Multi-Scale Manifold Alignment: A Unified Framework for Enhanced Explainability of Large Language Models

Anonymous ACL submission

Abstract

Recent advances in Large Language Models (LLMs) have achieved strong performance, yet 004 their internal reasoning remains opaque, limiting interpretability and trust in critical applications. We propose a novel Multi-Scale Manifold Alignment framework that decomposes 007 800 the latent space into global, intermediate, and local semantic manifolds-capturing themes, context, and word-level details. Our method 011 introduces cross-scale mapping functions that jointly enforce geometric alignment (e.g., Pro-012 crustes analysis) and information preservation (via mutual information constraints like MINE or VIB). We further incorporate curvature regularization and hyperparameter tuning for stable optimization. Theoretical analysis shows that 017 018 alignment error, measured by KL divergence, can be bounded under mild assumptions. This 019 framework offers a unified explanation of how LLMs structure multi-scale semantics, advancing interpretability and enabling applications such as bias detection and robustness enhancement

1 Introduction

037

041

1.1 Background and Motivation

Large language models (LLMs) such as GPT-4 (OpenAI, 2023), LLaMA (Touvron et al., 2023) and PaLM (Chowdhery et al., 2022) have achieved human-level performance across a range of NLP tasks (Brown et al., 2020), yet their growing complexity renders them opaque and limits trust in high-stakes applications like healthcare and finance (Bommasani et al., 2021). Prior interpretability efforts—attention visualization (Vig, 2019), neuron activation analysis (Gurnee et al., 2023) and probing tasks (Tenney et al., 2019)—provide layer-specific insights but fail to capture how semantic information is transmitted and integrated across layers (Geva et al., 2023). While studies have traced inter-layer information flow (Hahn and Jurafsky, 2023; Belinkov and Riedl, 2022), a unifying theoretical framework remains absent. 042

043

044

047

048

053

054

056

059

060

061

062

063

064

065

066

067

069

070

071

072

073

074

075

077

078

Empirical and theoretical work shows that Transformer representations organize hierarchically: lower layers encode lexical and syntactic details, intermediate layers capture local semantic relations, and higher layers model global discourse (Zhang et al., 2023; Seo et al., 2023), mirroring stages in human language processing (Friederici, 2011). Manifold-based analyses (Daxberger et al., 2023) and alignment techniques (Wang et al., 2023; Li et al., 2023), grounded in information geometry (ichi Amari and Nagaoka, 2007) and representation learning principles (Bengio et al., 2013), suggest modeling these semantic strata as nested manifolds. Building on these insights, we propose Multi-Scale Manifold Alignment, a unified framework that learns cross-scale geometric and informationtheoretic mappings to analyze and control LLM internals.

1.2 Contributions

We propose a multi-scale manifold alignment theory that provides a unified framework for analyzing LLM information processing across semantic scales. Our key contributions include:

- Hierarchical Decomposition: We decompose LLM hidden spaces into three semantic manifolds—global (document-level), intermediate (sentence-level), and local (word-level)—demonstrating this structure emerges naturally across architectures.
- Cross-Scale Alignment: We develop novel mapping functions combining geometric alignment (Procrustes analysis) with information-theoretic constraints (mutual information), enabling precise tracking of information flow.
- Theoretical Guarantees: We establish error 079

bounds for alignment quality using KL divergence, grounded in information geometry principles.

> • Practical Framework: We provide complete implementation including layer identification, cross-model adaptation, and optimization strategies, with applications in bias detection and robustness enhancement.

Compared to single-scale approaches, our framework offers a comprehensive view of LLM information organization, advancing interpretability and control.

The remainder of the paper is organized as follows. Section 2 reviews related work. Section 3 presents the proposed multi-scale manifold alignment framework. Section 4 reports experimental results. Section 5 concludes the paper and outlines future directions.

2 Related Work

096

100

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

2.1 Explainability of Large Language Models

Understanding the internal mechanisms of large language models (LLMs) has become a major research topic in natural language processing. As model scale and complexity grow, developing effective interpretability methods is increasingly crucial.

2.1.1 Attention Mechanism Analysis

The attention mechanism of Transformer architectures offers an intuitive lens on model internals. Vig (2019) pioneered visualising attention as an explanatory tool, spurring this area of research. Recently, Sarti et al. (2023) analysed whether attention in large language models reflects linguistic structure, finding clear correspondence to syntactic dependency relations. Focusing on the reliability of attention weights, the Attention Attribution method proposed by Chefer et al. (2021) combines gradients and attention to more accurately identify key input components driving model decisions. The latest work also explores functional differences among multi-head attentions: Xie et al. (2023) show that distinct attention heads in GPT-3 and PaLM form specialised clusters responsible for particular linguistic tasks.

2.1.2 Representation Analysis and Probing

Neural probes allow researchers to detect what information is encoded in model representations. Recently, Meng et al. (2022) introduced a locatingand-extracting method that pinpoints specific neurons and attention heads storing knowledge in large language models, markedly improving explainability. In representation-space analysis, Li et al. (2023) explored how fine-tuning and prompt learning reshape model topology, finding that even minimal fine-tuning can significantly reconfigure representational geometry. Liu et al. (2023) analysed the "thought process" of LLMs, proposing hidden-state tracing to follow reasoning paths. Notably, Gurnee et al. (2023) recently introduced a representationspace decomposition method that projects hidden states onto specific semantic features, offering a new tool for understanding representation organisation.

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

153

154

155

156

157

158

159

160

161

162

163

164

165

166

168

169

170

171

172

173

174

175

176

2.1.3 Information-theoretic and Causal Analyses

An information-theoretic perspective provides a principled framework for understanding information flow in language models. Ghandeharioun et al. (2023) proposed *PatchScores*, which quantifies mechanism importance via causal interventions, revealing how different components contribute to performance. For mutual-information analysis, Xu et al. (2023) introduced a new method that explains compositionality by measuring information transfer among components. Hahn and Jurafsky (2023) recently proposed *information-trajectory analysis* to track how specific information fragments propagate through model hierarchies, offering fresh insight into inter-layer information exchange.

2.2 Multi-scale Representations and Hierarchical Analysis

2.2.1 Hierarchical Representation Learning

Language is inherently hierarchical—from characters to words, syntax, semantics, and discourse. Recent studies show that LLMs form similar hierarchies internally. Zhang et al. (2023) revealed layered representation patterns in Transformers: shallow layers process word-level features, mid-layers capture syntax, and deep layers integrate semantics. In this area, Seo et al. (2023) demonstrated through large-scale experiments that features extracted at different layers align closely with stages in the traditional NLP pipeline, reflecting a shallow-to-deep processing logic. Singh and III (2023) systematically evaluated hierarchical capability across model sizes, finding that larger models reinforce this hierarchy.

2.2.2 Cross-scale Information Transfer

177

178

179

181

185

186

189

190

191

193

194

195

198

199

203

204

208

210

211

213

214

215

216

217

218

219

220

221

223

224

Understanding information transfer across representation levels is vital for explaining model reasoning. Hernandez et al. (2023) used careful experimental design to show how information flows from shallow to deep layers—from local to global—forming a coherent information network. Recently, manifold alignment has made notable progress in cross-modal learning. Wang et al. (2023) proposed an alignment framework that connects semantic structures across modalities or representation spaces, offering valuable ideas for aligning different semantic scales.

2.3 Information Geometry and Manifold Theory

2.3.1 Applications of Information Geometry to Neural Networks

Information geometry offers a rigorous mathematical framework for analysing neural networks.
Raghu et al. (2022) pioneered using the Fisher information matrix to analyse LLMs, demonstrating how pre-training shapes representational geometry. In theoretical advances, Gromov–Monge optimal transport was applied to representation-space analysis by Chuang et al. (2023), providing new metrics for geometric similarity between representations—laying a solid foundation for studying mappings among manifolds.

2.3.2 A Manifold View of Language-model Representations

Viewing language-model representations as manifolds has yielded several breakthroughs. Recently, Daxberger et al. (2023) used manifold analysis to reveal semantic organisation in representation space, finding separable sub-manifolds for different concepts. Of particular note, the seminal work of Elhage et al. (2022) proposed treating internal representations as nested manifolds, offering theoretical tools for feature decomposition and information transfer. Grover et al. (2023) further combined the manifold view with geometric deep learning, introducing methods to analyse curvature and topology of embedding spaces.

2.4 Research Gaps

Despite rich research on LLM explainability, critical gaps remain. First, there is a lack of a unified multi-scale theoretical framework. Existing studies often focus on a single semantic level or merely analyse inter-layer differences; a formal framework describing and analysing crossscale semantic organisation and transfer is missing. Second, there is a separation between geometric and information-theoretic methods. Although geometric (distance/structure-based) and informationtheoretic (mutual-information/entropy-based) approaches have advanced separately, no framework unifies them. Third, cross-scale mapping mechanisms remain under-explored. Prior work has not systematically investigated mapping functions across semantic scales—especially how to preserve geometric structure and semantic information simultaneously.

226

227

228

229

230

231

232

233

234

235

236

237

239

240

241

243

244

245

246

247

248

249

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

269

270

3 Theory and Framework

This section develops the theoretical foundation for multi-scale manifold alignment in large language models (LLMs), systematically unifying geometric, information-theoretic, and practical aspects. We first motivate our approach with information geometry, then formalize multi-scale semantic decomposition, introduce cross-scale mappings, and present the overall optimization framework.

3.1 Theoretical Foundations: Multi-Layered Nature of LLM Representations

Large language models naturally develop hierarchical internal representations as a result of both model architecture and the intrinsic layered structure of human language. Attention patterns and hidden feature distributions reveal that each model layer progressively aggregates and abstracts semantic information, giving rise to distinct strata in representation space.

In the lens of **information geometry**, each hidden state can be viewed as a point in a highdimensional space, collectively forming a *statistical manifold*—the space of parameterized probability distributions associated with each state.

Definition 1 (Statistical Manifold). Given a family of distributions $\{p(x \mid \theta)\}$ parameterized by $\theta \in \Theta$ with observed data $x \in \mathcal{X}$, the statistical manifold \mathcal{M} is the set of all such distributions, each corresponding to a unique representation.

The **Fisher information matrix** gives a Riemannian metric for \mathcal{M} , quantifying how infinitesimal changes in parameters affect the probability distri-

272

277

278

279

281

284

287

290

291

293

296

297

299

301

305

309

312

bution:

$$F_{ij}(\theta) = \mathbb{E}_{p(x|\theta)} \left[\frac{\partial \log p(x|\theta)}{\partial \theta_i} \frac{\partial \log p(x|\theta)}{\partial \theta_j} \right].$$
(1)

This metric enables us to quantify geometric and semantic relationships between representations. For a small parameter shift $d\theta$, the KL divergence is locally quadratic:

$$D_{KL}(p(x|\theta)||p(x|\theta+d\theta)) \approx \frac{1}{2} d\theta^{\top} F(\theta) d\theta.$$
 (2)

Thus, the representation space of an LLM can be understood and analyzed through the lens of Riemannian geometry, with the Fisher metric capturing the structure and sensitivity of its internal representations.

Theorem 1 (Layered Submanifolds). Representations at different semantic levels in LLMs form nested submanifolds within the statistical manifold, each induced by the Fisher information metric and corresponding to a specific granularity of semantic abstraction.

3.2 Multi-Scale Semantic Decomposition

Building on this analysis, we propose that the hidden space of an LLM can be decomposed into three interrelated semantic levels: *global*, *intermediate*, and *local* semantics.

• Global semantic level: Encodes macro-level features such as document topic, overall sentiment, discourse structure, and writing style. This level is typically captured by deep model layers and supports the generation of coherent and consistent long-form text.

- Intermediate semantic level: Focuses on inter-sentential relationships, logical transitions, and mid-range contextual dependencies. Usually represented in the middle layers, this level bridges the global context with local details, supporting logical flow and structured exposition.
- Local semantic level: Captures micro-level details, including word choice, phrase structure, and fine-grained syntax. Primarily handled by shallow layers, this level determines fluency, grammatical correctness, and lexical accuracy.

Although the precise boundaries may vary across architectures, this hierarchical semantic decomposition is ubiquitous. In practice, we identify these levels by analyzing attention spans, inter-layer mutual information, and diagnostic probing tasks.Each manifold encodes language information at a distinct granularity, enabling both fine-grained and high-level understanding.

• Global manifold (\mathcal{M}_G) : Captures documentlevel semantics, discourse structure, and abstract concepts, typically represented by deep layers:

$$\mathcal{M}_G = \{ h_G \in \mathbb{R}^d \mid h_G = f_G(x_{1:T}, c) \}.$$
(3)

• Intermediate manifold (M_I) : Encodes paragraph/sentence-level relationships, logical links, and local discourse, often found in middle layers:

$$\mathcal{M}_{I} = \{ h_{I} \in \mathbb{R}^{d} \mid h_{I} = f_{I}(x_{1:T}, c) \}.$$
 (4)

• Local manifold (\mathcal{M}_L) : Focuses on lexical/syntactic information and micro-level structure, primarily in shallow layers:

$$\mathcal{M}_L = \{ h_L \in \mathbb{R}^d \mid h_L = f_L(x_{1:T}, c) \}.$$
 (5) 334

These manifolds are hierarchically nested: $\mathcal{M}_L \subset \mathcal{M}_I \subset \mathcal{M}_G$, with increasing dimensionality $k_G < k_I < k_L$, reflecting the principle that more abstract representations are often lower-dimensional.

Theorem 2 (Emergent Layer Stratification). For deep Transformer models, there exist boundaries $1 \le l_1 < l_2 \le L$ such that:

- 1. Layers [1, l₁] primarily encode local semantics (M_L);
- 2. Layers $(l_1, l_2]$ encode intermediate semantics (\mathcal{M}_I) ;
- 3. Layers $(l_2, L]$ encode global semantics (\mathcal{M}_G) . 346

These boundaries can be identified by sharp347changes in attention span, mutual information be-
tween layers, and performance in targeted probing
tasks.349

314 315

313

316 317

318 319 320

321

322

323

324

325

326

327

328

329

331

332

335

337

338

339

340

341

343

344

345

054

356

361

363

370

371

373

375

388

390

Information-Theoretic Perspective. To mathematically characterize information flow and semantic organization, we use mutual information between representations:

$$I(h_1; h_2) = \int p(h_1, h_2) \log \frac{p(h_1, h_2)}{p(h_1)p(h_2)} dh_1 dh_2.$$
(6)

For a cross-scale mapping $f : h_1 \mapsto h_2$, we seek to maximize target-relevant information while minimizing redundancy:

$$\max_{f} I(h_2; y) - \beta I(h_1; h_2), \tag{7}$$

where y is the target output and β balances retention and compression.

3.3 Cross-Scale Mapping: Bridging Semantic Layers

The central challenge of multi-scale alignment is to construct mappings between semantic manifolds that *faithfully preserve both geometric structure and semantic content*. These mappings elucidate how LLMs transform micro-level details into macro-level abstractions, revealing information flow and reasoning processes within the model.

Definition 2 (Cross-Scale Mapping). We define two key mappings: $f_{GI} : \mathcal{M}_G \to \mathcal{M}_I$ (global to intermediate), and $f_{IL} : \mathcal{M}_I \to \mathcal{M}_L$ (intermediate to local). Each mapping consists of a *geometric component* (f_{geo}), which maintains topological relationships and minimizes distortion, and an *information component* (f_{info}), which maximizes the retention of critical semantic information (often via mutual information maximization). The overall mapping is given by $f = f_{geo} \circ f_{info}$.

For practical construction, we offer three realizations of increasing expressiveness: (1) Linear projection (solve for W via least squares), (2) Orthogonal mapping (Procrustes analysis to preserve distances/angles), and (3) Nonlinear alignment (multi-layer neural networks for complex relationships). The information component may be instantiated by maximizing mutual information (MINE), applying a variational bottleneck (VIB), or enforcing contrastive learning objectives.

391 3.4 Optimization Framework for Alignment

392To achieve robust cross-scale alignment, we pro-
pose a multi-objective optimization framework
that balances all desired properties. The total loss

is defined as:

 $\mathcal{L}_{\text{total}} = \lambda_{\text{geo}} \,\mathcal{L}_{\text{geo}} + \lambda_{\text{info}} \,\mathcal{L}_{\text{info}} + \lambda_{\text{curv}} \,\mathcal{L}_{\text{curv}},$ 396

395

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

(8)
$$\mathcal{L}_{\text{geo}} = \|f_{GI}(h_G) - h_I\|^2 + \|f_{IL}(h_I) - h_L\|^2,$$
(9)

$$\mathcal{L}_{info} = -I(h_G; f_{GI}(h_G)) - I(h_I; f_{IL}(h_I)), \tag{10}$$

$$\mathcal{L}_{\rm curv} = \int_{\mathcal{M}} K^2 dV \approx \sum_i K_i^2 \Delta V_i.$$
(11)

Here, the hyperparameters λ_{geo} , λ_{info} , and λ_{curv} control the trade-off among structural preservation, information fidelity, and geometric regularity.

Theorem 3 (Bound on Alignment Error). If the mapping functions are Lipschitz-continuous with geometric and information errors bounded by ε_{geo} and ε_{info} , then the total KL divergence satisfies:

$$D_{KL}(p_{\text{true}} \| p_{\text{aligned}}) \le C(\varepsilon_{\text{geo}} + \varepsilon_{\text{info}})$$
 (12)

where C is a constant depending on the manifold's dimension and the local Lipschitz constant.

3.5 Summary

Our framework reveals LLMs' hierarchical information processing across lexical, syntactic and discourse levels, enabling both interpretation and control. Unlike flat approaches, it captures LLMs' true multi-scale nature, with applications in model optimization and safety.

4 **Experiments**

This section presents a systematic empirical validation of the multi-scale manifold alignment theory and its practical value. We design three main experimental groups to assess: (1) the existence and architecture-dependence of semantic stratification; (2) the alignment quality and representational improvements of our multi-scale mapping method; and (3) the effects of interventions and downstream applications. Results confirm the theory's effectiveness and reveal new insights into the internal mechanisms of large language models (LLMs).

4.1 Empirical Analysis of Semantic Stratification

Models and Experimental Setup. We evaluate representative LLMs with varying architectures, In our experiments, we compare four prominent pretrained models: GPT-2 (an autoregressive decoder



Figure 1: Multi-Scale Manifold Alignment Framework

435 with 1.5 B parameters), BERT (a bidirectional encoder with 340 M parameters), RoBERTa (an en-436 hanced encoder with 355 M parameters), and T5 437 (an encoder-decoder architecture with 11B param-438 eters). Experiments use 20,000 documents from 439 the Brown and Reuters corpora, covering various 440 genres and topics. Analyses integrate three met-441 rics: attention span, inter-layer mutual information, 442 and functional probing. All experiments are re-443 peated five times with statistically significant re-444 sults (p<0.05). 445

Table 1: Semantic Layer Distribution across Models

Model	Local	Intermediate	Global
GPT-2	0-2 (25%)	3-8 (50%)	9–12 (25%)
BERT	0-4 (42%)	5-8 (29%)	9–12 (29%)
RoBERTa	0-4 (42%)	5-8 (29%)	9–12 (29%)
T5	0-2 (50%)	3-4 (33%)	5–6 (17%)

Layer Distribution and Architectural Features. 446 As Table 1 shows, autoregressive models (GPT-2) 447 devote half their layers to intermediate scales, while 448 bidirectional models (BERT/RoBERTa) emphasize 449 local processing (>40%). Average attention span 450 grows monotonically with depth, mutual informa-451 tion heatmaps show block structure, and probing 452 453 tasks reveal sharp layer specialization. In BERT, local layers (0-4) excel at POS tagging (F1=0.77), 454 intermediate layers (5-8) peak in sentence relation 455 tasks, and global layers (9-12) dominate topic clas-456 sification (accuracy >0.82). 457

Stability of Semantic Boundaries. Crossvalidation and perturbation tests confirm boundary stability: semantic boundary locations shift minimally (std<0.5 layers) across datasets, input lengths, and injected noise. All three detection methods (attention, mutual information, probing) are consistent, with GPT-2 showing clear boundaries at layers $2\rightarrow3$ (local \rightarrow intermediate) and $8\rightarrow9$ (intermediate \rightarrow global); BERT exhibits similar breaks at $4\rightarrow5$ and $8\rightarrow9$. Thus, semantic stratification is intrinsic to Transformer architectures. 458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

4.2 Cross-Scale Intervention Experiments

Intervention Methods and Metrics. We design four intervention types at each scale: (1) translation $(\mathbf{h}' = \mathbf{h} + \Delta)$, (2) amplification/scaling $(\mathbf{h}' = \alpha \mathbf{h})$, (3) Gaussian noise $(\mathbf{h}' = \mathbf{h} + \epsilon)$, and (4) attention modification. Metrics include lexical diversity, sentence count, mean sentence length, max dependency depth, coherence, and sentiment. Each model-scale-intervention is repeated 30 times, with Wilcoxon tests and Cliff's Delta for effect size.

Scale-Specific Response Patterns. Findings (see Table 3) reveal strong scale-specific effects: *local* interventions shift lexical choices (δ =+0.342); *intermediate* interventions alter sentence structure (sentence count +25%, mean length -19%); *global* interventions impact both lexical diversity (+7.39%) and discourse coherence (δ =-0.238). These patterns confirm functional specialization 488 across scales.

489 490

491

492

493

494

496

497 498

499

502

503

504

506

507

508

510

511

513

514

515

Architecture Dependency and Nonlinear Effects. GPT-2 is highly sensitive to interventions, BERT displays structural robustness, and XLM-R shows unique resilience in sentiment. Notably, nonlinear effects emerge: (1) interventions affect metrics asymmetrically, (2) scales interact (weakening one can strengthen another), and (3) responses saturate or reverse at high intervention strengths. This demonstrates intricate cross-scale regulatory mechanisms.

4.3 Evaluation of Multi-Scale Alignment Methods

Ablation and Setup. Our MSMA framework combines geometric alignment, information alignment, and curvature regularization. We conduct ablation with baselines and component removals (see Table 2). Adam optimizer, $lr=2 \times 10^{-5}$, batch=128, 15 epochs, on GPT-2/BERT.

Table 2: Ablation Settings

Group	Geom.	Info.	Curv.	$\lambda_{ m geo}$	$\lambda_{ m info}$	$\lambda_{ ext{curv}}$
baseline	×	Х	×	0	0	0
full_msma	\checkmark	\checkmark	\checkmark	0.1	0.1	0.01
no_geo	×	\checkmark	\checkmark	0	0.1	0.01
no_info	\checkmark	×	\checkmark	0.1	0	0.01
no_curv	\checkmark	\checkmark	×	0.1	0.1	0
only_geo	\checkmark	×	×	0.1	0	0
only_info	×	\checkmark	×	0	0.1	0
only_curv	×	×	\checkmark	0	0	0.01

Alignment Quality Results. We report KL divergence (distributional difference), mutual information, and distance correlation (geometry preservation) in Table 4. Geometric alignment is crucial for structure preservation, information alignment for content, and curvature for optimization stability. Single components alone are insufficient; multiobjective optimization is essential. BERT, under MSMA, achieves lower KL than GPT-2, indicating a more alignable representation space.

4.4 Summary

The experimental results provide comprehensive validation for the three central hypotheses of the multi-scale manifold alignment theory:**Semantic Stratification:** Large language models spontaneously organize their internal representations into local, intermediate, and global semantic layers, each exhibiting distinct functional specialization;

Table 3: Significant Intervention Effects (p < 0.05, $|\delta| > 0.10$)

Model	Scale	Interv.	Metric	Median $\Delta\%$	Cliff δ	p
GPT-2	Global	Amplify	LexDiv	+7.4	+0.23	0.020
		Amplify	Coher.	0.00	-0.24	0.007
	Inter.	Translate	LexDiv	+6.6	+0.32	0.014
		Amplify	SentCt	+25	+0.24	0.028
		Amplify	MeanSL	-19	-0.27	0.004
		Amplify	MaxDep	-11	-0.20	0.030
	Local	Amplify	LexDiv	+7.3	+0.34	0.005
		Amplify	Sentim	-72	-0.21	0.020
BERT	Inter.	Attn.	SentCt	0.00	+0.27	0.003
XLM-R	Global	Noise	Sentim	-14	+0.24	0.005

Architecture-Dependent Characteristics: differ-525 ent model architectures show unique layer distribu-526 tions and intervention response patterns, reflecting 527 the influence of pre-training objectives and architec-528 tural design choices; and Benefits of Multi-Scale 529 Alignment: integrating geometric and information-530 theoretic constraints within multi-scale alignment 531 leads to significant improvements in model per-532 formance, robustness, and interpretability.Beyond 533 offering a new lens for understanding the inner 534 workings of large language models, the multi-scale 535 manifold alignment theory also provides practical 536 tools for enhancing model capability and reliability. 537 The methods and findings in this study open new 538 pathways for developing more transparent and con-539 trollable language models, representing an impor-540 tant step toward trustworthy artificial intelligence. 541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

5 Conclusion

Key Contributions and Insights. This work presents the Multi-Scale Manifold Alignment (MSMA) framework, a unified theory for interpreting and controlling large language models (LLMs) by decomposing their internal representations into local, intermediate, and global semantic manifolds. Our key findings include:

- Hierarchical Semantic Organization: LLMs inherently structure their representations into three distinct semantic scales—local (word-level), intermediate (sentence-level), and global (discourse-level)—each governing different aspects of language understanding and generation.
- Universal Yet Architecture-Dependent: While
 semantic stratification emerges universally
 across models (GPT-2, BERT, RoBERTa, T5),
 the distribution of layers across scales varies

Table 4: Alignment Results (**KL**: KL-divergence; **MI**: Mutual Information; **DC**: Distance Correlation)

			(a) GPT-2			
Group	$\mathrm{KL}_{g \to m}$	$\mathrm{KL}_{m \to l}$	$\mathrm{MI}_{g \to m}$	$\mathrm{MI}_{m \to l}$	$\mathrm{DC}_{g \rightarrow m}$	$DC_{m \rightarrow l}$
baseline	6955	1.5e4	0.23	0.20	0.97	0.91
full-msma	33	35	1.25	1.49	1.00	1.00
no-curv	39	35	1.35	1.35	1.00	1.00
no-geo	3.4e4	4.2e6	1.29	0.36	0.99	0.97
no-info	57	36	0.80	0.87	1.00	1.00
only-curv	8132	11694	0.24	0.23	0.97	0.90
only-info	5.7e4	5.5e6	1.37	0.38	1.00	0.99
geo-0.1	52	44	0.89	1.08	1.00	1.00
geo-0.2	113	131	0.62	0.78	1.00	1.00
geo-0.3	57	37	1.07	0.80	1.00	1.00
geo-0.4	52	44	0.90	0.74	1.00	1.00
geo-0.5	54	47	0.93	0.84	1.00	1.00
geo-0.6	51	39	0.75	1.07	1.00	1.00
geo-0.7	52	45	0.89	1.09	1.00	1.00
geo-0.8	51	48	0.87	0.84	1.00	1.00
geo-0.9	48	42	1.11	0.79	1.00	1.00
geo-1	70	43	0.92	0.87	1.00	1.00
			(b) BERT			
Group	ИI	KI .	М		-	
-	$\mathbf{KL}_{g \to m}$	$\mathbf{KL}_{m \to l}$	$MI_{g \rightarrow m}$	$MI_{m \rightarrow l}$	$DC_{g \to m}$	$DC_{m \rightarrow l}$
baseline	$KL_{g \rightarrow m}$ 403	$KL_{m \rightarrow l}$ 3840	$MI_{g \rightarrow m}$ 0.06	$\frac{MI_{m \rightarrow l}}{0.13}$	$\frac{\text{DC}_{g \to m}}{0.87}$	$\frac{\text{DC}_{m \to l}}{0.82}$
baseline full-msma	$\frac{\text{KL}_{g \rightarrow m}}{403}$ 0.51	$\frac{\text{KL}_{m \to l}}{3840}$ 1.29	$\frac{\text{MI}_{g \to m}}{0.06}$ 2.89	$\frac{\text{MI}_{m \rightarrow l}}{\begin{array}{c} 0.13\\ 2.64 \end{array}}$	$\frac{\text{DC}_{g \to m}}{0.87}$ 1.00	$\frac{\text{DC}_{m \to l}}{\begin{array}{c} 0.82\\ 1.00 \end{array}}$
baseline full-msma no-curv	$\begin{array}{c} \text{KL}_{g \to m} \\ 403 \\ 0.51 \\ 0.83 \end{array}$	$\frac{\text{KL}_{m \to l}}{3840}$ 1.29 1.04	$MI_{g \to m}$ 0.06 2.89 2.79	$ \begin{array}{r} \mathrm{MI}_{m \to l} \\ 0.13 \\ 2.64 \\ 2.63 \end{array} $	$\begin{array}{c} \mathrm{DC}_{g \to m} \\ \hline 0.87 \\ 1.00 \\ 1.00 \end{array}$	$\frac{\text{DC}_{m \to l}}{\begin{array}{c} 0.82\\ 1.00\\ 1.00 \end{array}}$
baseline full-msma no-curv no-geo	$\begin{array}{c} \text{KL}_{g \to m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \end{array}$	$\frac{\text{KL}_{m \to l}}{3840}$ 1.29 1.04 12367	$ \begin{array}{c} \mathrm{MI}_{g \to m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \end{array} $	$ \begin{array}{c} \text{MI}_{m \to l} \\ \hline 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \end{array} $	$DC_{g \to m}$ 0.87 1.00 1.00 0.82	$ \begin{array}{r} DC_{m \to l} \\ \hline 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ \end{array} $
baseline full-msma no-curv no-geo no-info	$ \begin{array}{c} \text{KL}_{g \to m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \end{array} $	$\begin{array}{c} \text{RL}_{m \to l} \\ \hline 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \end{array}$	$\begin{array}{c} \text{MI}_{g \to m} \\ \hline 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \end{array}$	$ \begin{array}{c} \text{MI}_{m \to l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \end{array} $	$\begin{array}{c} DC_{g \to m} \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \end{array}$	$ \begin{array}{c} \text{DC}_{m \to l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ \end{array} $
baseline full-msma no-curv no-geo no-info only-curv	$\begin{array}{c} \text{KL}_{g \to m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \end{array}$	$\begin{array}{c} \operatorname{RL}_{m \to l} \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \to m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \to l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \end{array}$	$\begin{array}{c} DC_{g \to m} \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \end{array}$	$ \begin{array}{c} DC_{m \to l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \end{array} $
baseline full-msma no-curv no-geo no-info only-curv geo-0.1	$\begin{array}{c} \text{KL}_{g \to m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \end{array}$	$\begin{array}{c} \text{RL}_{m \to l} \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \end{array}$	$\begin{array}{c} \text{MI}_{g \to m} \\ \hline 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \end{array}$	$\begin{array}{c} \text{MI}_{m \rightarrow l} \\ \hline 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \end{array}$	$\begin{array}{c} DC_{g \to m} \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \end{array}$	$\begin{array}{c} \text{DC}_{m \to l} \\ \hline 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2	$\begin{array}{c} \text{KL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \end{array}$	$\begin{array}{c} RL_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.62 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ \hline 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \end{array}$	$\frac{\text{DC}_{g \to m}}{0.87}$ $\frac{0.87}{1.00}$ $\frac{0.82}{1.00}$ 0.87 1.00 1.00	$\frac{\text{DC}_{m \to l}}{0.82}$ 1.00 1.00 0.86 1.00 0.86 1.00 0.86 1.00 1.00
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.3	$\begin{array}{c} \text{KL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \end{array}$	$\begin{array}{c} \text{RL}_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.62 \\ 2.55 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\frac{\text{DC}_{m \to l}}{0.82}$ 1.00 1.00 0.86 1.00 0.86 1.00 1.00 1.00 1.00
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.3 geo-0.4	$\begin{array}{c} \text{KL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \\ 0.50 \end{array}$	$\begin{array}{c} \text{RL}_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.62 \\ 2.55 \\ 2.61 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} {\rm DC}_{m \to l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.3 geo-0.4 geo-0.5	$\begin{array}{c} \text{KL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \\ 0.50 \\ 0.70 \end{array}$	$\begin{array}{c} \text{RL}_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \\ 1.07 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.62 \\ 2.55 \\ 2.61 \\ 2.69 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \\ 2.48 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} {\rm DC}_{m \to l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.3 geo-0.4 geo-0.5 geo-0.6	$\begin{array}{c} \text{RL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \\ 0.50 \\ 0.70 \\ 0.65 \end{array}$	$\begin{array}{c} RL_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \\ 1.07 \\ 0.85 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.65 \\ 2.65 \\ 2.61 \\ 2.69 \\ 2.67 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.49 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} \mathrm{DC}_{m \to l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.3 geo-0.4 geo-0.5 geo-0.6 geo-0.7	$\begin{array}{c} \text{RL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \\ 0.50 \\ 0.70 \\ 0.65 \\ 0.39 \end{array}$	$\begin{array}{c} RL_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \\ 1.07 \\ 0.85 \\ 0.75 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.65 \\ 2.65 \\ 2.65 \\ 2.61 \\ 2.69 \\ 2.67 \\ 2.58 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.49 \\ 2.36 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} \mathrm{DC}_{m \to l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.3 geo-0.4 geo-0.5 geo-0.6 geo-0.7 geo-0.8	$\begin{array}{c} \text{RL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \\ 0.50 \\ 0.70 \\ 0.65 \\ 0.39 \\ 0.39 \end{array}$	$\begin{array}{c} RL_{m \rightarrow l} \\ \hline \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \\ 1.07 \\ 0.85 \\ 0.75 \\ 1.86 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ \hline 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.65 \\ 2.55 \\ 2.61 \\ 2.69 \\ 2.67 \\ 2.58 \\ 2.71 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.49 \\ 2.36 \\ 2.50 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} \text{DC}_{m \rightarrow l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo no-info only-curv geo-0.1 geo-0.2 geo-0.2 geo-0.3 geo-0.4 geo-0.5 geo-0.6 geo-0.7 geo-0.8 geo-0.9	$\begin{array}{c} \text{RL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.47 \\ 0.50 \\ 0.70 \\ 0.65 \\ 0.70 \\ 0.65 \\ 0.39 \\ 0.39 \\ 0.48 \end{array}$	$\begin{array}{c} \mathrm{RL}_{m \rightarrow l} \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \\ 1.07 \\ 0.85 \\ 0.75 \\ 1.86 \\ 0.98 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ \hline 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.65 \\ 2.62 \\ 2.55 \\ 2.61 \\ 2.69 \\ 2.67 \\ 2.58 \\ 2.71 \\ 2.67 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.49 \\ 2.36 \\ 2.50 \\ 2.52 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} \text{DC}_{m \rightarrow l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$
baseline full-msma no-curv no-geo only-curv geo-0.1 geo-0.2 geo-0.3 geo-0.4 geo-0.5 geo-0.6 geo-0.7 geo-0.8 geo-0.9 geo-1	$\begin{array}{c} \text{RL}_{g \rightarrow m} \\ 403 \\ 0.51 \\ 0.83 \\ 3146 \\ 0.42 \\ 423 \\ 0.43 \\ 0.49 \\ 0.37 \\ 0.50 \\ 0.70 \\ 0.65 \\ 0.39 \\ 0.39 \\ 0.48 \\ 0.50 \end{array}$	$\begin{array}{c} \mathrm{RL}_{m \rightarrow l} \\ 3840 \\ 1.29 \\ 1.04 \\ 12367 \\ 1.30 \\ 4310 \\ 1.61 \\ 0.80 \\ 0.87 \\ 1.59 \\ 1.07 \\ 0.85 \\ 0.75 \\ 1.86 \\ 0.98 \\ 0.89 \end{array}$	$\begin{array}{c} \mathrm{MI}_{g \rightarrow m} \\ 0.06 \\ 2.89 \\ 2.79 \\ 0.03 \\ 2.75 \\ 0.07 \\ 2.65 \\ 2.62 \\ 2.55 \\ 2.61 \\ 2.69 \\ 2.67 \\ 2.58 \\ 2.71 \\ 2.67 \\ 2.78 \end{array}$	$\begin{array}{c} \mathrm{MI}_{m \rightarrow l} \\ 0.13 \\ 2.64 \\ 2.63 \\ 0.05 \\ 2.51 \\ 0.11 \\ 2.48 \\ 2.64 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.48 \\ 2.49 \\ 2.36 \\ 2.50 \\ 2.52 \\ 2.53 \end{array}$	$\begin{array}{c} {\rm DC}_{g \rightarrow m} \\ \\ 0.87 \\ 1.00 \\ 1.00 \\ 0.82 \\ 1.00 \\ 0.87 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$	$\begin{array}{c} \text{DC}_{m \rightarrow l} \\ 0.82 \\ 1.00 \\ 1.00 \\ 0.86 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \\ 1.00 \end{array}$

with architecture (e.g., autoregressive models prioritize intermediate semantics, while bidirectional models emphasize local features).

561

562

565

566

567

569

573

574

575

576

577

578

580

581

- Alignment Enables Control: Our framework successfully bridges semantic scales via geometric preservation, information retention, and manifold smoothness, achieving nearperfect alignment (99% KL reduction, 5–7× mutual information gain) and enabling precise interventions (e.g., editing lexical choice without disrupting coherence).
- Functional Specialization Proven: Interventions confirm scale-specific roles—local manipulations alter word choice, intermediate adjustments reshape sentence structure, and global modifications impact both discourse and fine-grained features.

Broader Implications. The MSMA framework bridges the gap between theoretical interpretability and practical control in LLMs by elucidating cross-scale information flow, enabling three key applications: (1) bias mitigation through targeted manifold editing of stereotypical associations, (2) robustness enhancement via curvature-constrained regularization that preserves model stability, and (3) controlled generation with fine-grained manipulation of output properties such as formality and discourse coherence. This unified approach transforms theoretical insights into actionable model improvement strategies. 582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

6 Limitations

Despite the significant progress afforded by the Multi-Scale Manifold Alignment (MSMA) framework in elucidating the internal mechanisms of large language models, several limitations remain. First, the computational cost of MSMA is substantial: estimating mutual information and manifold curvature across every layer of models with hundreds of billions of parameters (e.g., GPT-4, PaLM) demands considerable resources. Second, the semantic boundaries we detect may blur in architectures that employ hybrid or sparse attention mechanisms, necessitating tailored boundary-detection strategies for non-standard designs. Third, although our experiments used general-purpose text corpora, the layerwise semantic organization may differ in highly specialized domains (e.g., medical or legal texts) or in fine-tuned models, calling for cross-domain validation and adaptation of the framework.

Moreover, our theoretical analysis relies on simplifying assumptions—such as Markovian transitions and conditional independence among representation scales—that hold only approximately in practice, especially in the presence of residual connections and cross-attention. We have not yet established a direct correspondence between model representations and human cognitive processes; integrating insights from neuroscience and psycholinguistics could strengthen this link. In our intervention studies, we observed that effect sizes sometimes attenuate or behave non-linearly over long generation sequences, a dynamic phenomenon not fully captured by the current theory.

Finally, while we evaluated alignment quality using KL divergence, mutual information, and distance-based metrics, these measures may not fully reflect the richness of semantic content or downstream task performance. Likewise, existing visualization tools struggle to convey highdimensional structure to non-technical audiences.

734

735

736

737

738

682

683

684

632Developing more comprehensive evaluation met-633rics and interactive visual interfaces will be critical634for broadening MSMA's applicability and inter-635pretability.

7 Acknowledgements

636

659

663

673

674

675

676

677

678

679

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 641 in screening research literature in related fields. All AI-polished text content has been strictly reviewed 643 by the author to ensure that it complies with academic standards and is accompanied by accurate citations. The core research ideas, method design 646 and conclusion derivation of this article were independently completed by the author, and the AI 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 651 academic rigor, data authenticity and citation integrity of the full text, and hereby declares that the generative AI tool is not a co-author of this study.

References

- Yonatan Belinkov and Mark O. Riedl. 2022. Towards mechanistic transparency of language models. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5327–5343.
- Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. 2013. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35:1798–1828.
- 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 A. Creel, Jared Davis, Dora Demszky, and 66 others. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901.

- Hila Chefer, Shir Gur, and Lior Wolf. 2021. Transformer interpretability beyond attention visualization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 782–791.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, and 48 others. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Chia-Yi Chuang, Ananya Kumar, Percy Liang, and Christopher Re. 2023. Tldr: Transfer learning via distillation of pre-trained representations. In Advances in Neural Information Processing Systems, volume 36.
- Erik Daxberger, Sarthak Mittal, Leonard Berrada, Milad Alizadeh, Stephen Roller, Pierre H. Richemond, Samuel L. Smith, Aaron Sisto, Jared Kaplan, Dzmitry Bahdanau, Arthur Conmy, Andreas Steiner, Andrej Karpathy, Antonio Norelli, Ellie Pavlick, Marco Baroni, Ethan Perez, Jared Tanner, Luke Metz, and 6 others. 2023. Representation manifolds of images in generative models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 18799–18825.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Joseph, Nova DasSarma, Tom Henighan, Scott Sievert, Ben Mann, Sam McCandlish, Tom Brown, Dario Amodei, Jared Kaplan, Jack Clark, and Chris Olah. 2022. Toy models of superposition. *Transformer Circuits Thread*, 1.
- Angela D. Friederici. 2011. The brain basis of language processing: from structure to function. *Physiological Reviews*, 91:1357–1392.
- Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2023. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pages 1049–1077.
- Asma Ghandeharioun, Avi Caciularu, Adam Kalai, Rosina Weber, Han Liu, and Mor Geva. 2023. Patchscopes: A unifying framework for inspecting hidden representations of language models. In *Advances in Neural Information Processing Systems*, volume 36.
- Aditya Grover, Johannes Gasteiger, Petar Veličković, Stephan Günnemann, Frédéric Ohme, Charles Harris, Pietro Liò, Yoshua Bengio, and José Miguel Hernández-Lobato. 2023. Geometric deep learning on molecular representations. In *Advances in Neural Information Processing Systems*, volume 36.
- William Gurnee, Alyssa Loo, Cassidy Laidlaw, Cathy Wu, and Jacob Andreas. 2023. Language models as agent models: Mechanistic analysis beyond the

844

845

information stream. In *Proceedings of the 40th International Conference on Machine Learning*, pages 11920–11935.

739

740

741

742

743

744

745

746

747

748

749

750

751

753

754

757

762

765 766

769

770 771

772

773

774

775

776

778

790

791

- Michael Hahn and Dan Jurafsky. 2023. Tracking information flow in large language models. *Transactions of the Association for Computational Linguis tics*, 11:1225–1242.
- David Hernandez, Sabrina J. Mielke, Marta Recasens, James Shen, Ilker Kesen, Ethan Aoyama Liu, and Benjamin Van Durme. 2023. Lingoqa: Multi-step reasoning for language models through natural language constraints. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16334–16347.
- Shun ichi Amari and Hiroshi Nagaoka. 2007. Methods of information geometry. *Translations of Mathematical Monographs*, 191.
- Tianyu Li, Colin Raffel, and Julian Michael. 2023. Decoding representations with semantic classifiers: Bridging fine-tuning and prompt approaches. *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pages 4023– 4039.
- Xiaoshi Liu, Jingbin Cao, Chang Liu, Quan Chen, Zhen Chen, Jun Wang, Wang Zhou, Li Guanbin, Saining Xie, and Michael Jordan. 2023. Mind the gap: Understanding the modality gap in multi-modal contrastive learning. In *Advances in Neural Information Processing Systems*, volume 36.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. In *Advances in Neural Information Processing Systems*, volume 35.
- OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Maithra Raghu, Thomas Unterthiner, Shibani Santurkar, Xiaohua Zhai, Basil Mustafa, Simon Kornblith, Alexey Dosovitskiy, and Neil Houlsby. 2022. Vision models are more robust and fair when pretrained on uncurated images without supervision. In *International Conference on Learning Representations*.
- Gabriele Sarti, Nora Kassner, and Gemma Boleda. 2023. Latent understanding of actions in attention heads of large language models. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 781–802.
- Jinyoung Seo, Annie Chen, Ian Covert, Anna Huang, Yejin Choi, and Xiang Lisa Li. 2023. Do language models have coherent mental models of everyday things? In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11336–11350.

- Arvind Singh and Hal Daumé III. 2023. I would ask if that makes sense but i'm hallucinating: Language models analyze and document their decision making.
 In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13773–13790.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. *International Conference on Learning Representations*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timoth'ee Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Jesse Vig. 2019. Analyzing the structure of attention in a transformer language model. *Proceedings of the* 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, pages 63–76.
- Feng Wang, Huaping Liu, Di Hu, Baoxing Qiu, Guocheng Niu, and Fengwei Zhou. 2023. Align, manipulate and learn: Exploiting geometric approaches for self-supervised representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1383–1392.
- Kayo Xie, Yudong Zhang, Jiahui Zhang, Sean Frey, Qingyan Guo, and Ding Zhao. 2023. A mechanistic dissection of the attention function in large language models. *Advances in Neural Information Processing Systems*, 36.
- Yixuan Xu, Peihao Wang, Mingyang You, Jingbo Shang, and Julian McAuley. 2023. Understanding hidden information in bert through mutual information. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 747–761.
- Ziqian Zhang, Cheng Lu, Adrian Weller, Matthew J. Johnson, Wenbo Gong, Da Yu, Hongzhou Lin, and Xuechen Li. 2023. White-box transformers via selfdistillation: History-free, analytically tractable, and certified. In *Advances in Neural Information Processing Systems*, volume 36.

A Experimental Setup and Analysis for Semantic-Scale Identification

A.1 Experimental Design

Research Questions: We test three central hypotheses of the Multi-Scale Manifold Alignment (MSMA) theory: (1) Do Transformer layers form identifiable local/intermediate/global semantic scales? (2) How do architecture and pre-training

849 850

851

852 853

854

858

864

869

870

objectives influence these scales? (3) Do targeted interventions yield the scale-specific effects predicted by MSMA?

A.1.1 Models

We evaluate representative large language models as shown in Table 5:

Table 5: Evaluated models.

Model	Architecture	Params	Pretrain Objec- tive
GPT-2	Autoregressive Decoder	1.5B	Next- token Predic- tion
BERT	Bidirectional Encoder	340M	Masked LM
RoBERTa	Enhanced BERT Encoder	355M	Dynamic Masked LM
Τ5	Encoder–Decoder	11B	Sequence- to- Sequence

A.1.2 Data Resources

We construct a balanced corpus of 20,000 samples from three sources (Table 6):

Table 6: Corpus composition and average sample length.

Source	# Samples	Avg. Length (tokens)
Brown (15 genres)	6,667	293.5
Reuters (8 topics)	6,667	318.2
GPT-2 academic synth	6,666	352.8

Brown: 15 genres, classic written English; **Reuters:** 8 topic categories, global news; **GPT-2:** Academic-style synthetic documents generated from 68 field prompts and manually filtered for quality.

A.1.3 Feature Hierarchies

Global: Genre, source, LDA topic, stylistic markers.

Intermediate: Mean sentence length, clause count, lexical complexity, topic coherence.

Local: Token length variance, function word ratio, POS/dependency distribution, sentiment score.

A.1.4 Scale Identification Methods

• Attention Patterns: Compute mean span $d_{\text{attn}}^{(\ell)} = \frac{1}{H} \sum_{h} \sum_{i,j} A_{i,j} |i - j|$ and entropy $H_{\text{attn}}^{(\ell)}$.



(a) Mean attention span by layer across models.



(b) Attention span distance heatmap.



(c) Attention entropy by layer.

Figure 2: Comprehensive attention profile analysis for four Transformer models: (a) Layerwise mean attention span, (b) Attention span heatmap, (c) Entropy of attention by layer.

• **Representation Similarity:** KL divergence, mutual information (k-NN, PCA to 50D).

871

872

873

874

875

876

877

878

879

881

882

884

- **Probing Tasks:** Layerwise SVMs for POS/dependency (local), nextsentence/paragraph (intermediate), topic/genre (global).
- Voting Integration: $S_{\text{scale}} = 0.4 \text{ Probe} + 0.4 \text{ Attn} + 0.2 \text{ MI}$, followed by continuity smoothing.

A.2 Layered Structure Revealed by Attention Patterns

Fig. 2b(a) shows the mean attention span by layer. In GPT-2, span rises from 12.5 (layer 0) to 36.2 (layer 12), clustering as local (0–2, median <15),



(b) KL divergence across models.

Figure 3: Comparative analysis of information metrics: (a) mutual information and (b) KL divergence for different Transformer models.

intermediate (3–8, 15–30), global (9–12, >30). BERT/RoBERTa show a smooth span rise, from 17.3 (layers 0–4) to above 30 (layers 9–12). T5 (six layers) exhibits clear separation: encoder spans grow from 12.4 to 27.8; decoder from 14.2 to 31.5. Spearman correlations (span vs. depth) all exceed 0.85 (p < 0.01), confirming span as a semantic scale indicator.

889

896

900

Fig. 2c(b) plots attention entropy per layer. GPT-2 shows a "U-shaped" curve: peak entropy in layers 0–1, sharp drop at 7, then global expansion. BERT/RoBERTa have entropy dips in 5–8, matching intermediate layers. T5's curve is flatter but shows encoder dip. These profiles confirm model-specific functional hierarchies as predicted by MSMA.



(b) Probing performance by layer across models.

Figure 4: Comparative analysis of probing performance: (a) overall results across models, (b) results by layer.

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

A.3 Representation Similarity Confirms Semantic Boundaries

Fig. 3a(b): Layerwise KL divergence for each model, with light colors (low KL) marking high similarity, dark (high KL) marking sharp transitions. GPT-2 shows three clear blocks (local/intermediate/global): KL jumps from 9.1 to 19.6 (layers $2\rightarrow 3$), and from 6.7 to 17.9 ($8\rightarrow 9$). BERT and RoBERTa display similar boundaries. All jumps are significant (Z > 2.0, p < 0.01).

Fig. 3(a): Layerwise MI, quantifying shared information. BERT's MI matrix forms three modules $\{0-4, 5-8, 9-12\}$, with within-module MI $\sim 40\%$ higher than between-module MI. RoBERTa/T5 are similar; GPT-2's MI estimates are noisier but consistent with its KL blocks. These results confirm three functional modules per model.

A.4 Probing Tasks Validate Functional Specialization

Fig. 4a(a) shows layerwise probing. BERT exhibits three regimes: layers 0–4 excel on local tasks (F1 rises $0.18 \rightarrow 0.77$), 5–8 peak on intermediate, 9– 12 on global (acc. >0.82). GPT-2 achieves nearperfect local F1 (~0.99), but lower global accu-

993

994

995

996

997

998

999

1001

1002

1003

1004

1005

1007

1009

1010

1011

1012

1013

1014

1015

1017

1018

974

975

976

977

978

925 926

930

931

932

933

934

936

937

938

941

942

943

944

945

950

951

952

955

957

960

961

962

963

965

967

968

969

970

971

973

racy (~ 0.53), reflecting its autoregressive nature. RoBERTa and T5 show architecture-specific stratification. Across all models, probing peaks align 927 closely with attention/MI boundaries, verifying that 928 each semantic scale fulfills its predicted function.

A.5 Intervention Experiments

We test MSMA's causal predictions by perturbing hidden representations at three scales (local/intermediate/global) in each model, using: Translation $(\mathbf{h}^{\prime(\ell)} = \mathbf{h}^{(\ell)} + \Delta)$, Scaling $(\mathbf{h}^{\prime(\ell)} =$ $\alpha \mathbf{h}^{(\ell)}$), Noise $(\mathbf{h}^{\prime(\ell)} = \mathbf{h}^{(\ell)} + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2 I)),$ Attention modification $(A'_{i,j}^{(\ell,h)} = f_{\text{att}}(A_{i,j}^{(\ell,h)}))$. We measure effects on: lexical diversity, sentence count, mean sentence length, max dependency depth, coherence, and sentiment.

A.5.1 Statistical Analysis

Each model-scale-intervention is repeated 30 times (over 5,000 samples). We use Wilcoxon signed-rank tests (p < 0.05, FDR-corrected) and Cliff's delta (small effect $|\delta| > 0.147$). Bootstrap (1,000), leave-one-out, and power analysis confirm robustness.

Note: *p<0.05, **p<0.01 (FDR). Cliff's delta: +=increase, -=decrease.

A.5.2 Intervention Effect Analysis

Multi-dimensional interventions reveal unique responses by architecture. GPT-2 shows marked lexical sensitivity: local scaling gives largest diversity effect ($\delta_{\text{max}} = +0.342, p < 0.01$); global scaling increases diversity by +7.39% but reduces coherence ($\delta = -0.238$). Intermediate translation increases diversity +6.60%, scaling increases sentence count +25%, and shortens mean sentence length -19%. All are as MSMA predicts: local controls lexicon, intermediate controls sentence structure, global controls discourse. Even small perturbations shift GPT-2's output, showing its autoregressive nature and reliance on precise representations.

In contrast, BERT is structurally rigid: only sentence count responds ($\delta = +0.269, p < 0.01$), while other metrics stay constant, reflecting stable bidirectional encoding. XLM-R is sentimentrobust—global noise shifts sentiment by -13.6% $(\delta = +0.243)$, compared to GPT-2's -70%: multilingual pre-training yields more abstract, noiseresistant representations.

Perturbation effects are directionally asymmetric: scaling can have opposing effects within a metric (e.g., global scaling increases diversity, lowers syntactic complexity); scaling down at one scale can enhance another's properties; increasing attention may suppress some attributes, revealing nonmonotonic attention-content relationships.

Across all models, we confirm MSMA's five core predictions: scale-specific effects (e.g., local diversity $\delta = +0.342$, intermediate structure $\delta = +0.239$, global coherence $\delta = -0.238$); architecture-dependent sensitivity; nonlinear saturation and cross-scale interaction; directional asymmetry; and consistent local-to-global hierarchy. These convergent findings validate MSMA as an explanatory and predictive framework for Transformer language generation.

A.6 MSMA Method Implementation Details

We detail implementation and hyperparameters for multi-scale manifold alignment. The process is multi-stage: first, semantic boundaries are detected; next, cross-scale mappings are constructed and optimized.

A.6.1 Layer Identification Algorithm

We employ an ensemble approach, integrating attention, mutual information, and probing evidence. For model M with L layers:

Multi-Scale Manifold Alignment B **Theory Proofs**

This appendix provides the complete mathematical proofs for the multi-scale manifold alignment theory. Proofs are organized into six main parts: information geometry preliminaries, KL divergence upper bound, mutual information lower bound, local convergence, mapping implementation, and hierarchical Markov properties with error decomposition.

B.1 Preliminaries and Assumptions

Information Geometry and Statistical B.1.1 Manifolds

Definition B.1.1 (Statistical Manifold). Given a family of probability distributions $\{p(x|\theta)\}$ parameterized by $\theta \in \Theta$, with $x \in \mathcal{X}$, the statistical manifold \mathcal{M} is defined as:

$$\mathcal{M} = \{ p(x|\theta) : \theta \in \Theta \}$$
 1016

Definition B.1.2 (Fisher Information Matrix). For $p(x|\theta)$, the Fisher information matrix is:

$$g_{ij}(\theta) = \mathbb{E}_{p(x|\theta)} \left[\frac{\partial \log p(x|\theta)}{\partial \theta_i} \frac{\partial \log p(x|\theta)}{\partial \theta_j} \right]$$
 1019

Model	Scale	Intervention	Metric	Median Change (%)	Cliff's δ	<i>p</i> -value	Sig.
GPT-2	Global	Scale up	Lexical diversity	+7.39	0.232	0.020	*
GPT-2	Global	Scale up	Coherence score	0.00	-0.238	0.007	**
GPT-2	Global	Scale down	Lexical diversity	+6.78	0.272	0.017	*
GPT-2	Intermed.	Translate	Lexical diversity	+6.60	0.316	0.014	*
GPT-2	Intermed.	Scale up	Sentence count	+25.00	0.239	0.028	*
GPT-2	Intermed.	Scale up	Mean sent. length	-19.04	-0.266	0.004	**
GPT-2	Intermed.	Scale up	Max dep. depth	-11.11	-0.203	0.030	*
GPT-2	Intermed.	Scale down	Lexical diversity	+5.84	0.211	0.016	*
GPT-2	Intermed.	Scale down	Max dep. depth	-11.11	-0.192	0.037	*
GPT-2	Intermed.	Attn	Lexical diversity	+4.55	0.195	0.028	*
GPT-2	Intermed.	Attn	Sentiment score	-80.09	-0.246	0.004	**
GPT-2	Local	Translate	Coherence score	0.00	-0.180	0.020	*
GPT-2	Local	Scale up	Lexical diversity	+7.27	0.342	0.005	**
GPT-2	Local	Scale up	Sentiment score	-71.84	-0.206	0.020	*
GPT-2	Local	Scale down	Lexical diversity	+5.62	0.276	0.015	*
GPT-2	Local	Scale down	Coherence score	0.00	-0.180	0.037	*
BERT	Global	Noise	Sentence count	0.00	0.154	0.046	*
BERT	Intermed.	Translate	Sentence count	0.00	0.154	0.033	*
BERT	Intermed.	Attn	Sentence count	0.00	0.269	0.003	**
XLM-R	Global	Noise	Sentiment score	-13.58	0.243	0.005	**
XLM-R	Intermed.	Scale up	Sentiment score	-1.03	0.104	0.046	*
XLM-R	Local	Attn	Sentiment score	-10.79	0.149	0.043	*

Table 7: Significant intervention effects across models (p < 0.05, $|\delta| > 0.10$). Median changes (%) are relative to baseline.

The Fisher matrix induces a Riemannian metric on \mathcal{M} , enabling distances, geodesics, and curvature. For infinitesimal $d\theta$, the KL divergence is locally quadratic:

Lemma B.1.1. For parameter θ and small $d\theta$,

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1042

$$D_{\mathrm{KL}}(p(x|\theta)||p(x|\theta+d\theta)) = \frac{1}{2}d\theta^T g(\theta)d\theta + O(||d\theta|)$$

Proof sketch. By Taylor expansion and $\mathbb{E}_{p(x|\theta)} \left[\frac{\partial \log p(x|\theta)}{\partial \theta_i} \right] = 0$, this follows from the Fisher matrix definition and KL divergence Taylor expansion.

B.1.2 Multi-Scale Representation in Transformers

Assumption B.1.1 (Representation Hierarchy). For a Transformer with L layers, there exist $1 \le l_1 < l_2 \le L$ such that:

- Layers $[1, l_1]$: local semantics, manifold \mathcal{M}_L
- Layers $(l_1, l_2]$: intermediate, \mathcal{M}_I
- Layers $(l_2, L]$: global, \mathcal{M}_G

Assumption B.1.2 (Hierarchical Information Flow). Information primarily flows $\mathcal{M}_L \rightarrow \mathcal{M}_I \rightarrow \mathcal{M}_G$, with local computation at each layer, consistent with residual-based Transformer design and confirmed experimentally. Assumption B.1.3 (Conditional Independence).1043Given h_G , intermediate representation h_I is conditionally independent of unrelated factors; likewise,
given h_G and h_I , local h_L is conditionally independent:1044104510461046104610471047

 $\|\theta\|^3)^{p(h_I|h_G, z)} \approx p(h_I|h_G), \quad p(h_L|h_G, h_I, z) \approx p(h_L|h_G, h_I)$

where z denotes external nuisance variables. 1049

B.2 Proof of KL Divergence Upper Bound

Consider mappings f_{GI} (global-to-intermediate) 1051 and f_{IL} (intermediate-to-local). 1052

Lemma B.2.1 (Local Mapping Error Decomposition). For f_{GI} , total error decomposes as:

$$\mathcal{E}_{G \to I} = \mathcal{E}_{G \to I}^{\text{geo}} + \mathcal{E}_{G \to I}^{\text{info}}$$
1055

1053

1054

1058

1059

1060

1061

with
$$\mathcal{E}_{G \to I}^{\text{geo}} = \|f_{GI}(h_G) - h_I\|^2$$
, $\mathcal{E}_{G \to I}^{\text{info}} = D_{\text{KL}}(p(h_I|h_G)\|p(f_{GI}(h_G)|h_G))$.

Proof sketch. By the chain rule of KL and Fisher norm local approximation, as in Lemma A.1, the total error splits into a geometric and an information-theoretic part.

Assumption B.2.1 (Lipschitz Continuity). Mappings f_{GI} , f_{IL} are Lipschitz: $||f_{GI}(h_G^1) - 1063$ $f_{GI}(h_G^2)|| \le L_{GI} ||h_G^1 - h_G^2||$, and similarly for 1064 f_{IL} .

1087

1089

1091

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1116

1117

1118

1119

1120

1121

1122

1123

1124

B.3.2 Information Preservation in Cross-Scale Mapping

where q(z|x), q(z) are variational approximations.

(Proof: KL non-negativity and standard VIB deriva-

tion.)

Theorem B.4 (Mutual Information Preservation). *Minimizing information loss* $\mathcal{L}_{info} = -I(h_G; f_{GI}(h_G)) - I(h_I; f_{IL}(h_I))$ ensures:

- Conditional entropy $H(h_G|f_{GI}(h_G))$, $H(h_I|f_{IL}(h_I))$ minimized;
- Critical information for predicting y is preserved across mappings.

Proof sketch. By mutual information definition, maximizing $I(h_G; f_{GI}(h_G))$ minimizes $H(h_G|f_{GI}(h_G))$. Data-processing inequality shows $I(h_G; y) \ge I(f_{GI}(h_G); y)$; minimizing their difference ensures $f_{GI}(h_G)$ preserves h_G 's information about y.

B.4 Proof of Local Convergence

B.4.1 Existence of Local Minimum

Theorem B.5 (Existence of Local Minimum). For total loss $\mathcal{L}_{total} = \lambda_{geo} \mathcal{L}_{geo} + \lambda_{info} \mathcal{L}_{info} + \lambda_{curv} \mathcal{L}_{curv}$, if \mathcal{L}_{total} is smooth with bounded second derivatives, stochastic gradient descent with proper step size converges to a local minimum with high probability.

Proof sketch. By standard stochastic optimization analysis: for parameter θ_t , learning rate $\eta_t = \eta_0/\sqrt{t}$, bounded gradient variance, and Lipschitz gradients, we have

$$\mathbb{E}[\|\nabla_{\theta} \mathcal{L}_{\text{total}}(\bar{\theta}_T)\|^2] \to 0 \text{ as } T \to \infty.$$
 1115

B.4.2 Effect of Curvature Regularization

Theorem B.6 (Stability of Curvature Regularization). The curvature regularization $\mathcal{L}_{curv} = \int_{\mathcal{M}} K^2 dV$ improves loss landscape smoothness and bounds total alignment distortion by controlling the maximum curvature K_{max} via λ_{curv} .

Proof sketch. By Rauch comparison, for points $p, q \in \mathcal{M}$ with geodesic γ ,

$$d(f(p), f(q)) \le d(p, q) \exp\left(\int_{\gamma} K(s) ds\right).$$
 1125

tion **Input:** Model *M*, number of layers *L*, test corpus \mathcal{D} **Output:** Boundaries l_1 (local \rightarrow intermediate), l_2 (intermediate→global) for each layer $l \in \{1, \ldots, L\}$ do Compute mean attention span S_l ; for each $l \in \{1, ..., L-1\}$ do Compute difference $\Delta S_l = S_{l+1} - S_l$; for each pair (i, j) of layers do Compute mutual information I_{ij} ; Build MI matrix *I*; for each layer l and each task t do Evaluate task accuracy P_{I}^{t} ; Compute gradient $\nabla P_l^t = P_{l+1}^t - P_l^t$; for each l do Compute boundary score $B_l = \alpha \Delta S_l + \beta \Delta I_l + \gamma \sum_t w_t \nabla P_l^t;$ Identify two highest B_l as boundaries l_1, l_2 ; **Parameters:** $\alpha = 0.4, \beta = 0.4, \gamma = 0.2,$ w_t is task-specific weight. Apply smoothing and 5-fold cross-validation for stability;

Algorithm 1: Semantic Boundary Detec-

Theorem B.1 (KL Divergence Upper Bound). *Un*der the above, for true and aligned distributions,

1066

1067

1068

1069

1070

1073

1074

1075

1076

1078

1079

$$D_{\mathrm{KL}}(p_{\mathrm{true}} \| p_{\mathrm{aligned}}) \le C(\varepsilon_{\mathrm{geo}} + \varepsilon_{\mathrm{info}})$$

where $\varepsilon_{\text{geo}}, \varepsilon_{\text{info}}$ sum geometric and information errors; C depends on manifold dimension and Lipschitz constants.

Proof sketch. Apply KL chain rule, triangle inequality, error propagation under Lipschitz continuity, and Lemma A.2 to bound each mapping's KL by geometric and information terms. \Box

B.3 Mutual Information Lower Bound

B.3.1 MINE and VIB Variational Bounds

Theorem B.2 (MINE Lower Bound). For X, Y,

$$I(X;Y) \ge \mathbb{E}_{p_{XY}}[T_{\phi}(x,y)] - \log \mathbb{E}_{p_X p_Y}[e^{T_{\phi}(x,y)}]$$

1080with T_{ϕ} a neural network. (Proof: Donsker-1081Varadhan representation for KL divergence.)

1082**Theorem B.3** (VIB Lower Bound). Given encoder1083p(z|x),

1084
$$I(X;Z) \ge \mathbb{E}_{p(x)p(z|x)}[\log q(z|x)] - \mathbb{E}_{p(z)}[\log q(z)]$$

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1140

1141

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1161

1164

Cauchy-Schwarz gives

$$|\int_{\gamma} K(s)ds| \le L^{1/2} \left(\int_{\mathcal{M}} K^2 dV\right)^{1/2}$$

where L is geodesic length. Thus, minimizing \mathcal{L}_{curv} tightens distortion bounds and improves convergence by conditioning the Hessian.

Curvature regularization is especially important for generalization and stability in cross-scale mapping, as demonstrated by smoother training curves and improved robustness.

B.5 Proof of Mapping Function Implementation (Continued)

Corollary A.1 (MINE Implementation). Us-1137 ing the MINE framework, information mapping 1138 is achieved by maximizing: 1139

$$\begin{aligned} \max_{\theta,\phi} & \mathbb{E}_{p(h_G, f_{\text{info}}(h_G; \theta))} \left[T_{\phi}(h_G, f_{\text{info}}(h_G; \theta)) \right] \\ & - \log \mathbb{E}_{p(h_G)p(f_{\text{info}}(h_G; \theta))} \left[e^{T_{\phi}(h_G, f_{\text{info}}(h_G; \theta))} \right] \end{aligned}$$

where T_{ϕ} is a neural network estimator for mutual 1142 information. 1143

> *Proof:* This is a direct application of Theorem A.2, using MINE's variational lower bound with our representations and mappings. Both θ and ϕ are optimized jointly to preserve maximal information.

B.6 AHierarchical Markov Properties and Error Decomposition

Theorem A.9 (Hierarchical Markov Property). Suppose the joint distribution of Transformer representations decomposes as:

$$p(h_G, h_I, h_L | C) = p(h_G | C) \cdot p(h_I | h_G, C) \cdot p(h_L$$

where C is the context. Then, given h_G , h_I is conditionally independent of irrelevant factors; similarly, given h_G , h_I , h_L is conditionally independent of other factors.

Proof: By information-theoretic conditional independence and the hierarchical processing structure, 1160 information mainly flows along layers, with each abstracting its input. 1162

For irrelevant factors Z, 1163

$$I(h_I; Z|h_G) = H(h_I|h_G) - H(h_I|h_G, Z) \approx 0$$

since h_G is an information bottleneck; thus, Z con-1165 tributes little to h_I . Similarly, $I(h_L; Z | h_G, h_I) \approx$ 1166

0. Hence, the Markov structure enables decompo-1167 sition into local mappings, simplifying alignment. 1168 1169

Theorem A.10 (Error Accumulation Theo-1170 **rem**). Let mapping errors at each level be ε_G , ε_I , 1171 ε_L . Under the hierarchical Markov assumption, 1172 total KL divergence error is: 1173

$$\mathcal{E}_{\text{total}} \approx \varepsilon_G + \varepsilon_I + \varepsilon_L$$
 1174

1175

1177

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1202

1203

1204

1205

1206

1207

1209

Proof: Consider the total mapping KL error:

$$\mathcal{E}_{\text{total}} = D_{\text{KL}}(p(h_G, h_I, h_L) \| p(h_G, f_{GI}(h_G), f_{IL}(f_{GI}(h_G))))$$

By the chain rule and Markov property:

$$\mathcal{E}_{\text{total}} = D_{\text{KL}}(p(h_G) \| p(h_G))$$
1178

$$+ \mathbb{E}_{h_G}[D_{\mathrm{KL}}(p(h_I|h_G) \| p(f_{GI}(h_G)|h_G))]$$
1179

$$+ \mathbb{E}_{h_G,h_I}[D_{\mathrm{KL}}(p(h_L|h_G,h_I) \| p(f_{IL}(h_I)|h_I))]$$
 118

The first term is 0, the second is ε_G , and the third 1181 simplifies to ε_I by conditional independence. Error 1182 from f_{GI} propagates through f_{IL} , but is bounded 1183 by Lipschitz continuity, and can be absorbed into 1184 ε_L . Hence, total error is approximately additive. \Box 1185

B.7 Theoretical Summary and Discussion

Main Results Our theoretical analysis yields:

- KL upper bound (Thm. A.1): Alignment KL error is bounded by a weighted sum of geometric and informational errors, supporting multi-objective optimization.
- Mutual information preservation (Thm. A.4): Maximizing mutual information ensures that critical semantic information for prediction is retained across scales.
- $|h_I, h_{C}, C \partial cal$ convergence (Thm. A.5, A.6): Multiobjective optimization converges locally; curvature regularization improves stability.
 - Optimal mapping construction (Thm. A.7, A.8): Theoretically optimal constructions for geometric and information mappings, with practical implementation.
 - Error decomposition (Thm. A.10): Under the Markov structure, total error decomposes into the sum of scale-wise mapping errors.

These provide a rigorous mathematical foundation for multi-scale manifold alignment.

Key Assumptions and Limitations Our proofs rely on several key assumptions:

 Lipschitz continuity: Assumed for mappings, usually satisfied locally for neural networks, reinforced via regularization and gradient clipping.

1214

1215

1216 1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1245

1246

1247

1248

1249

1250

1251

1252

1253

- *Hierarchical Markov assumption*: Conditional independence is assumed; real models may have residual dependencies, but experiments show the approximation is sufficiently accurate.
 - *Curvature regularization*: The choice of λ_{curv} is crucial. Over-regularization may cause underfitting, under-regularization may not improve stability. Empirically, we tune this via validation.

Future work may relax these assumptions or extend the theory to richer dependency structures.

Experimental Correspondence Our theoretical predictions closely match empirical results:

- KL divergence scales linearly with geometric/information errors (Thm. A.1); full MSMA (multi-objective) outperforms singleobjective baselines.
- Curvature regularization improves optimization stability, especially early in training. Methods without it show higher oscillation.
- Different architectures exhibit varying hierarchical boundaries and mappings, but all are consistent with the basic Markov structure, explaining MSMA's robustness.

B.8 Automatic Detection of Hierarchical Boundaries

We provide a practical algorithm to detect semantic hierarchy boundaries, critical for applying MSMA. Algorithm A.1 (Semantic Boundary Detection):

- 1244 1. Input: Model M with L layers, corpus \mathcal{D} .
 - 2. Compute attention span: For each layer l, calculate mean span S_l and difference $\Delta S_l = S_{l+1} S_l$.
 - 3. Compute inter-layer mutual information: For each pair (i, j), compute I_{ij} and construct the matrix I.
 - 4. Functional probing: For each l, evaluate linguistic task accuracy P_l^t and compute performance gradient $\nabla P_l^t = P_{l+1}^t - P_l^t$.

5. Boundary integration: Integrate evidence into a boundary score $B_l = \alpha \Delta S_l + \beta \Delta I_l + 1255$ $\gamma \sum_t w_t \nabla P_l^t$. Identify two peaks as boundaries l_1, l_2 . 1256

1258

1259

1260

1261

1263

1264

1265

1266

1267

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1300

6. *Output:* Boundaries l_1 (local \rightarrow intermediate) and l_2 (intermediate \rightarrow global).

This robustly identifies semantic boundaries for subsequent manifold alignment. Experiments show this ensemble method is more reliable than any single metric.

B.9 Conclusion

This appendix gives a complete mathematical foundation for multi-scale manifold alignment, from information geometry to KL bounds, mutual information preservation, and error decomposition. Our results support the main paper's conclusions and provide new theoretical insights.

Key innovations include: (1) explicit KL connection to geometric/information errors; (2) proof of mutual information retention across mappings; (3) theoretical role of curvature regularization; (4) how hierarchical Markov structure enables error decomposition. These results are consistent with experiments, validating MSMA as a unified and broadly applicable LLM interpretability framework.

C Experimental Setup and Analysis for Multi-Scale Alignment Methods

C.1 MSMA Model Architecture

The Multi-Scale Semantic Alignment (MSMA) framework integrates hierarchical feature extraction with joint optimization, consisting of three main components:

Local Layers: Shallow Transformer blocks capture token-level semantics and syntactic patterns, mainly handling lexical choice, part-of-speech features, and local dependencies, laying the foundation for higher-level semantic abstraction.

Intermediate Layers: These model phraselevel compositionality via mid-depth attention mechanisms, focusing on inter-sentence relations, logical transitions, and local discourse structure, thereby connecting micro-level word features to macro-level topics.

Global Layers: Deep Transformer modules aggregate document-level context, handling topic consistency, discourse structure, and global stylistic coherence, ensuring overall textual fluency.

1304

- 1306
- 1307
- 1308
- 1309
- 1310
- 1311

1312

1313

1314

1315

1316

1317

1320

1321

1322

1325

1326

1327

1329

1330

1331

1332

1333

1334

1335

1336

Each scale produces a semantic vector by mean pooling and layer aggregation, reflecting the empirical disentangling of information observed in intervention studies.

Parallel classifiers operate on hierarchical representations:

Global: $f_{\text{global}} : \mathbb{R}^{\text{hidden_size}} \to \mathbb{R}^{62} \text{ (softmax)}$

Intermediate: $f_{\text{mid}} : \mathbb{R}^{\text{hidden}_{\text{size}}} \to \mathbb{R}^3$ (softmax w/ temperature)

 $f_{\text{local}}: \mathbb{R}^{\text{hidden}_{\text{size}}} \to \mathbb{R}^3$ (label smoothing) Local:

The joint classification loss combines weighted cross-entropy:

$$L_{\rm cls} = \frac{1}{3} \left(H(y_{\rm global}, \hat{y}_{\rm global}) + H(y_{\rm mid}, \hat{y}_{\rm mid}) + H(y_{\rm local}, \overline{\hat{y}_{\rm local}}, \overline{\hat{$$

where H denotes cross-entropy, and y, \hat{y} are ground truth and predictions. This encourages the model to learn effective representations at all semantic scales.

C.2 Semantic Alignment Optimization

Three complementary methods are used in MSMA, 1318 each targeting a different aspect of alignment: 1319

Geometric Alignment. Enforces structural consistency by minimizing the Euclidean distance between representations at different scales. For global-to-intermediate mapping f_{GI} and intermediate-to-local mapping f_{IL} :

$$\mathcal{L}_{\text{geo}} = \|f_{GI}(h_G) - h_I\|^2 + \|f_{IL}(h_I) - h_L\|^2$$

Both linear (least-squares) and nonlinear (MLP) mappings were explored; linear suffices in most cases.

Information Alignment. Maximizes mutual information (MI) between source and mapped representations:

$$\mathcal{L}_{info} = -I(h_G; f_{GI}(h_G)) - I(h_I; f_{IL}(h_I))$$

MI is estimated via MINE:

$$I(X;Y) \approx \mathbb{E}_{p(x,y)}[T_{\theta}(x,y)] - \log \mathbb{E}_{p(x)p(y)}[e^{T_{\theta}(x,y)}]$$

where T_{θ} is a neural network scoring joint vs. marginal samples.

Curvature Regularization. Penalizes high-1337 curvature regions on the representation manifold 1338 for smoother optimization: 1339

1340
$$\mathcal{L}_{curv} = \int_{\mathcal{M}} K^2 dV \approx \sum_i K_i^2 \Delta V_i$$

K is Riemannian curvature; computed via finite differences in practice.

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

The regularization coefficients $\lambda_{geo} = 0.1$, $\lambda_{info} = 0.1, \ \lambda_{curv} = 0.01$ are tuned by grid search. Empirically, geometric alignment is most critical for output quality, so λ_{geo} was varied in $\{0.1, 0.2, \dots, 1.0\}$ for further study.

Default configuration uses Adam (lr=2e-5), batch size 128, 15 epochs, with the multi-scale classifier (output dims: 62/3/3).

Table 8: Ablation Group Configurations

$\hat{y}_{\text{local}}^{\text{Name}}$	Geo	Info	Curv	$\lambda_{ m geo}$	$\lambda_{ ext{info}}$	$\lambda_{ ext{curv}}$
baseline	×	×	×	0	0	0
full_msma	\checkmark	\checkmark	\checkmark	0.1	0.1	0.01
no_geo	×	\checkmark	\checkmark	0	0.1	0.01
no_info	\checkmark	×	\checkmark	0.1	0	0.01
no_curv	\checkmark	\checkmark	×	0.1	0.1	0
only_info	×	\checkmark	×	0	0.1	0
only_curv	×	×	\checkmark	0	0	0.01
only_geo_0.1	\checkmark	×	×	0.1	0	0
only_geo_0.2	\checkmark	×	×	0.2	0	0
	• • •					• • •
only_geo_1	\checkmark	×	×	1.0	0	0

Metrics: KL divergence (lower is better), Mutual Information (MI) (higher is better), Distance **Correlation (D-Corr)** (closer to 1 is better).

C.4 Results and Analysis

Training Loss Analysis. Figures 5 and 6 (not shown here for brevity) compare loss trajectories for each group, confirming: (1) geometric alignment is critical for stability; (2) BERT is more stable overall; (3) curvature regularization is effective early in training; (4) groups with geometry converge faster.

Hyperparameter Sensitivity

Effect of λ_{geo} . On GPT-2, KL is stable for $0.1 \leq$ $\lambda_{\text{geo}} \leq 0.9$ but increases slightly at 1.0. MI peaks at intermediate values. D-Corr remains above 0.999 for all values.

On BERT, KL is minimized at $\lambda_{geo} = 0.3$ or 0.7, while MI follows a U-shape, peaking at 1.0. Default $\lambda_{\text{geo}} = 0.1$ works well for most cases; BERT may benefit from higher weights.

Other Hyperparameters. λ_{info} is stable in [0.05, 0.2], with higher values harming KL. λ_{curv} is optimal in [0.005, 0.02]; too small gives little regularization, too large restricts flexibility. Learning rate 2e-5 is best—higher values destabilize training, lower values slow convergence.

Training and Validation Loss Comparison



Figure 5: Training Loss Curves of Different Experimental Groups for GPT2

Training and Validation Loss Comparison



Figure 6: Training Loss Curves of Different Experimental Groups for BERT

1378	These analyses guide robust MSMA application
1379	across models and tasks.