Dynamic Manifold Evolution Theory: Modeling and Stability Analysis of Latent Representations in Large Language Models

Anonymous ACL submission

Abstract

We introduce Dynamic Manifold Evolution Theory (DMET), a unified framework that models large-language-model generation as a controlled dynamical system evolving on a low-dimensional semantic manifold. By casting latent-state updates as discrete-time Euler approximations of continuous dynamics, we map intrinsic energy-driven flows and contextdependent forces onto Transformer components (residual connections, attention, feedforward networks). Leveraging Lyapunov stability theory We define three empirical metrics (state continuity, clustering quality, topological persistence) that quantitatively link latenttrajectory properties to text fluency, grammaticality, and semantic coherence. Extensive experiments across decoding parameters validate DMET's predictions and yield principled guidelines for balancing creativity and consistency in text generation.

1 Introduction

003

011

012

014

027

034

042

Large Language Models (LLMs) have achieved revolutionary advances in recent years, from GPT-4 (OpenAI, 2023), LLaMA-3 (Touvron et al., 2024) to Claude (Anthropic, 2024), demonstrating unprecedented capabilities in language understanding and generation. However, despite their abundant applications, their internal mechanisms remain largely opaque, functioning as "black boxes" (Bommasani et al., 2022). This opacity constitutes a critical barrier to further improving the reliability, interpretability, and safety of LLMs. In particular, our understanding of how models organize and evolve their latent representations during the generation process remains limited—a knowledge gap that hinders our ability to effectively address core challenges such as hallucinations (Huang et al., 2023), inconsistencies (Zheng et al., 2023), and semantic drift (Shi et al., 2024).

Recent research has attempted to unveil the internal mechanisms of LLMs through various analytical approaches. Methods such as attention visualization (Vaswani et al., 2023), feature attribution (Sundararajan et al., 2022), and probing techniques (Liu et al., 2023) have revealed static properties of model representations, while the residual stream analysis by Elhage et al. (2021) and mechanistic interpretability research by Anthropic (2022) have begun to explore the dynamic aspects of cross-layer information propagation. However, these efforts largely provide localized or fragmented perspectives, lacking a unified theoretical framework that can describe the temporal evolutionary characteristics of the generation process. Traditional views simplify LLM generation as a concatenation of discrete token predictions, neglecting the continuous evolutionary dynamics in the latent space, which limits our understanding of how models gradually refine initial concepts into coherent text.

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

This paper proposes the **Dynamic Manifold Evolution Theory** (DMET), an innovative mathematical framework that reconceptualizes the LLM generation process as a dynamical system evolving on high-dimensional manifolds. Our key insight is that LLM generation is essentially a continuous process of latent representation evolution along semantic manifold trajectories, gradually refining macroscopic semantic concepts into specific textual expressions. By integrating dynamical systems theory, manifold geometry, and deep learning, DMET provides a rigorous mathematical foundation for understanding and optimizing the dynamics of internal representations in LLMs.

The main contributions of this paper are as follows: We propose the Dynamic Manifold Evolution Theory, which, for the first time, conceptualizes the LLM generation process as a dynamical system evolving on high-dimensional manifolds and establishes a rigorous mathematical link between latent representation evolution and generated text quality. We develop both continuous-time and discretized dynamical system models, provide ex-

182

183

184

plicit mappings to Transformer architectures, and introduce a comprehensive toolkit for analyzing internal state evolution. Leveraging Lyapunov the-086 ory, we prove convergence conditions for latent dynamics and establish a theoretical foundation connecting semantic consistency of generated text with dynamical stability, thereby offering princi-090 pled strategies for mitigating hallucinations and inconsistencies. Furthermore, we pioneer geometric and topological optimization approaches-such as curvature regularization and topology simplification-to address the challenges posed by complex high-dimensional geometry and improve the stability of generation. Finally, we validate our theoretical framework through extensive experiments and advanced visualizations, demonstrating strong correlations between latent trajectory properties and 100 output quality, and highlighting the critical role of 101 attractor structures in the generative process. 102

> The remainder of this paper is organized as follows: Section 2 reviews related work and introduces necessary preliminaries; Section 3 details the mathematical foundations and implementation methods of the Dynamic Manifold Evolution Theory; Section 4 describes the experimental design and analyzes results; and finally, Section 5 summarizes the main findings of this research and discusses directions for future work. Through this comprehensive framework, we not only deepen our understanding of internal LLM mechanisms but also provide theoretical guidance for designing more reliable and controllable next-generation language models.

2 Related Work

103

104

107 108

109

110

111

112

113

114

115

116

117

118

119

120

121

We ground DMET at the intersection of three research strands: dynamical systems in deep learning, manifold-based representation analysis, and latent trajectory modeling in language models.

Dynamical Systems in Neural Networks In-122 terpreting deep networks as discretized contin-123 uous systems has gained traction since Neural 124 ODEs (Chen et al., 2018a), which view residual 125 connections as Euler steps. Extensions include 126 augmented neural differential equations (Dupont 127 et al., 2019) and stability analyses for recurrent and 129 feed-forward architectures (Miller and Hardt, 2019; Santos et al., 2023; Li et al., 2023). In the language 130 domain, Lu et al. (2023) and Patil et al. (2024) ana-131 lyze Transformer dynamics, while Zhang and Xiao 132 (2024) frame decoding as a Markov decision pro-133

cess. Unlike these localized or task-specific studies, DMET provides a unified mapping from continuous dynamics (with Lyapunov stability) to all core Transformer components.

Manifold Geometry and Topology The manifold hypothesis posits that high-dimensional representations lie on low-dimensional structures (Roweis and Saul, 2000; Tenenbaum et al., 2000). Deep manifold learning methods include Riemannian metric estimation (Arvanitidis et al., 2018) and neural tangent space analysis (Chen et al., 2018b). In NLP, latent geometry has been explored via syntactic probes (Hewitt and Manning, 2019), linear subspace visualizations (Reif et al., 2019), and hierarchical manifold discovery in GPT (McCoy et al., 2022). Topological tools such as persistent homology (Liu et al., 2024; Dai et al., 2023) reveal global structural features. DMET leverages these geometric and topological insights to define dynamic trajectory metrics that directly link manifold structure to generation quality.

Latent Trajectory Analysis in Language Models Examining how hidden states evolve during text generation has illuminated RNN behavior (Li et al., 2016; Mardt et al., 2018) and Transformer residual streams (Elhage et al., 2021). Recent work investigates trajectory bifurcations (Rajamohan et al., 2023) and "thought manifold" evolution (Hernandez-Garcia et al., 2024). However, these analyses typically focus on visualization or specific phenomena. In contrast, DMET systematically models latent evolution as a dynamical system, quantifies its properties, and empirically correlates them with text fluency, grammaticality, and coherence.

By unifying continuous-time theory, manifold geometry, and trajectory analysis, DMET offers the first end-to-end framework for interpreting and controlling LLM generation dynamics. Our Dynamic Manifold Evolution Theory (DMET) uniquely integrates dynamical systems, Lyapunov stability, and manifold learning for a unified interpretation of LLM latent representation evolution. DMET differs from previous works by: (1) modeling representations as evolving points on dynamic manifolds rather than static vectors; (2) applying differential equations to model continuous evolution; (3) using Lyapunov stability to link representational stability and text quality; and (4) proposing geometric regularization for active optimization of latent manifold geometry and topology.



Dynamic Manifold Evolution Theory (DMET)

Figure 1: Overview of the DMET framework: latent trajectories evolve on a low-dimensional semantic manifold under intrinsic energy gradients and context-driven forces, with discrete Transformer layers implementing Euler steps of this continuous dynamics.

Dynamic Manifold Evolution Theory 3 (DMET): Methodology

3.1 Framework Overview

185

187

190 191

192

193

195

197

207

In this section, we give an overview of the Dynamic Manifold Evolution Theory (DMET): we begin by stating its three foundational assumptions, then show how these lead to a continuous-time dynamical model, and finally explain how residual connections, self-attention, and feed-forward layers implement that model in a Transformer.

3.2 Three Core Assumptions

Dynamic Manifold Evolution Theory (DMET) fun-196 damentally reinterprets the generation process of large language models (LLMs). Unlike traditional perspectives that simplify text generation as a se-199 quential prediction of discrete tokens, DMET conceptualizes it as a continuous trajectory evolution within a structured semantic space. This section elaborates on the three core assumptions that underpin this theoretical framework. The Dynamic 204 Manifold Evolution Theory (DMET) is grounded in three core assumptions that shape our understanding of LLM internal dynamics.We summarize

DMET's theoretical foundation in three concise pillars:

208

209

210

211

212

213

214

215

216

217

218

219

220

221

- 1. Manifold Structure. The hidden state $\mathbf{h} \in \mathbb{R}^d$ always lies on a much lowerdimensional semantic manifold $\mathcal{M} \subset \mathbb{R}^d$, with $\dim(\mathcal{M}) \ll d$.
- 2. Continuous Evolution. Text generation corresponds to a continuous trajectory $\mathbf{h}(t)$ smoothly traversing \mathcal{M} over time.
- 3. Attractor Landscape. The manifold \mathcal{M} contains multiple attractor basins $\{A_i\}$ —each representing a coherent semantic state-and $\mathbf{h}(t)$ naturally converges into one of these basins.

Core Assumptions of DMET. We build DMET 222 on three intertwined hypotheses. First, although 223 LLMs operate in a high-dimensional hidden space 224 \mathbb{R}^d , their meaningful representations lie on a much 225 lower-dimensional semantic manifold $\mathcal{M} \subset \mathbb{R}^d$, 226 capturing the regularities of language. Second, text 227 generation is not a series of independent jumps but 228 a smooth trajectory $\mathbf{h}(t)$ that continuously traverses 229



Figure 2: The curved surface represents the lowdimensional semantic manifold where latent representations exist. Rather than occupying the entire highdimensional space, meaningful linguistic representations are constrained to this manifold.

 \mathcal{M} from an initial state $\mathbf{h}(0)$ to a final state $\mathbf{h}(T)$, much like a hiker following a well-defined ridge rather than a random path. Third, \mathcal{M} is shaped by multiple attractor basins $\{\mathcal{A}_i\}$ —each corresponding to a coherent semantic frame—so that once the latent state enters an attractor's domain, it naturally converges to a stable region, explaining why LLMs generate focused, logically connected passages instead of disjoint word sequences.

3.3 Dynamical System Modeling and Transformer Mapping

236

238

240

241

242

245

246

247

248

250

251

254

We model latent evolution as a controlled dynamical system on the semantic manifold:

$$\frac{d\mathbf{h}(t)}{dt} \;=\; -\nabla V \big(\mathbf{h}(t) \big) \;+\; g \big(\mathbf{h}(t), \mathbf{u}(t) \big),$$

where V is an energy potential encoding semantic coherence, and g is a context-driven force from input $\mathbf{u}(t)$. Discretizing via the explicit Euler method with step size Δt gives

$$\mathbf{h}_{t+1} = \mathbf{h}_t + \Delta t \left[-\nabla V(\mathbf{h}_t) + g(\mathbf{h}_t, \mathbf{u}_t) \right].$$

Remarkably, each term aligns with a core Transformer component:

3.4 Dynamic Metrics (Aligned to Assumptions)

To quantify latent trajectories, we define three metrics directly reflecting our core assumptions:

Theoretical Mapping Visualization

Transformer Architecture and Dynamical System Mapping



Figure 3: Mapping DMET dynamics to Transformer layers. Detailed derivations and error bounds are provided in Appendix A.

1. State Continuity (smoothness):

$$C = \frac{1}{T} \sum_{t=1}^{T} \|\mathbf{h}_{t} - \mathbf{h}_{t-1}\|_{2}.$$
 250

255

257

259

260

261

264

266

267

269

270

271

272

273

274

275

277

278

279

282

2. Attractor Clustering Quality (structure):

$$Q = \frac{1}{N} \sum_{i=1}^{N} \frac{b(i) - a(i)}{\max\{a(i), b(i)\}},$$
 25

where a(i) and b(i) are intra- and nearestneighbor cluster distances.

3. Topological Persistence (global stability):

$$P = \sum_{\alpha} \left| d_{\alpha} - b_{\alpha} \right|, \qquad 262$$

summing the lifespans of topological features (birth b_{α} , death d_{α}).

Note: Implementation details—PCA for dimension reduction, k-means clustering for Q, and Vietoris–Rips filtration for P—are described in Appendix B.

3.5 Theory–Quality Correspondence

To bridge the gap between theory and practical application, we posit a central hypothesis: *the dynamic properties of latent trajectories directly determine the quality of generated text*, manifesting in three critical correspondences—**state continuity** leads to fluency, **clustering quality** underpins grammaticality, and **topological persistence** ensures semantic coherence. These associations are grounded not merely in conjecture, but in rigorous theoretical analysis and linguistic intuition. We elaborate on the mechanisms underlying each relationship below. We formalize three propositions linking our dynamic metrics to generation quality:

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

284

290 291

294 295

301

302 303

305

306

307

311 312

315

316

317

319

321

323

324

325

327



$$P = \sum_{\alpha} |d_{\alpha} - b_{\alpha}|$$

Proposition 1 (Continuity–Fluency). Higher

state continuity C implies smoother information

 $C = \frac{1}{T} \sum_{t} \|\mathbf{h}_t - \mathbf{h}_{t-1}\|_2 \quad \text{is large,}$

then the KL divergence between successive token

 $D_{KL}(p(w_t|w_{\leq t}) \parallel p(w_{t+1}|w_{\leq t+1})),$

The link between state continuity and textual

fluency can be understood through the lens of infor-

mation flow: smooth transitions between adjacent

latent states enable gradual information blending

rather than abrupt changes. For example, consider

the sentence "The clouds drift slowly across the

sky." In latent space, the transition from "sky" to

"clouds" to "drift" is realized as a smooth semantic

evolution, with natural progression between tokens.

In contrast, abrupt state transitions may yield inco-

herent outputs such as "The computers drift slowly

across the sky." Thus, state continuity fosters natu-

ral word choice, syntactic flow, and overall fluency.

Proposition 2 (Clustering–Grammaticality).

Stronger attractor separation Q supports robust

 $Q = \frac{1}{N} \sum_{i} \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$

is high, then the model remains within consistent

syntactic attractors, reducing grammatical errors.

terpreted as stable representations of grammat-

that such grammatical regimes are clearly sepa-

rated, allowing the model to reliably identify and

maintain correct syntax. For instance, syntactic

rules-such as subject-verb agreement or tense con-

sistency-manifest as stable attractors; initiating

a "If...then..." structure triggers a specific attrac-

tor, guiding the model to complete a well-formed

conditional sentence. Distinct attractor boundaries

prevent the model from abruptly switching gram-

matical frameworks (e.g., from declarative to inter-

Greater topological persistence P ensures stable

global structure and semantic coherence. When

rogative), or mixing tenses within a sentence.

3

Attractor structures in latent space can be in-

High clustering quality indicates

(Persistence-Coherence).

grammatical regimes. If

ical states.

Proposition

flow and lower perplexity. Formally, if

remains small, yielding more fluent text.

distributions.

is large, thematic loops and topic transitions remain well-connected, preventing abrupt topic shifts.

Topological persistence captures the global stability of manifold organization, particularly the prominence of loops, connecting paths, and semantic "regions." Consider an article on climate change, which might cover diverse subtopics such as scientific evidence, policy responses, and societal impacts. In latent space, these subtopics form distinct "semantic zones," and high topological persistence ensures that stable paths connect these zones-allowing smooth thematic transitions and preserving overall document coherence. In contrast, low persistence may yield abrupt topic shifts or logical discontinuities. Notably, the H_1 homology group (representing persistent cycles) is closely tied to argumentative closure and logical completeness in text. Persistent cycles facilitate the return of arguments to central themes, enabling essays to form closed, self-consistent reasoning structures rather than fragmenting into disconnected parts.

3.6 Summary

DMET provides a unified framework by casting 350 LLM generation as a controlled dynamical system 351 on a semantic manifold; it introduces quantitative 352 metrics-state continuity, attractor clustering, and 353 topological persistence-that offer concrete, measurable lenses on latent evolution; it maps these 355 dynamics to Transformer components, where resid-356 ual connections implement inertia, self-attention 357 provides contextual force, and feed-forward layers 358 approximate gradient flow; and it delivers prac-359 tical value by demonstrating that correlations be-360 tween latent dynamics and text quality can guide 361 decoding parameter tuning to achieve improved 362 fluency, grammaticality, and coherence. Our mani-363 fold and attractor assumptions may fail under high-364 dimensional noise or in very long sequences where 365 semantic drift accumulates. In such cases, trajec-366 tories can wander off $\mathcal M$ or traverse spurious at-367 tractors. Future work includes adapting DMET 368 to non-Transformer architectures (e.g., diffusion-369 based generators), and modeling advanced decod-370 ing strategies (e.g., beam search, mixture sam-371 pling) by incorporating multi-step lookahead forces 372 $g(\mathbf{h}, \{\mathbf{u}_{t+1}, \dots\}).$ 373

- 374
- 376
- 378

391

4 **Experiments and Analysis**

4.1 **Experimental Setup**

4.1.1 **Models and Data**

We use the DeepSeek-R1 Transformer as our base model. For each decoding configuration, we generate 10 continuations of 100 tokens each from the prompt:

"The future of AI is"

This yields a total of 400 samples across all settings. Hidden states are extracted from every layer of the model at each token step.

4.1.2 Decoding Parameter Grid

We sweep the temperature τ over 10 values from 0.1 to 2.0 and the nucleus (top-p) threshold over $\{0.3, 0.6, 0.8, 1.0\}$, resulting in 40 unique configurations.

4.1.3 Validation Pipeline

Algorithm 1 summarizes our pipeline for computing the four latent-dynamics metrics from each generated sequence.

Algorithm 1 Latent Dynamical System Validation **Require:** Transformer model M, input text x**Ensure:** Dynamics metrics $\{\delta, \mathcal{J}, s, \rho\}$ 1: $\mathbf{H} \leftarrow \text{GetHiddenStates}(M, x)$ 2: $\delta \leftarrow \text{ComputeDistances}(\mathbf{H})$ 3: $\mathcal{J} \leftarrow \text{DetectJumps}(\delta)$ 4: $\mathbf{V} \leftarrow \text{ReduceDim}(\mathbf{H})$ 5: $s \leftarrow \text{ClusterStates}(\mathbf{V})$ 6: $\rho \leftarrow \text{ComputePersistence}(\mathbf{V})$ 7: return $\{\delta, \mathcal{J}, s, \rho\}$

Evaluation Metrics 4.2

We evaluate both latent dynamics and textquality using a three-tiered framework. For dynamic metrics, we measure state continuity as $C = \frac{1}{T} \sum_{t=1}^{T} \|\mathbf{h}_t - \mathbf{h}_{t-1}\|_2$, attractor clustering via the silhouette-inspired score Q = $\frac{1}{N}\sum_{i=1}^{N}\frac{b(i)-a(i)}{\max\{a(i),b(i)\}}$, and topological persistence $P = \sum_{\alpha} |d_{\alpha} - b_{\alpha}|$. For text-quality metrics, we use both intrinsic measures-perplexity (computed with GPT-2-XL) and lexical diversity (log type-token ratio)-and extrinsic measures, including grammar accuracy and topical coherence.

Correlation Analysis: We fit mixed-effects regression models predicting each text-quality metric

from the three dynamic metrics, treating temperature and top-p as random effects to isolate their influence.

4.3 **Experimental Results**



Figure 4: PCA projection of 400 sample trajectories, showing two robust clusters (silhouette = 0.76)

Attractor Structure Analysis. Figure 4 visualizes the latent-space attractor structure via PCA and k-means clustering. We observe two prominent clusters-one large, one smaller-indicating that hidden states converge to distinct, stable regions rather than spreading randomly. This empirical finding aligns with our DMET prediction of semantic attractors.



Figure 5: Aggregated trajectories of all 400 samples as a 3D surface, colored from purple (start) to red (end).

Collective Dynamics. Figure 5 presents the aggregate dynamics of all 400 samples, visualized as a surface in 3D latent space with time progression color-coded from purple (start) to red (end). Several key features emerge from this visualization: most trajectories originate in the purple region on the left, indicating similar initial semantic states; as generation proceeds, the trajectories disperse in different directions, forming a fan-shaped pattern;

412

394

406

407

428

420

421

422

423

424

409 410

411

	Dependent Variables				
Predictors	Log-PPL	Spelling	Lexical Diversity	Grammaticality	Coherence
State Continuity	-0.031^{***}	-0.000	-0.003^{***}	-0.001	0.002**
Clustering Quality	0.044	-0.010	-0.074	0.081^{*}	0.047
Topological Persistence	-0.000	0.000	-0.003	0.002	0.009^{***}
Random Effects (Var.)	0.962***	0.048	0.265	0.058	0.115^{*}
Observations (N)	120	120	120	120	120

Table 1: Mixed-effects Model Results for Text Quality Predictors

Note: Coefficients shown with significance levels: *p < 0.05, **p < 0.01, ***p < 0.001. Log-PPL = Log-Perplexity

multiple local clusters appear in space, corresponding to distinct semantic attractors; and the overall structure exhibits a coherent manifold rather than a random point cloud. These collective observations strongly support our hypothesis that LLM generation follows constrained dynamical evolution paths—analogous to fluid flow in physics—rather than exhibiting a random walk in representation space.

Correlation Between Dynamics and Text Quality. Table 1 and fig 6 summarizes the mixedeffects regression findings. State continuity correlates negatively with log perplexity ($\beta = -0.031$, p < 0.001) and lexical diversity ($\beta = -0.003$, p < 0.001), but positively with coherence ($\beta = 0.002$, p < 0.01), indicating a trade-off between fluency and creativity. Clustering quality (silhouette) is positively associated with grammatical accuracy ($\beta = 0.081$, p < 0.05). Topological persistence is strongly correlated with coherence ($\beta = 0.009$, p < 0.001), empirically validating the theoretical prediction that robust manifold topology underpins logical, coherent text.

Effect of Decoding Parameters. We observe that *low temperature* ($\tau \le 0.5$) yields highly deterministic, smooth trajectories converging on major attractors. *Moderate temperature* ($0.6 \le \tau \le 1.2$) enables a balance of exploration and convergence, maximizing both continuity and topological persistence. *High temperature* ($\tau \ge 1.3$) leads to more stochastic, jumpy trajectories and weaker clustering. Lower top-p values (0.3) constrain exploration and boost continuity, while higher values (0.8–1.0) support diversity at the expense of global coherence. Notably, a combination of moderate temperature ($\tau \approx 0.7$) and top-p (0.6–0.8) achieves the optimal balance between creativity and coherence.

4.4 Experimental Conclusion

Our experiments robustly validate the Dynamic Manifold Evolution Theory (DMET) by demonstrating that latent dynamics critically influence text quality. Clustering of 400 generated samples uncovers clear attractor structures-confirmed by multidimensional scaling to align with distinct semantic frames-showing that representations collapse to coherent regions rather than disperse randomly. Mixed-effects regression reveals that smoother trajectories (state continuity C) reduce perplexity ($\beta = -0.031, p < 0.001$) and boost coherence ($\beta = 0.002, p < 0.01$), stronger attractor separation (clustering quality Q) predicts grammatical accuracy ($\beta = 0.081, p < 0.05$), and greater topological persistence (P) enhances semantic coherence ($\beta = 0.009, p < 0.001$). Parameter sweeps further show that low sampling temperature ($\tau \leq 0.5$) yields overly deterministic paths, high temperature ($\tau \ge 1.3$) produces erratic trajectories with weak attractors, and moderate settings $(0.7 \le \tau \le 1.0, \text{ top-}p \in [0.6, 0.8])$ strike the optimal balance of creativity and coherence. Finally, increasing sequence length leads to decreased continuity and increased topological complexity-explaining semantic drift in long-form generation-while some tasks preserve stable clustering. These findings confirm DMET's predictions and offer practical decoding guidelines: by tuning sampling parameters to shape latent dynamics, one can systematically improve fluency, grammaticality, and semantic coherence.

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

5 Summary

In this work, we introduced Dynamic Manifold499Evolution Theory (DMET), a unified mathematical framework that conceptualizes LLM genera-500

450

451

452

453

454

455

456

457

458

459

460 461

462

463

464

465

429

430

431

432



Figure 6: Quality Metrics across Different Parameter.

tion as a dynamical system evolving on a high-502 503 dimensional semantic manifold. Our main contributions are: (1) establishing a formal mapping between continuous-time dynamical systems and the discrete Transformer architecture; (2) deriving rep-506 resentation stability conditions via Lyapunov the-507 ory; (3) defining quantifiable dynamic metrics; (4) 508 empirically validating strong correlations between these metrics and text quality; and (5) propos-510 ing theory-driven decoding parameter optimiza-511 tion strategies. Our experiments robustly support DMET's central predictions: state continuity en-513 hances fluency, attractor clustering improves gram-514 matical accuracy, and topological persistence en-515 sures semantic coherence. In particular, we demon-516 strate that tuning temperature and top-p thresholds can effectively shape latent-trajectory dynamics, 518 enabling fine-grained control over generation out-519 comes. From a broader theoretical perspective, 520 DMET reveals that language generation is driven jointly by an internal energy function (linguistic 522 knowledge) and an external input function (con-523 text), offering a principled basis for both interpret-524 ing current models and designing next-generation 525 architectures with improved consistency, reduced hallucination, and enhanced coherence. 527

6 Limitations

Despite these encouraging results, our study has several limitations. Firsty, computational complexity of manifold and topological analyses remains high for very large models; more efficient algorithms are needed for real-time or large-scale deployment. Second, while we demonstrate strong correlations, causal relationships between latent dynamics and text quality remain to be established; developing interventions to directly manipulate latent trajectories will be crucial. Fourth, our framework rests on the *idealized manifold assumption*; real LLM representations may exhibit complex folds and self-intersections, posing challenges for accurate manifold estimation. Finally, although we propose theory-based tuning strategies, practical control mechanisms for manipulating latent dynamics (e.g., optimized regularization or decoding algorithms) are yet to be developed.

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

7 Acknowledgements

During the writing of this article, generative artificial intelligence tools were used to assist in language polishing and literature retrieval. The AI tool helped optimize the grammatical structure and expression fluency of limited paragraphs, and assisted

in screening research literature in related fields. All 553 AI-polished text content has been strictly reviewed by the author to ensure that it complies with aca-555 demic standards and is accompanied by accurate citations. The core research ideas, method design and conclusion derivation of this article were in-558 dependently completed by the author, and the AI 559 tool did not participate in the proposal of any innovative research ideas or the creation of substantive content. The author is fully responsible for the 562 academic rigor, data authenticity and citation in-563 tegrity of the full text, and hereby declares that the 564 generative AI tool is not a co-author of this study. 565

References

570

571

573

574

583

584

588

589

590

591

593

594

595

598

601

- Anthropic. 2022. Mechanistic interpretability in language models: A survey. *arXiv preprint arXiv:2211.00593*.
- Anthropic. 2024. Claude 3 technical report. *arXiv* preprint arXiv:2401.10409.
- Georgios Arvanitidis, Lars Kai Hansen, and Søren Hauberg. 2018. Latent space oddity: on the curvature of deep generative models. *International Conference on Learning Representations*.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Davis, Dora Demszky, and 95 others. 2022. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

- Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K. Duvenaud. 2018a. Neural ordinary differential equations. *Advances in Neural Information Processing Systems*.
- Yuhang Chen, D. Kalashnikov, Pietro Perona, and P. Welinder. 2018b. Bn-nas: Neural architecture search with batch normalization. *arXiv preprint arXiv:1812.03443*.
- Xiao Dai, Ziling Song, Mikita Balesni, and Dani Yogatama. 2023. Representation engineering in large language models. *arXiv preprint arXiv:2310.01405*.
- Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. 2019. Augmented neural odes. *Advances in Neural Information Processing Systems*.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Maxwell Nova, David Dalrymple, Jared Kaplan, Sam McCandlish, Dario Amodei, and Chris Olah. 2021.
 A mathematical framework for transformer circuits. *Anthropic Technical Report.*

Alex Hernandez-Garcia, Daniel Y. Fu, Quan Vuong, Kartik Sreenivasan, Patrick Blöbaum, Jake Tyo, Shibani Santurkar, Ishita Dasgupta, L. Schmidt, Samuel J. Gershman, Peter W. Battaglia, and Benjamin A. Spector. 2024. Thought manifolds: Understanding llm reasoning through activation space geometry. arXiv preprint arXiv:2404.09813.

606

607

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

- John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations. *NAACL*.
- Sewon Huang, Guangxuan Xiao, Weijia Shi, Wenhan Xiong, J. Weston, William Cohen, Luke Zettlemoyer, and Tao Yu. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. *arXiv preprint arXiv:2305.14251*.
- Jiwei Li, Xinlei Chen, Eduard H. Hovy, and Dan Jurafsky. 2016. Visualizing and understanding neural models in nlp. *NAACL*.
- Kaiming Li, Ling Yang, Hongxu Yin, S. Levine, and Rene Vidal. 2023. On the stability of transformerbased models. *arXiv preprint arXiv:2306.00148*.
- Nelson F. Liu, Ananya Kumar, Percy Liang, and Robin Jia. 2023. Probes as instruments for causal understanding of neural networks. *arXiv preprint arXiv:2303.04244*.
- Siyu Liu, Yue Fan, Shujian Zhang, J. Weston, and L. Dinh. 2024. Topological analysis of language model representations. *arXiv preprint arXiv:2402.07622.*
- Yuanzhi Lu, Sebastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Dynamical systems perspective on transformer training. *arXiv preprint arXiv:2303.17504*.
- Andreas Mardt, Luca Pasquali, Hao Wu, and Frank Noé. 2018. Vampnets: Deep learning of molecular kinetics. *Nature Communications*.
- R. Thomas McCoy, Harsh Trivedi, Richard Socher, Tal Linzen, Naomi Saphra, Blerta Ferko Serif, and Alexander Rush. 2022. Linguistic structure inherent in gpt embeddings. *arXiv preprint arXiv:2211.00603*.
- John Miller and Moritz Hardt. 2019. Stable recurrent models. *International Conference on Learning Representations*.
- OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Shishir G. Patil, Tianjun Zhang, Tatsu Hashimoto, Tommi S. Jaakkola, Michael R. DeWeese, and Joseph Gonzalez. 2024. Time-evolution of transformer reasoning through dynamical systems. *arXiv preprint arXiv:2402.05633*.

Ashwin Rajamohan, Nate Gruver, Hattie Zhou,
Sorelle A. Friedler, Samy Bengio, and Been Kim.
2023. A focus in your hidden states: Evidence
for bifurcation in llm reasoning. *arXiv preprint*

661

667

670

671

672

673

675

676

677

678

679

687

691

697

698

702

707

710

- for bifurcation in llm reasoning. arXiv preprint arXiv:2402.03189.
 Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B. Viégas, Andy Coenen, Adam Pearce, and Been Kim.
- 2019. Visualizing and measuring the geometry of bert. Advances in Neural Information Processing Systems.Sam T. Roweis and Lawrence K. Saul. 2000. Nonlinear
- Sam 1. Rowers and Lawrence K. Saul. 2000. Nonlinear dimensionality reduction by locally linear embedding. *Science*.
 - Matheus Viana Santos, Yuan Gao, Aditi Krishnapriyan, and Michael W. Mahoney. 2023. A theory of learning dynamics in transformers with application to optimization. *arXiv preprint arXiv:2306.01129*.
 - Junjie Shi, Suhang Wang, and Dongwon Lee. 2024. Detailed evaluation of output stability in large language models. *arXiv preprint arXiv:2401.06706*.
 - Mukund Sundararajan, Xiaoyin Xie, Matej Zečević, and Matej Zemljic. 2022. Attribution methods in natural language processing: A survey. *arXiv preprint arXiv:2207.05585*.
 - J. B. Tenenbaum, V. de Silva, and J. C. Langford. 2000. A global geometric framework for nonlinear dimensionality reduction. *Science*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 50 others. 2024. Llama 3: Third-generation open foundation language models. *arXiv preprint arXiv:2404.14819*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. Advances in Neural Information Processing Systems.
- Xiang Zhang and Shunyu Xiao. 2024. Generative language modeling as feedforward markov decision process. *arXiv preprint arXiv:2402.17152*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

A Appendix Overview

This appendix provides all detailed derivations, implementation details, and additional results that support the main text.

- Appendix B: Mapping Transformer to Continuous Dynamics (corresponds to Sec. 3.3)
 Complete proofs and error-bound derivations for the correspondence between Transformer components (residual, MHSA, FFN, Layer-Norm) and the continuous-time dynamical system model.
- Appendix C: Experimental Method Details (corresponds to Sec. 4.1) Implementation specifics for computing dynamic metrics (state continuity, clustering, persistence) and text-quality metrics (perplexity, lexical diversity, grammar, coherence).

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

- Appendix D: Supplementary Results (corresponds to Sec. 4.2) Additional visualizations, including single-sequence trajectory phases and temperature ablation curves, that illustrate latent evolution and the "golden zone" for decoding parameters.
- Appendix E: Mathematical Proofs (corresponds to Sec. 3.5 and Sec. 3.6) Full statements and proofs of Lyapunov stability (Theorems 1–2), continuity–fluency, clustering–grammar, persistence–coherence (Theorems 3–5), and temperature effects (Propositions 3–4).

B Mapping between Transformer Architecture and Dynamical Systems

In the previous section we established the basic framework of Dynamic Manifold Evolution Theory. In this section, we delve into how this theory can be precisely mapped onto the concrete implementation of the Transformer architecture. We first provide rigorous mathematical proofs of the dynamical-system interpretation for each architectural component, then analyze the limitations and approximation errors of the mapping, and finally discuss the crucial modulatory role of Layer Normalization in this framework.

Table 2 summarizes the correspondence between DMET theoretical concepts and Transformer components. This table is intended to fit within a single column of a two-column layout.

Theory Concept Transformer Component Functional Role Latent state h(t)Hidden state Encodes current semantics Evolution time tLayer index l Discrete update step Energy function $V(\mathbf{h})$ Feed-forward network (FFN) Semantic optimization External function $g(\mathbf{h}, \mathbf{u})$ Multi-head self-attention (MHSA) Contextual integration Time step Δt Layer normalization + scaling Controls update magnitude Restricts representation to valid manifold Manifold constraint $c(\mathbf{h})$ Activation + Residual

Table 2: Mapping between DMET theoretical concepts and Transformer components.

B.1 Dynamical-System Interpretation of **Transformer Components**

754

755

756

758

759

761

765

767

770

772

773

775

776

777

779

B.1.1 Residual Connection as Inertia: Rigorous Proof

Residual connections are a key innovation in Transformers, allowing direct passage of the previous layer's output so that the network learns residual mappings. From a dynamical-system perspective, a residual connection implements "inertia," keeping the representation evolution continuous.

Theorem 3 (Equivalence of Residual Inertia). The Transformer residual update

$$\mathbf{h}_{t+1} = \mathbf{h}_t + F(\mathbf{h}_t)$$

is formally equivalent to the inertia term in the discrete Lagrangian system, where F denotes a nonlinear transformation.

Proof. Start from the continuous Lagrangian system:

$$\frac{d^2 \mathbf{h}}{dt^2} = \mathbf{f} \left(\mathbf{h}, \, \frac{d \mathbf{h}}{dt} \right). \tag{1}$$

Let $\mathbf{v} = \frac{d\mathbf{h}}{dt}$, yielding the first-order form:

$$\frac{d\mathbf{h}}{dt} = \mathbf{v},\tag{2}$$

$$\frac{d\mathbf{v}}{dt} = \mathbf{f}(\mathbf{h}, \mathbf{v}). \tag{3}$$

Applying the explicit Euler discretization with step Δt :

$$\mathbf{h}_{t+1} = \mathbf{h}_t + \Delta t \, \mathbf{v}_t,\tag{4}$$

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \Delta t \, \mathbf{f}(\mathbf{h}_t, \mathbf{v}_t). \tag{5}$$

In a highly damped regime, $\mathbf{v}_t \approx \Delta t \cdot \mathbf{f}(\mathbf{h}_t, \mathbf{0})$. Substituting into the update for h gives 781

$$\mathbf{h}_{t+1} = \mathbf{h}_t + (\Delta t)^2 \, \mathbf{f}(\mathbf{h}_t, \mathbf{0}). \tag{6}$$

Defining $F(\mathbf{h}_t) = (\Delta t)^2 \mathbf{f}(\mathbf{h}_t, \mathbf{0})$ yields 783

$$\mathbf{h}_{t+1} = \mathbf{h}_t + F(\mathbf{h}_t),$$

which matches the Transformer residual update.

This proof shows that residual connections discretely simulate physical inertia. The "strength" of inertia is controlled by Δt and affects trajectory smoothness.

Lemma 1 (Residual Strength-Continuity Relationship). Let α be a residual-strength parameter, with update

$$\mathbf{h}_{t+1} = \alpha \, \mathbf{h}_t + F(\mathbf{h}_t). \tag{79}$$

Then the continuity metric

$$C = \frac{1}{T} \sum_{t=1}^{T} \left\| \mathbf{h}_t - \mathbf{h}_{t-1} \right\|_2$$
 79

is monotonically decreasing in α .

B.1.2 Self-Attention as Contextual Force: **Functional Analysis**

Self-attention implements context-dependent dynamics in text generation, allowing representation evolution to respond to global information. We now prove how multi-head self-attention approximates an arbitrary Lipschitz-continuous "contextual force" function $g(\mathbf{h}, \mathbf{u})$.

Theorem 4 (Equivalence of Self-Attention and Contextual Force). Under suitable parameterization, multi-head self-attention (MHSA) can approximate any Lipschitz-continuous function

$$g: \mathbb{R}^d \times \mathbb{R}^{n \times d} \to \mathbb{R}^d$$
 809

arbitrarily well.

Proof Sketch. Standard MHSA computes

$$MHSA(\mathbf{h}_t, \mathbf{X}) = \sum_{i=1}^{h} W_i^O\left(\sum_{j=1}^{n} \alpha_{ij}(\mathbf{h}_t, \mathbf{X}) \cdot W_i^V \mathbf{x}_j\right), \qquad (7)$$

where the attention weights are defined as

$$\alpha_{ij}(\mathbf{h}_t, \mathbf{X}) = \frac{\exp\left((W_i^Q \mathbf{h}_t)^\top (W_i^K \mathbf{x}_j) / \sqrt{d_k}\right)}{\sum_{j'=1}^n \exp\left((W_i^Q \mathbf{h}_t)^\top (W_i^K \mathbf{x}_{j'}) / \sqrt{d_k}\right)}.$$
(8)

785

786

787

790

791

792

794

796

797

798

799

800

801

802

803

805

806

807

808

810

811

815 816

- 823
- 824 825
- 826

830 831

832

- 834 835

- 838
- 841
- 842

843

845

853

855

 $\mathbf{h}_{t+1} = \mathbf{h}_t + \text{FFN}(\text{LN}(\mathbf{h}_t)) + \text{MHSA}(\text{LN}(\mathbf{h}_t)),$ (10)

B.3 Layer Normalization as Time-Step

with

By Stone–Weierstrass, any continuous function

 $g(\mathbf{h}, \mathbf{u})$ can be approximated by finite sums of basis functions in h with context-dependent coefficients.

As the number of heads h grows, MHSA's paramet-

ric form can realize these sums to arbitrary preci-

sion. A detailed epsilon-net construction shows the

The position-wise feed-forward network (FFN) cor-

responds to a negative gradient field $-\nabla V(\mathbf{h})$,

driving the system toward local minima of a se-

Theorem 5 (Equivalence of FFN and Gradi-

ent Field). A two-layer ReLU FFN of sufficient

width can approximate any smooth energy gradient

field $-\nabla V(\mathbf{h})$ on compact sets to arbitrary accu-

Proof Sketch. By the universal approximation

theorem, the FFN can approximate any continuous

 $\mathcal{L}_{\text{curl}} = \left\| \nabla \times \text{FFN}(\mathbf{h}) \right\|_{F}^{2}$

drives the learned field to be (approximately) irrotational, hence a gradient of some scalar potential.

Although Transformers can approximate continu-

ous dynamical systems, the mapping incurs inher-

ent errors. We bound the global error between the

continuous trajectory $\mathbf{h}(t)$ and its discrete counter-

where ϵ_{FFN} and ϵ_{MHSA} are the FFN and MHSA

approximation errors respectively. Standard error-

LayerNorm not only stabilizes training but also

adapts the effective time-step. Consider the update

analysis yields the stated bound.

Modulator

Mapping Limitations and Approximation

vector field. Imposing a curl-penalty

Error Analysis

B.1.3 Feed-Forward Network as Gradient

Field: Expressivity Proof

sup-norm error decays as $O(1/\sqrt{h})$.

mantic energy function V.

racy.

B.2

part $\mathbf{h}_{|t/\Delta t|}$:

$$LN(\mathbf{h}_t) = \gamma \, \frac{\mathbf{h}_t - \mu_t}{\sigma_t} + \beta.$$

If FFN and MHSA are approximately homogeneous of degree one, then

$$\Delta t_{\rm eff} = \frac{\gamma}{\sigma_t} \tag{860}$$

858

859

862

864

865

868

869

871

872

874

875

876

877

878

879

880

882

883

884

885

886

889

890

891

acts as an adaptive step-size. A stability condition 861 follows:

$$\gamma < \sigma_{\min} \frac{2}{\lambda_{\max}},$$
 863

where λ_{\max} is the largest eigenvalue of the Jacobian.

This completes the appendix material for "Sec-866 tion 4 — Continuous-to-Discrete Mapping." 867

Experimental Details С

C.1 Latent Dynamics Computation

We compute the three dynamic metrics as follows: 870

• State Continuity: For each hidden-state sequence $\{\mathbf{h}_t\}_{t=0}^T$, we calculate

$$C = \frac{1}{T} \sum_{t=1}^{T} \|\mathbf{h}_t - \mathbf{h}_{t-1}\|_2.$$
 873

Hidden states are L2-normalized before differencing.

• Attractor Clustering: We apply PCA to reduce $\{\mathbf{h}_t\}$ to 2–3 dimensions (retaining > 85% variance), then run k-means++ with k chosen by silhouette analysis. We report the average silhouette score

$$Q = \frac{1}{N} \sum_{i=1}^{N} \frac{b(i) - a(i)}{\max\{a(i), b(i)\}},$$
881

 $\sup_{t \in [0,T]} \left\| \mathbf{h}(t) - \mathbf{h}_{\lfloor t/\Delta t \rfloor} \right\|_{2} \leq C_{1} \Delta t + C_{2} \epsilon_{\text{FFN}} + C_{3} \epsilon_{\text{MHSA}},$ where a(i) and b(i) are intra- and nearestcluster distances.

> • Topological Persistence: Using a Vietoris-Rips filtration over a range of radii, we compute H_1 barcodes via Ripser. Persistence is measured as

$$P = \sum_{\alpha} |d_{\alpha} - b_{\alpha}|, \qquad 888$$

summing lifespans of 1D cycles. We apply a significance threshold $\rho > 0.2$ determined by permutation testing.

(9)

- 897

- 900
- 901
- 902
- 903

- 904 905
- 906
- 907
- 908
- 909

910 911

C.2 Text Quality Evaluation

• Lexical Diversity:

type-token ratio:

and learned metrics:

model:

We evaluate generated text using both automated

• Perplexity: Computed with a GPT-2-XL

 $PPL = \exp\left(-\frac{1}{N}\sum_{i=1}^{N}\log p(w_i|w_{< i})\right).$

 $LTTR = \frac{\log(\#types)}{\log(\#tokens)}.$

• Grammar Accuracy: Detected via a fine-

 $1 - \frac{\# \text{errors}}{\# \text{tokens}}.$

• Coherence: Entity-grid model combining lo-

Coherence = $0.7 \times \text{Local} + 0.3 \times \text{Global}$.

tuned BERT classifier; score is

cal and global coherence:

Measured as log

• Factuality: Checked against a knowledge graph; factuality score is

 $\frac{\#\text{correct facts}}{\#\text{verifiable facts}}.$

C.3 Supplementary Results

Trajectory Evolution Analysis. Figure 7-10 depict the evolution of latent representations along 912 the generation trajectory of a single sequence. By 913 tracking the representations across time steps (to-914 kens), we identify a characteristic three-phase pat-915 tern: during the Initial Phase (first 10 tokens), the 916 trajectory explores a local neighborhood, reflecting 917 a search for initial semantic direction; in the Ex-918 pansion Phase (approximately tokens 30-60), the 919 trajectory expands into new regions, corresponding to topic development and elaboration; finally, in 921 the Convergence Phase (around tokens 70–100), 922 the trajectory moves toward a specific region, in-924 dicating the natural closure of content. This dynamic progression aligns closely with our theoreti-925 cal framework: a well-formed generation process exhibits a structured transition from exploration to convergence in latent space. 928

D Mathematical Proofs

This appendix provides detailed proofs for the the-930 orems and propositions referenced in the main text. 931 Throughout, let h denote the latent representation, 932 $V(\mathbf{h})$ the energy function, and $q(\mathbf{h}, \mathbf{u})$ the external 933 input function.

D.1 Lyapunov Stability Proof

Theorem 1 (Stability Condition). *If for all* $h \neq d$ h*,

$$(\mathbf{h}-\mathbf{h}^*)^T Q g(\mathbf{h},\mathbf{u}) \le (\mathbf{h}-\mathbf{h}^*)^T Q \nabla V(\mathbf{h}),$$
(11)

then the equilibrium \mathbf{h}^* is asymptotically stable.

Proof. Consider the Lyapunov function $L(\mathbf{h}) =$ $(\mathbf{h} - \mathbf{h}^*)^T Q(\mathbf{h} - \mathbf{h}^*)$, where Q is symmetric positive definite. $L(\mathbf{h}) > 0$ for $\mathbf{h} \neq \mathbf{h}^*$ and $L(\mathbf{h}^*) = 0$.

The time derivative is

$$\frac{dL}{dt} = 2(\mathbf{h} - \mathbf{h}^*)^T Q \frac{d\mathbf{h}}{dt}$$
(12) 944

$$= 2(\mathbf{h} - \mathbf{h}^*)^T Q \left[-\nabla V(\mathbf{h}) + g(\mathbf{h}, \mathbf{u}) \right]$$
(13)

$$= -2(\mathbf{h} - \mathbf{h}^*)^T Q \nabla V(\mathbf{h}) + 2(\mathbf{h} - \mathbf{h}^*)^T Q g(\mathbf{h}, \mathbf{u})$$
(14)

By the assumption, the second term is no greater 947 than the first, so $\frac{dL}{dt} \leq 0$. If the inequality is strict, 948 then $\frac{dL}{dt} < 0$. By Lyapunov's direct method and 949 LaSalle's invariance principle, all trajectories con-950 verge to h^* . 951

D.2 Discrete System Stability

Theorem 2 (Discrete Stability). For the discrete system $\mathbf{h}_{t+1} = \mathbf{h}_t + \Delta t [-\nabla V(\mathbf{h}_t) + g(\mathbf{h}_t, \mathbf{u}_t)],$ suppose:

1.
$$\nabla V(\mathbf{h}_t)^T g(\mathbf{h}_t, \mathbf{u}_t) \le \|\nabla V(\mathbf{h}_t)\|^2$$
 95

2. $\Delta t < 2/\lambda_{\max}(H_V)$, where H_V is the Hessian of V

Then the discrete system is stable, and $V(\mathbf{h})$ decreases approximately monotonically.

Proof. The change in V per update is:

$$\Delta V_t = V(\mathbf{h}_{t+1}) - V(\mathbf{h}_t) \tag{15}$$

$$\approx \nabla V(\mathbf{h}_t)^T(\mathbf{h}_{t+1} - \mathbf{h}_t)$$
 (16) 96

$$+\frac{1}{2}(\mathbf{h}_{t+1}-\mathbf{h}_t)^T H_V(\mathbf{h}_{t+1}-\mathbf{h}_t)$$
 (17) 964

$$= \Delta t [-\|\nabla V(\mathbf{h}_t)\|^2 + \nabla V(\mathbf{h}_t)^T g(\mathbf{h}_t, \mathbf{u}_t)]$$
(18)

$$+\frac{1}{2}\Delta t^2 \lambda_{\max}(H_V) \| -\nabla V(\mathbf{h}_t) + g(\mathbf{h}_t, \mathbf{u}_t) \|^2$$
(19)

934

929

935 936

937

938 939

940

941

942

943

952

953

954

955

957

958

959

960

Trajectory Evolution (temp=0.1, top_p=0.3)



Figure 7: Dynamic evolution along a generation trajectory in 2D latent space (temperature=0.1 and top_K=0.3).

Trajectory Evolution (temp=0.1, top_p=0.6)



Figure 8: Dynamic evolution along a generation trajectory in 2D latent space (temperature=0.1 and top_K=0.6).

Trajectory Evolution (temp=0.1, top_p=1.0)



Figure 9: Dynamic evolution along a generation trajectory in 2D latent space (temperature=0.1 and top_K=1.0).

Trajectory Evolution (temp=2.0, top_p=0.6)



Figure 10: Dynamic evolution along a generation trajectory in 2D latent space (temperature=2.0 and top_K=0.6).

Condition 1 ensures the first term is nonpositive; condition 2 ensures the second (quadratic) term 968 does not dominate. Thus, $\Delta V_t \leq 0$ if both conditions hold.

967

969

970

973

974

975

976

977

978

981

985

990

991

993

997

998

999

1001

1002

1003

1004

1005

1006

1007

1009

D.3 Continuity–Fluency, Clustering–Grammar, Persistence–Coherence (Theorems 3–5)

Theorem 3 (Continuity–Fluency). Higher state continuity

$$C = \frac{1}{T} \sum_{t} \|\mathbf{h}_t - \mathbf{h}_{t-1}\|_2$$

implies smaller KL divergences between successive token distributions,

$$D_{KL}(p_t \| p_{t+1}),$$

and thus smoother, more fluent text.

Proof Sketch. Since \mathbf{h}_t parametrizes $p(w_t | \mathbf{h}_t)$ smoothly, small latent steps yield small changes in logits and hence low D_{KL} . Empirically this correlates with lower perplexity.

4 Theorem (Clustering–Grammaticality). Higher silhouette score Q indicates well-separated latent attractors corresponding to distinct syntactic regimes, thereby reducing grammatical errors.

Proof Sketch. Well-defined clusters imply distinct syntactic states; transitions remain within a single cluster, preventing abrupt grammatical shifts.

Theorem 5 (Persistence–Coherence). Greater topological persistence

$$P = \sum_{\alpha} |d_{\alpha} - b_{\alpha}|$$

reflects robust global manifold structure, supporting coherent theme transitions.

Proof Sketch. Persistent homology captures longlived loops correlating with thematic cycles; higher P ensures stable topic connectivity and logical flow.

D.4 Parameter Sensitivity Analysis and **Optimization Strategies**

A central practical question is how generation parameters shape the dynamic evolution of latent trajectories. Our theoretical framework provides fresh insight into the influence of temperature (τ) and sampling threshold (top-p), and yields actionable strategies for controllable, high-quality generation. **Temperature Effects: Dynamics of Randomness.** 1010 Temperature (τ) serves as a primary dial for con-1011 trolling stochasticity during generation. Our theory 1012 predicts that temperature fundamentally reshapes 1013 latent dynamics: low temperatures ($\tau \rightarrow 0$) en-1014 hance state continuity, reduce topological complex-1015 ity, and reinforce dominant attractors, while high 1016 temperatures ($\tau \rightarrow \infty$) decrease continuity, am-1017 plify topological diversity, and weaken attractor 1018 structure. This is closely analogous to physical sys-1019 tems, where low temperatures yield ordered states 1020 (e.g., crystalline solids) and high temperatures in-1021 duce disorder (e.g., gases). In language modeling, 1022 low temperatures produce highly deterministic, po-1023 tentially rigid text that closely follows the "energy-1024 minimizing" path, whereas higher temperatures in-1025 troduce more exploratory, creative, but potentially 1026 less coherent content. Mathematically, tempera-1027 ture scales the logits in the sampling distribution, 1028 $p(w|\mathbf{h}) \propto \exp(z_w/\tau)$; as $\tau \to 0$, the distribution 1029 collapses to the argmax, while large τ makes the 1030 distribution uniform, directly modulating trajectory 1031 determinism and diversity. 1032

Sampling Threshold Effects: Dynamic Constraints on Feasible Space. Beyond temperature, top-p (nucleus) sampling imposes a dynamic constraint on the allowable state space. Our framework predicts: lower top-p restricts trajectories to a narrow feasible region, increasing continuity but potentially limiting topological complexity; higher top-pexpands the search space, potentially reducing continuity but enriching manifold structure and output diversity. This can be likened to path planning in traffic systems: low top-p is akin to only permitting travel on main highways, ensuring smooth but constrained trajectories, while high top-p opens up all roads, allowing for more exploration-at the cost of potential complexity or detours. By limiting the set of next-token candidates, low top-p "smooths out" suboptimal paths, whereas high top-p retains more manifold detail, shifting the balance between exploration and exploitation.

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1049

1050

1051

1052

1053

1054

1055

Optimization Strategies: Theory-Grounded Practical Guidance. Leveraging these insights, we propose three core optimization strategies for practical generation control:

1. Balanced Parameters: To trade off coher-1056 ence and creativity, set moderate tempera-1057 ture ($\tau \approx 0.7-0.8$) and top-*p* ($p \approx 0.7-0.9$), yielding text that is both inventive and well-1059



Figure 11: Effects of decoding parameters on dynamic metrics and recommended optimization strategy.

structured. For tasks like story or essay writing, this balance prevents the system from being overly deterministic while maintaining sufficient continuity to avoid abrupt logical jumps.

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1074 1075

- 2. Task-Adaptive Tuning: Adjust parameters based on task requirements—use lower temperature ($\tau \approx 0.3-0.5$) for highly coherent, technical documents (e.g., API documentation, legal texts), and higher temperature ($\tau \approx 0.9-1.2$) for creative or poetic tasks, where exploration and novelty are valued. In the former, strict adherence to grammatical attractors is vital; in the latter, higher temperature encourages novel expressions, while persistent topology ensures overall theme cohesion.
- 3. Dynamic Adjustment: Modulate parameters 1077 during generation-begin with higher temper-1078 ature ($\tau \approx 0.8$ –1.0) to encourage exploration 1079 ("brainstorming" phase), then lower $\tau \approx 0.5$ – 1080 0.7) for convergence and refinement ("polish-1081 ing" phase). For example, in academic writ-1082 ing, initial high temperature facilitates idea 1083 diversity; subsequently decreasing τ enables 1084 the system to consolidate around optimal ar-1085 gumentative structure. 1086