, steerability, or final accuracy would isolate effects attributable to altering the availability of alternate reasoning paths, rather than to generic entropy reweighting [r70]. Establishing this entropy-control baseline therefore enables a clean causal test of whether engineering the “road not taken” in representation space uniquely modulates uncertainty beyond what can be achieved by decoding-time entropy control alone.

Trajectory Sources

Trajectory r70: Yes—the surveyed entropy-aware control literature provides explicit, mechanistic objectives and update rules (e.g., max-entropy RL, SAC, MME, f-divergence proximal control) that specify how to compute and inject an entropy term into policy/value losses, and these are directly suggestive of steering ...

Trajectory r77: The literature supports the concept of dual temperature tuning for entropy control, but none of the provided excerpts specify an explicit per-step gradient update rule for α nor a concrete equation for modifying logits using this parameter (huang2025a pages 11-14, zhang2403e...)