

# Time Travel Engine: A Shared Latent Chronological Manifold Enables Historical Navigation in Large Language Models

Anonymous ACL submission

## Abstract

Time functions as a fundamental dimension of human cognition, yet the mechanisms by which Large Language Models (LLMs) encode chronological progression remain opaque. We demonstrate that temporal information in their latent space is organized not as discrete clusters but as a continuous, traversable geometry. We introduce the Time Travel Engine (TTE), an interpretability-driven framework that projects diachronic linguistic patterns onto a shared chronological manifold. Unlike surface-level prompting, TTE directly modulates latent representations to induce coherent stylistic, lexical, and conceptual shifts aligned with target eras. By parameterizing diachronic evolution as a continuous manifold within the residual stream, TTE enables fluid navigation through period-specific "zeitgeists" while restricting access to future knowledge. Furthermore, experiments across diverse architectures reveal topological isomorphism between the temporal subspaces of Chinese and English—indicating that distinct languages share a universal geometric logic of historical evolution. These findings bridge historical linguistics with mechanistic interpretability, offering a novel paradigm for controlling temporal reasoning in neural networks. The code and datasets are available at <https://anonymous.4open.science/r/TTE-FF57>.

## 1 Introduction

*"Time is the substance I am made of. Time is a river which sweeps me along, but I am the river; it is a tiger which destroys me, but I am the tiger; it is a fire which consumes me, but I am the fire."*

— Jorge Luis Borges, Labyrinths

Time is a dimension: in Minkowski spacetime, it is a continuous coordinate that governs physical dynamics. Yet time is more than physics. In philosophy, Kant (1908) treats time as an a priori form

of intuition—an inner frame that binds disparate sensations into a single stream of consciousness. Neuroscience echoes this view: hippocampal "time cells" represent the flow of time within specific memories Eichenbaum (2014), and the lateral entorhinal cortex encodes elapsed time in the overall population state across scales from seconds to hours Tsao et al. (2018). Such a manifold supports "mental time travel": memories are retrieved not as snapshots but as trajectories through structured chronology Buzsáki and Llinás (2017).

A paradox arises with large language models (LLMs). Trained on centuries of text, they are often cast as "stochastic parrots"—diachronic data compressed into static, synchronic weights. Unlike the continuous temporal stream in the brain, an LLM first encounters time as scattered tokens in a time-agnostic corpus. Yet mechanistic interpretability suggests these systems may organise historical variation into ordered internal structure. Gurnee and Tegmark (2023) show that LLMs learn linear representations of space and time that are robust across scales. So we ask: is temporal flexibility mere memorisation, or does the model build an implicit "time flow" aligned with the directionality of human narrative? If a latent chronological manifold exists, can we move beyond probing and actively steer the model along the historical axis?

Recent work reveals that LLMs encode time geometrically: Gurnee and Tegmark (2023) map internal activations to uncover space-time embeddings across scales and identify "time neurons" that encode historical coordinates, implying temporal metadata is intrinsic rather than memorized. Papadopoulos et al. (2024) similarly observe a causal asymmetry—the perplexity gap between forward-backward prediction reveals an arrow of time. Yet Dhingra et al. (2022) and Zhang and Choi (2023) note that language models trained on static corpora often cannot determine when facts become outdated, highlighting temporal misalignment. Taken

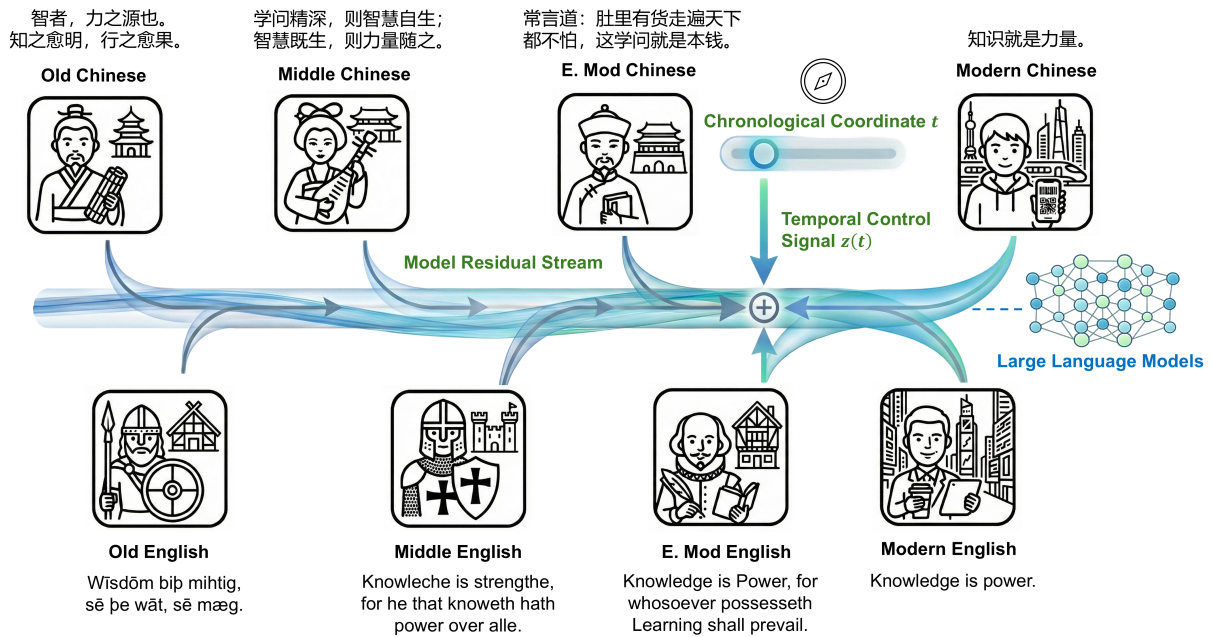


Figure 1: The Time Travel Engine (TTE). This framework navigates a traversable latent chronological manifold within the model. By integrating manifold-derived temporal control signals  $z(t)$  into the residual stream, the mechanism dynamically modulates cognitive states and linguistic registers. This approach facilitates fluid navigation to specific historical eras (e.g., *Old*, *Middle*, *Early Modern* [*E. Mod*]) and highlights a topological isomorphism between Chinese and English temporal subspaces.

084 together, time is structured in these models but has  
 085 been treated as a static property rather than a di-  
 086 mension to be traversed.

087 This raises the question: can we navigate time?  
 088 Research on internal control offers clues. Tang et al.  
 089 (2024) find that a small set of language-specific  
 090 neurons can be toggled to switch output languages.  
 091 At a higher level, Lindsey et al. (2025) use attribu-  
 092 tion graphs to map circuits for reasoning, planning  
 093 and multilinguality, and Zou et al. (2023) together  
 094 with Rimsky et al. (2024) show that steering vec-  
 095 tors can modulate properties such as honesty or  
 096 harmlessness.

097 Despite these successes, the linear representation  
 098 hypothesis may not fully account for the non-linear,  
 099 high-dimensional complexity of historical time.  
 100 Specifically, Tan et al. (2024) characterize static  
 101 steering vectors as brittle, and Huang et al. (2025)  
 102 observe that complex behaviors reside on mani-  
 103 fold geometries where direct vector addition intro-  
 104 duces interference. Consequently, understanding  
 105 how these temporal directions intersect and evolve  
 106 within the latent manifold becomes essential for  
 107 precise, diachronic navigation.

108 We address this gap by drawing parallels be-  
 109 tween multilingualism and historical semantics.

110 Multilingual models align different languages in  
 111 a shared semantic space Conneau et al. (2020), al-  
 112 beit at the risk of cultural erasure Han et al. (2025).  
 113 We hypothesize that historical eras behave simi-  
 114 larly: *Old* and *Modern* English are distinct re-  
 115 gions of a shared semantic manifold. This perspec-  
 116 tive resonates with Hamilton et al. (2016), who  
 117 found statistical laws governing semantic change. e  
 118 thus leverage this topological isomorphism to map  
 119 chronological progression, reconceptualizing eras  
 120 not as discrete points, but as continuous regions on  
 121 a traversable semantic manifold.

122 In this paper, we introduce Time Travel Engine  
 123 (TTE) (Figure 1), a novel framework treating the  
 124 latent space of the model as a traversable temporal  
 125 map. By integrating both static steering vectors and  
 126 continuous trajectories, we effectively modulate the  
 127 temporal signal governing the generation process.  
 128 Our contributions are threefold:

- 129 1. **Geometric Discovery of the Latent Chrono-**  
 130 **logical Manifold.** We identify a control-  
 131 lable chronological manifold within LLMs,  
 132 validating that temporal progression is en-  
 133 coded as a continuous, curvilinear geometry.  
 134 This moves beyond discrete linear approxi-  
 135 mations, establishing that historical evolution

136 follows a traversable, high-dimensional trajec-  
137 tory within the latent space.

138 2. **Dual-Mechanism Temporal Steering.** We  
139 demonstrate that model states can be steered  
140 via dual mechanisms: static vectors for tar-  
141 geted eras and dynamic manifolds for fluid  
142 navigation. This enables precise modulation  
143 of the stylistic and conceptual “zeitgeist” with-  
144 out compromising general reasoning.

145 3. **Cross-Lingual Topological Universality.**  
146 We reveal that the geometric encoding of time  
147 exhibits topological isomorphism across lan-  
148 guages (e.g., Chinese and English). Our find-  
149 ings suggest that the mechanism of chronolog-  
150 ical progression is a language-agnostic prim-  
151 itive, structurally distinct from surface-level  
152 linguistic realizations.

## 153 2 Related Work

154 **Temporal Cognition in Language Models.** Early  
155 temporal processing relied on deterministic rule-  
156 based tagging [Chang and Manning \(2012\)](#), but  
157 LLMs have shifted focus toward internal geometric  
158 world models. While [Dhingra et al. \(2022\)](#) note that  
159 static pretraining hinders dynamic fact updating,  
160 [Gurnee and Tegmark \(2023\)](#) identify “time neu-  
161 rons” encoding linear coordinates, and [Park et al.  
162 \(2025\)](#) pinpoint causal “temporal heads” represent-  
163 ing non-numerical temporal dimensions. [Li et al.  
164 \(2025\)](#) refine this geometric view, demonstrating  
165 that LLM temporal perception follows the Weber-  
166 Fechner law, exhibiting logarithmic compression  
167 relative to a reference point.

168 **Representation Engineering (RepE).** Control  
169 methodologies have evolved from training-time  
170 constraints [Hu et al. \(2017\)](#) to inference-time Rep-  
171 resentation Engineering [Zou et al. \(2023\)](#). Tech-  
172 niques like Function Vectors [Todd et al. \(2023\)](#)  
173 extract task vectors from attention heads, while  
174 Activation Addition [Turner et al. \(2023\)](#); [Rimsky  
175 et al. \(2024\)](#) injects steering directions to modu-  
176 late behavior without weight updates. To address  
177 intervention instability [Tan et al. \(2024\)](#); [Im and  
178 Li \(2025\)](#), recent work employs ensemble vectors  
179 [Siddique et al. \(2025\)](#) and In-Context Vectors [Liu  
180 et al. \(2023\)](#) for robust extraction. Furthermore,  
181 [Singh et al. \(2024\)](#) propose constrained affine steer-  
182 ing to shift generation properties while preserving  
183 semantic integrity.

**Cross-Lingual Isomorphism.** Control univer-  
184 sality relies on cross-lingual isomorphism, orig-  
185 inating from linear embedding mappings [Mikolov  
186 et al. \(2013\)](#) and unsupervised alignment [Con-  
187 neau et al. \(2017\)](#). However, modern analyses  
188 show syntax–semantics disentanglement is hard  
189 [Chen et al. \(2019\)](#), alongside language-specific  
190 shallow processing [Tang et al. \(2024\)](#). [Wendler  
191 et al. \(2024\)](#) find intermediate representations tran-  
192 siently pass through an English-anchored space,  
193 while [Ifergan et al. \(2024\)](#) show that cross-lingual  
194 output consistency does not imply shared represen-  
195 tations. This complicates style transfer; English-  
196 derived interventions may amplify English-centric  
197 knowledge transfer [Lim et al. \(2025\)](#), and neuron-  
198 activation modulation guides style transfer with  
199 content preservation [Kong et al. \(2025\)](#). How-  
200 ever, prior steering is confined to discrete mode-  
201 switching, neglecting continuous temporal geom-  
202 etry, while cross-lingual methods lack structural  
203 validation. We address this by formalizing latent  
204 time as a traversable manifold, enabling precise  
205 navigation isomorphic across languages. 206

## 207 3 Methods

208 To achieve precise temporal intervention and in-  
209 vestigate the internal temporal geometry of LLMs,  
210 we introduce the TTE. The core premise of TTE  
211 is that chronological progression is encoded as a  
212 continuous, navigable trajectory within the high-  
213 dimensional activation space. This section details  
214 the complete methodology: beginning with the cu-  
215 ration of diachronic corpora, we proceed to the  
216 extraction of discrete temporal anchors and their  
217 projection onto a continuous manifold, conclud-  
218 ing with the mechanisms for adaptive modulation,  
219 chronotope disentanglement, and cross-lingual val-  
220 idation.

### 221 3.1 Models

222 We run all experiments on four open-weight  
223 decoder-only LLMs, including Qwen2.5-14B-Base,  
224 Qwen2.5-14B-Instruct, Gemma-2-9B-Instruct and  
225 Llama-3.1-8B-Instruct; key specifications are re-  
226 ported in Appendix A.3 (Table 16).

### 227 3.2 Datasets

228 Our experiments rely on three distinct categories of  
229 datasets, constructed to support vector extraction,  
230 alignment, and multi-dimensional evaluation.

Period	English	Chinese
<i>Old</i>	450–1150	12th c. BC–3rd c. AD
<i>Middle</i>	1150–1500	4th–12th c. AD
<i>Early Modern</i>	1500–1700	13th–Early 20th c.
<i>Modern</i>	1700–Present	Early 20th c.–Present

Table 1: Chronological periodization for English and Chinese used in our TTE framework.

**Diachronic Steering Corpora.** To facilitate extracting temporal features, we compiled a parallel diachronic corpus covering four distinct eras for each language: *Old*, *Middle*, *Early Modern*, and *Modern*, following standard periodizations in English and Chinese historical linguistics Gelderen (2014); Dong (2020) (see Table 1). The corpus contains paired stylistic samples ( $\sim 60\text{K}$  prose and  $\sim 30\text{K}$  verse tokens per language-era pair), curated from historical archives such as York-Toronto-Helsinki Parsed Corpus of Old English prose (YCOE) Taylor (2003), the Penn Parsed Corpora of Historical English (PPCHE) Kroch (2020) and the Chinese Text Project. This data serves as the source of authentic historical features for our ensemble steering strategies.

**Evaluation Datasets.** To rigorously assess the steered models, we constructed three distinct benchmarks: the Epistemic Cutoff Dataset, the Causal Remodeling Dataset, and the Mismatch Entanglement Dataset. These benchmarks are designed to evaluate the chronological epistemic integrity of the model, zeitgeist-consistent reasoning, and the disentanglement between stylistic surface forms and deep cognitive representations, respectively. Comprehensive details and sample prompts for each dataset are provided in Appendix A.6.

### 3.3 Constructing the Chronological Manifold

We hypothesize that chronological progression is encoded as a specific direction—or more accurately, a manifold—within the residual stream of LLMs. We verify and navigate this latent chronological manifold through a hierarchy of four strategies, ranging from discrete anchor sampling to continuous trajectory modeling.

Let  $\mathcal{M}$  be the LLM with  $L$  layers. For a given input  $x$ , let  $\mathbf{h}_l(x) \in \mathbb{R}^d$  denote the hidden state at layer  $l$ . We define the set of target eras as  $\mathcal{T} = \{t_1, t_2, t_3\}$ , representing *Old*, *Middle*, and *Early*

*Modern*, with  $t_0$  denoting the *Modern* era (anchor).

**Method I: Contrastive Activation Addition (CAA).** To isolate the intrinsic representation of a specific era from disparate linguistic features, we employ a self-generated Contrastive Activation Addition (CAA) strategy (Rimsky et al., 2024). We utilize "Era Charters"—system prompts that strictly constrain the temporal persona of the model—paired with a diverse set of immersive tasks, ranging from introspective descriptions to epistemic reasoning. Crucially, these tasks are identical across target eras and the *Modern* anchor. By computing the centroid of activations for the same set of tasks under different temporal constraints, we effectively marginalize out task-specific semantics (e.g., topic, sentence structure) and surface-level covariates. The resulting time vector  $\mathbf{v}_{caa}^{(l,t)}$  is thus defined as the contrastive difference between the expected representation of the target era and the *Modern Stats* baseline:

$$\mathbf{v}_{caa}^{(l,t)} = \mathbb{E}_{x \sim \text{Gen}(\mathcal{T}_{\text{tasks}}|C_t)}[\mathbf{h}_l(x)] - \mathbb{E}_{x \sim \text{Gen}(\mathcal{T}_{\text{tasks}}|C_0)}[\mathbf{h}_l(x)]. \quad (1)$$

where  $\text{Gen}(\mathcal{T}_{\text{tasks}}|C_t)$  denotes responses generated by the model to the standardized task set under the specific "Era Charter" constraint.

**Method II: Ensemble CAA (EnsCAA).** While CAA provides clean, isolated temporal signals, synthetic prompts may lack the textural complexity of authentic literature. To bridge this gap, we introduce an ensemble strategy. We compute a parallel set of authentic centroids, denoted as  $\mathbf{v}_{real}^{(l,t)}$ , by averaging the activations across the real-world Diachronic Steering Corpora (Sec. 3.2). The ensemble vector  $\mathbf{v}_{ens}^{(l,t)}$  is formulated as a convex combination of the synthetic and authentic estimators:

$$\mathbf{v}_{ens}^{(l,t)} = \alpha \cdot \mathbf{v}_{caa}^{(l,t)} + (1 - \alpha) \cdot \mathbf{v}_{real}^{(l,t)} \quad (2)$$

where  $\mathbf{v}_{real}$  represents the mean shift calculated over authentic historical texts, and  $\alpha$  is a hyperparameter governing the mixing ratio.

**Method III: Chronological Manifold Projection (CMP).** While Method I effectively isolates discrete temporal states, it leaves the transitional spaces undefined. To bridge these gaps, we propose Chronological Manifold Projection (CMP). We treat the discrete synthetic centroids from Method I,  $\mathcal{V} = \{\mathbf{v}_{caa}^{(l,t)} | t \in \mathcal{T} \cup \{t_0\}\}$ , as anchor points

defining a high-dimensional trajectory. We first identify a principal temporal subspace  $\mathbf{U}_l \in \mathbb{R}^{d \times k}$  via PCA. Within this latent geometry, we fit a polynomial spline  $\mathcal{S}(\cdot)$  to the projected coordinates, modeling time as a continuous evolution:

$$\mathbf{z}(t) = \mathcal{S}(t; \mathcal{V}), \quad \text{where } t \in \mathbb{R} \quad (3)$$

The continuous time vector for any arbitrary point  $t$  (e.g.,  $t = 1.5$ , representing a transitional period) is reconstructed as:

$$\mathbf{v}_{man}^{(l)}(t) = \mu_l + \mathbf{U}_l \cdot \mathbf{z}(t) \quad (4)$$

This formulation extends the isolated time vectors into a differentiable temporal function, facilitating continuous modulation across the chronological spectrum and granting access to transitional eras between discrete anchors.

**Method IV: Ensemble Manifold Projection (En-sCMP).** To achieve the optimal synthesis of geometry and robustness, we apply the fitting algorithm CMP directly to the set of ensemble centroids  $\{\mathbf{v}_{ens}^{(l,t)}\}$  derived in Method II. This produces a trajectory  $\mathbf{v}_{ens-man}(t)$  that preserves the smoothness of the temporal curve while anchoring it to the statistically robust features of genuine historical text.

### 3.4 Temporal Intervention Mechanism

We modulate the temporal signals of LLMs during inference using adaptive relative strength intervention. Unlike fixed-norm addition, this method scales the intervention based on the inherent activation magnitude of the current context, ensuring stability across different layers. The modified hidden state  $\tilde{\mathbf{h}}_l$  is computed as:

$$\tilde{\mathbf{h}}_l = \mathbf{h}_l + \lambda \cdot \|\mathbf{h}_l\|_2 \cdot \frac{\mathbf{v}^{(l)}}{\|\mathbf{v}^{(l)}\|_2} \quad (5)$$

where  $\lambda$  is the steering strength coefficient. We empirically adopt  $\lambda \in [0.05, 0.15]$  based on extensive evaluation.

### 3.5 Chronotope Disentanglement

A critical question in temporal modeling is whether “time” functions merely as a stylistic filter (e.g., the use of archaic pronouns) or as a cognitive constraint. To probe this distinction, we propose a disentanglement mechanism to separate the cognitive temporal signal from the stylistic surface form.

We first construct a Contrastive Style Dataset containing semantic equivalents across diachronic

registers. Detailed samples from this dataset are provided in Appendix A.5. We define a style subspace  $\mathbf{U}_{style}$  by applying Principal Component Analysis (PCA) to the difference in mean activations between these stylistic pairs. While acknowledging that the interaction between style and semantics involves non-linear complexities, we isolate the top principal components as the dominant linear representation of the linguistic register. The raw time vector  $\mathbf{v}_{time}$  is then decomposed via orthogonal projection:

$$\mathbf{v}_{cog} = \mathbf{v}_{time} - \text{proj}_{\mathbf{U}_{style}}(\mathbf{v}_{time}) \quad (6)$$

where  $\mathbf{v}_{cog}$  denotes the residual *cognitive vector*. This vector represents the components of the temporal state strictly orthogonal to the explicit stylistic direction. By steering the model using  $\mathbf{v}_{cog}$ , we establish a robust lower bound for epistemic persistence, empirically verifying whether temporal reasoning endures even when the primary stylistic cues are mathematically neutralized.

### 3.6 Cross-Lingual Temporal Transfer

To verify the universality of these temporal representations, we first perform direct cross-lingual modulation, where the time vector  $\mathbf{v}_{src}^{(t)}$  derived from the source language (e.g., English) is directly added to the residual stream of target language prompts (e.g., Chinese) to test raw feature transferability. Second, to account for rotational misalignments between language subspaces, we employ Procrustes analysis to compute an orthogonal rotation matrix  $\mathbf{R}$  that maps the source temporal trajectory to the target:

$$\min_{\mathbf{R} \in \mathcal{O}(d)} \sum_t \|\mathbf{v}_{tgt}^{(t)} - \mathbf{R} \cdot \mathbf{v}_{src}^{(t)}\|_F^2 \quad (7)$$

This facilitates manifold-aligned modulation—steering the model on Chinese prompts using transformed time vectors derived purely from English history (and vice versa)—thereby rigorously validating the structural isomorphism of the chronological manifold across languages.

Beyond stylistic shifts, we evaluate the engine by quantifying its impact on epistemic integrity, style-cognition disentanglement, and cross-lingual isomorphism via the following indicators.

**Epistemic Integrity.** To measure if the knowledge boundary of the model successfully recedes to the target era, we analyze the generated entities

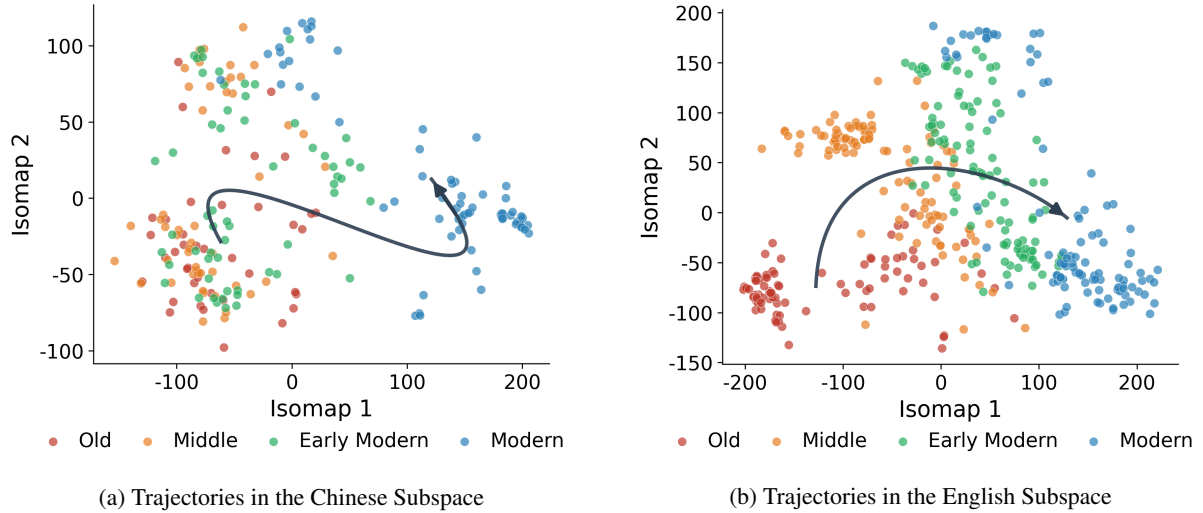


Figure 2: Isomap projection of the latent chronological manifold extracted from Qwen2.5-14B-Base. Each data point represents the activation centroid of a self-generated CAA sample, color-coded by era. The curves illustrate the continuous geometric progression from *Old* to *Modern* periods across both Chinese and English subspaces.

$E_{gen}$ . We classify each entity  $e \in E_{gen}$  into an era  $time(e)$  using an external knowledge base.

- **Future Leakage Rate (FLR):** The fraction of generated entities that belong to the future relative to the target era  $t$ .

$$FLR(t) = \frac{\sum_{e \in E_{gen}} \mathbb{I}(time(e) > t)}{|E_{gen}|} \quad (8)$$

- **Precision Rate (PR):** The fraction of entities that are chronologically accurate (in-scope).

$$PR(t) = \frac{\sum_{e \in E_{gen}} \mathbb{I}(time(e) \leq t)}{|E_{gen}|} \quad (9)$$

## 4 Results

In this section, we present the experimental findings derived from Qwen2.5-14B-Base utilizing the EnsCMP method. We demonstrate the efficacy of the framework through geometric visualization, epistemic boundary verification, chronotope disentanglement, and cross-lingual transfer. We extend our empirical validation in Appendix A.1, providing a comprehensive scaling analysis across diverse architectures (including Qwen2.5-Instruct, Gemma-2, and Llama-3.1). Furthermore, we conduct a comparative evaluation of discrete versus continuous steering paradigms, isolating the specific geometric benefits of the manifold approach.

### 4.1 The Geometry of Time: From Qualitative Modulation to Continuous Manifold

We initiate our analysis by examining the tangible effects of temporal steering on the generation logic

of the model. Tables 2 and 3 present the qualitative results of projecting the hidden states of *Modern* prompts onto historical subspaces.

The results reveal a profound diachronic register modulation. For instance, the *Old Chinese* signal (Table 2) steers standard vernacular generation toward *Classical Chinese* (文言文), characterized by classical vocabulary (e.g., “夫”, “之”) and period-specific syntactic structures. Similarly, the *Early Modern* English signal (Table 3) induces Shakespearean syntax and orthography (e.g., “thou”, “hath”) spontaneously. Crucially, this transition transcends lexical substitution, manifesting as a deeper shift toward the moral and cosmological worldviews of the target era.

Having established the qualitative efficacy of these vectors, we scrutinize the geometric structure of these representations. We collect activation centroids across four canonical eras—*Old*, *Middle*, *Early Modern*, and *Modern*—and project them into a two-dimensional space using Isomap Tenenbaum et al. (2000). As illustrated in Figure 2, the projected embeddings do not form disjoint clusters but rather organize into a smooth, curvilinear trajectory. This intrinsic chronological trajectory is observable in both Chinese and English subspaces. The distinct curvature validates the hypothesis that time is encoded not as a discrete categorical variable, but as a continuous manifold, necessitating the use of the chronological manifold method to navigate the transitional spaces that discrete steering vectors fail to faithfully resolve.

<b>PROMPT</b> 人应该如何度过一生? <i>How should one live a good life?</i>
<b>Baseline (Modern State)</b> 人生是一个复杂而独特的过程，每个人都有自己独特的经历和价值观... <i>Life is a complex and unique journey, where everyone possesses their own distinct experiences and values...</i>
<b>Signal: Old Chinese</b> 有二事为大，一曰孝，二曰敬。《礼记》云：夫孝，德之本也... <i>Two matters are paramount: Filial Piety and Reverence. The Book of Rites states: "Now, filial piety is the root of all virtue..."</i>
<b>Signal: Middle Chinese</b> 惟有三事可为：一曰孝，二曰悌，三曰敬。此三者乃为人之本也... <i>There are but three pursuits worthy of action: Filial Piety, Fraternal Duty, and Reverence. These three constitute the very foundation of being human.</i>
<b>Signal: Early Modern Chinese</b> 人生苦短，如何度日方好？《菜根谭》曰：天欲祸人，必先以微福骄之... <i>Life is bitterly short; how then should one pass one's days? The Vegetable Root Discourse says: "When Heaven intends to bring calamity upon a man, it first puffs him up with trifling blessings..."</i>

Table 2: Qualitative demonstration of diachronic register modulation via TTE. The unmodulated baseline is compared against outputs modulated by *Old*, *Middle*, and *Early Modern Chinese* signals, where text in gray represents English translations.

## 4.2 Diachronic Adaptation Verification

To verify that temporal signals capture intrinsic chronological priors rather than stylistic noise, we conduct a Perplexity (PPL) analysis. We compute time vectors using an 80% split of the Diachronic Steering Corpora and evaluate their effectiveness in minimizing perplexity on the held-out 20% validation split.

Figure 3 presents the PPL heatmaps for Chinese (left) and English (right). The rows represent the steered temporal signal, and the columns represent the test corpus era. A strong diagonal dominance is observed: the model achieves relatively low perplexity when the steered signal matches the era of the text. For example, on the *Old English* test corpus, modulating with the *Old English* vector yields substantially lower PPL compared to vectors from later periods. This validates that the EnsCMP mechanism effectively realigns the predictive distribution with the diachronic patterns of the target historical period.

<b>PROMPT</b> How should one live a good life?
<b>Baseline (Modern State)</b> This is the question that has been asked by philosophers for thousands of years. The answer to..
<b>Signal: Old English</b> Se þe wille gepyldlice lifian sceal healdan godne gewunum, swa swa se gecweneð on his gereorde... <i>He who wishes to live steadfastly must maintain good customs, just as the speaker declared in his speech...</i>
<b>Signal: Middle English</b> In the olde tymes of yore, whan lewys and lyf were grete, men soughte to fynde that which maketh a good lyfe.. <i>In the olden times of yore, when laws and life were great, men sought to find that which maketh for a good life...</i>
<b>Signal: Early Modern English</b> Thou shalt know thyselfe first, then seeke vertue, befiege it with thyne endeauours. To liue well is to liue vertuouoslie. <i>First, know thyself; then seek virtue, and pursue it with all thy might. To live well is to live virtuously.</i>

Table 3: Qualitative demonstration of diachronic register modulation via TTE. The unmodulated baseline is compared against outputs modulated by *Old*, *Middle*, and *Early Modern English* signals.

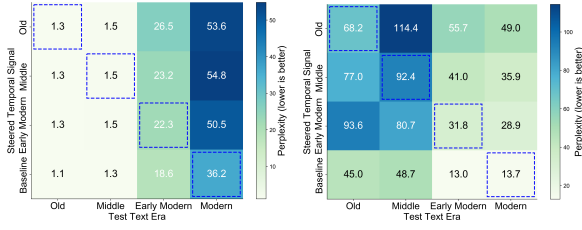
## 4.3 Epistemic Integrity and Knowledge Boundaries

Beyond stylistic adaptation, a robust temporal representation must impose accurate epistemic boundaries. We evaluate this using the Epistemic Cutoff Dataset, measuring the ability of models to retain era-appropriate knowledge (Retention, PR) while suppressing anachronisms (Leakage, FLR).

Figure 4 visualizes the epistemic integrity across target eras. While the unsteered baseline remains anchored in the *Modern State*, TTE induces a distinct trade-off between knowledge retention and anachronism suppression. Across the *Old*, *Middle*, and *Early Modern* trajectories, the engine successfully reconstructs chronological boundaries, maintaining low FLR by filtering future concepts while preserving high PR for era-valid entities. These results indicate that the temporal signal successfully constrains the latent knowledge manifold to historical epistemic boundaries while preserving essential factual integrity.

## 4.4 The Chronotope: Entanglement of Style and Cognition

We evaluate the independence of temporal features using the disentanglement mechanism in Sec. 3.5. Figure 5 presents the entanglement vector field, visualizing the interplay between stylistic precision and future leakage during cognitive isolation.



(a) Chinese PPL Matrix (b) English PPL Matrix

Figure 3: Perplexity matrices for temporal signal modulation. Rows indicate the steered temporal signal; columns indicate the test corpora era.

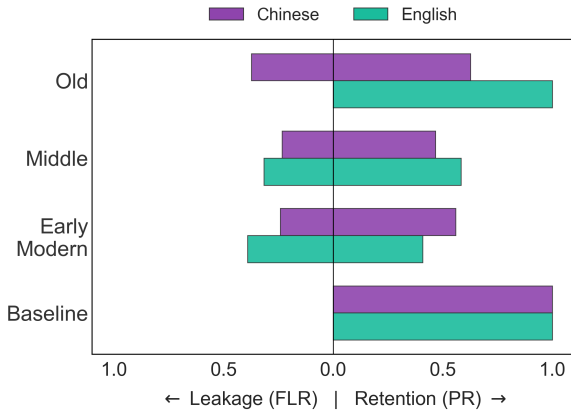
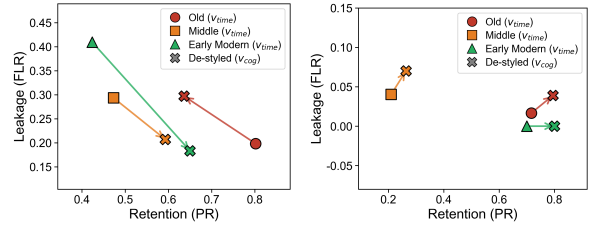


Figure 4: Epistemic detection visualizing knowledge boundary integrity. The unsteered Baseline (*Modern State*) shows minimal leakage.

Results reveal a distinct, era-dependent chronotope entanglement. While aggregate metrics (Appendix A.1.2) suggest that disentanglement improves performance across the timeline, a granular analysis exposes a specific divergence in the *Old* era. Here, removing dominant stylistic components notably increases FLR, suggesting that for deep historical horizons, the representation of “when” is intrinsically bound to “how.” Neutralizing the linguistic register weakens the enforcement of chronological constraints. Conversely, for *Middle* and *Early Modern* eras, the trajectory stabilizes and even improves, driving the overall robustness observed in the EnsCMP mechanism.

#### 4.5 Cross-Lingual Temporal Transfer

Finally, we investigate the universality of these diachronic representations. We postulate that if the geometric encoding of time is universal, the chronological manifold of one language should possess a structural alignment with that of another. To verify this, projection experiments are performed where the internal state of the model is steered us-



(a) Vector Field for Chinese (b) Vector Field for English

Figure 5: Visualization of the entanglement vector fields. The arrows depict the shift in performance metrics when orthogonalizing the time vector to remove stylistic features ( $\mathbf{v}_{time} \rightarrow \mathbf{v}_{cog}$ ).

ing foreign historical signals (e.g., English vectors modulating Chinese prompts). Qualitative illustrations of such cross-lingual register modulation are presented in Table 14 and Table 15.

Analyses in Sec. A.2 and Appendix A.1.3 uncover two governing mechanisms. Qualitatively, we observe representation override: temporal signals suppress prompt-induced priming to steer generation across languages, validating a robust, language-agnostic “chronotope.” Quantitatively, cross-architecture benchmarks show that while discrete alignment is brittle—prone to catastrophic forgetting—continuous manifold projections preserve semantic structure. This topological isomorphism enables zero-shot temporal transfer where discrete matching proves inadequate.

## 5 Conclusion

This work introduces the TTE, a framework demonstrating that chronological time in LLMs functions as a shared, continuous manifold rather than fragmented metadata. By uncovering a topological isomorphism between Chinese and English temporal subspaces, we establish that LLMs organize history through a universal geometric logic transcending linguistic boundaries. Our methodology unifies discrete steering with continuous trajectory modeling, enabling precise historical navigation and zero-shot transfer of temporal cognition. These findings reveal a “chronotope entanglement,” where linguistic registers serve as structural anchors for deep epistemic constraints. Ultimately, our evidence suggests that LLMs construct a coherent internal world model governed by a structured temporal dimension. By bridging historical linguistics with geometric interpretability, we provide a new paradigm for decoding the evolution of human discourse within machine intelligence.

## 6 Limitations

We identify three constraints regarding the scope, precision, and granularity of the TTE framework:

- **Latent Density Dependence.** The efficacy of the engine is strictly bounded by the pre-trained latent territory of the model. While our framework successfully *navigates* existing historical representations, it cannot *synthesize* linguistic competence absent from the training data. Consequently, resource-scarce periods may suffer from extraction noise, yielding vectors that capture superficial lexical markers rather than deep grammatical structures.
- **Chronotope Entanglement.** While we demonstrate a “Shared” manifold, complete disentanglement between stylistic surface forms and deep cognitive constraints remains elusive. Our findings indicate that for architectures of the current generation, the representation of “when” is inextricably bound to “how.” This limits the feasibility of purely semantic temporal steering, as neutralizing the archaic style often degrades the enforcement of epistemic boundaries.
- **Manifold Approximation Granularity.** Our framework constructs a continuous spline trajectory to enable fluid navigation. This geometric modeling inherently prioritizes global coherence over local granularity, serving as a low-dimensional approximation of historical evolution. Consequently, the trajectory may smooth over high-frequency, non-linear discontinuities (such as abrupt linguistic revolutions), capturing the prevailing zeitgeist rather than granular, event-level shifts.

**Potential Risks.** Chronological steering introduces dual risks. **Safety Bypass:** navigating to eras predating modern norms may inadvertently circumvent safety filters tuned on contemporary data. **Historical Essentialism:** compressing complex history into single control signals risks amplifying latent biases, reinforcing stereotypes rather than factual reconstruction.

## References

György Buzsáki and Rodolfo Llinás. 2017. Space and time in the brain. *Science*, 358(6362):482–485.

- Angel X Chang and Christopher D Manning. 2012. SUTime: A library for recognizing and normalizing time expressions. In *Lrec*, volume 12, pages 3735–3740.
- Mingda Chen, Qingming Tang, Sam Wiseman, and Kevin Gimpel. 2019. A multi-task approach for disentangling syntax and semantics in sentence representations. *arXiv preprint arXiv:1904.01173*.
- Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
- Alexis Conneau, Shijie Wu, Haoran Li, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Emerging cross-lingual structure in pretrained language models. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 6022–6034.
- Bhuwan Dhingra, Jeremy R Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W Cohen. 2022. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273.
- Hongyuan Dong. 2020. *A history of the Chinese language*. Routledge.
- Howard Eichenbaum. 2014. Time cells in the hippocampus: a new dimension for mapping memories. *Nature Reviews Neuroscience*, 15(11):732–744.
- Elly Gelderen. 2014. A history of the english language.
- Wes Gurnee and Max Tegmark. 2023. Language models represent space and time. *arXiv preprint arXiv:2310.02207*.
- William L Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic word embeddings reveal statistical laws of semantic change. *arXiv preprint arXiv:1605.09096*.
- HyoJung Han, Sweta Agrawal, and Eleftheria Briakou. 2025. Rethinking cross-lingual alignment: Balancing transfer and cultural erasure in multilingual llms. *arXiv preprint arXiv:2510.26024*.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. 2017. Toward controlled generation of text. In *International conference on machine learning*, pages 1587–1596. PMLR.
- Yao Huang, Huanran Chen, Shouwei Ruan, Yichi Zhang, Xingxing Wei, and Yinpeng Dong. 2025. Mitigating overthinking in large reasoning models via manifold steering. *arXiv preprint arXiv:2505.22411*.
- Maxim Ifergan, Leshem Choshen, Roei Aharoni, Idan Szpektor, and Omri Abend. 2024. Beneath the surface of consistency: Exploring cross-lingual knowledge representation sharing in llms. *arXiv preprint arXiv:2408.10646*.



## A Appendix

### A.1 Performance Scaling across Architectures and Steering Mechanisms

In this appendix, we expand our analysis beyond the specific case study of Qwen2.5-14B-Base model presented in the main text. By benchmarking four distinct architectures (Qwen2.5-14B-Base (Qwen-14B), Qwen2.5-14B-Instruct (Qwen-14B-Inst), Llama-3.1-8B-Instruct (Llama-8B-Inst), and Gemma-2-9B-Instruct (Gemma-9B-Inst) across three interpretability dimensions—World Knowledge, Style-Cognition Disentanglement, and Cross-Lingual Transfer—we uncover the fundamental geometric laws governing temporal representations in LLMs.

In Tables 4–13, all quantitative results are reported as mean  $\pm$  standard deviation, aggregated across three historical eras (*Old*, *Middle*, *Early Modern*). The reported values represent the arithmetic mean of the era-specific scores and standard deviations, respectively, quantifying performance variability across evaluation samples. We benchmark the four steering mechanisms defined in Sec. 3.3: CAA (Contrastive Activation Addition), Ens (Ensemble CAA), CMP (Chronological Manifold Projection), and EnsCMP (Ensemble Manifold Projection).

#### A.1.1 Principle I: The Geometry of Epistemic Boundaries

We first examine whether the geometric smoothing of temporal signals enhances the ability of the model to discern historical facts. Tables 4 and 5 present the Future Leakage Rate (FLR) and Precision Rate (PR).

The data reveals a distinct alignment-geometry trade-off. For the unaligned Qwen-14B base model, the manifold-based strategy (EnsCMP) significantly outperforms the discrete CAA baseline. Specifically, EnsCMP reduces the FLR to  $0.260 \pm 0.319$  compared to  $0.338 \pm 0.357$  for CAA, while simultaneously boosting the PR from 0.362 to 0.607. This suggests that pre-trained latent spaces contain stochastic noise; the manifold fitting process acts as a geometric filter, smoothing out irregularities to reconstruct a coherent temporal progression.

Conversely, instruction-tuned models exhibit a diminished reliance on geometric smoothing. For Qwen-14B-Inst, the discrete ensemble method (Ens) achieves the optimal balance (FLR 0.299,

PR 0.684), indicating that alignment techniques may cluster temporal concepts, rendering complex spline fitting less critical. Furthermore, Llama-8B-Inst and Gemma-9B-Inst show a preference for the raw CAA baseline, with Llama achieving its lowest FLR (0.187) via discrete modulation. This implies that strong instruction tuning linearizes the latent space, potentially making the "Time" feature sufficiently discrete such that manifold projection offers diminishing returns or introduces unnecessary interpolation noise.

#### A.1.2 Principle II: The Robustness of the Chronotope

To verify whether the steering relies on superficial style or deep cognition, we compare the performance of the full time vector ( $v_{time}$ ) against the orthogonalized cognitive vector ( $v_{cog}$ ).

Tables 6–9 reveal the superior cognitive persistence of manifold-based methods. For Qwen-14B, when moving from  $v_{time}$  to  $v_{cog}$  (removing style), the EnsCMP method not only maintains an extremely low FLR ( $0.193 \rightarrow 0.169$ ) but actually improves its PR ( $0.707 \rightarrow 0.747$ ). In contrast, discrete methods like CAA exhibit stagnation or degradation (FLR  $\approx 0.30$ ) when style is stripped. This confirms that the manifold projection method captures a stable "chronotope" that enforces epistemic boundaries even without stylistic cues.

This robustness extends across architectures with varying degrees of entanglement. For Llama-8B-Inst, EnsCMP demonstrates remarkable stability, maintaining a virtually constant PR ( $0.665 \rightarrow 0.668$ ) during disentanglement, whereas discrete methods show greater fluctuation. Conversely, Gemma-9B-Inst exhibits stronger intrinsic entanglement, where removing style generally triggers an increase in future leakage across all methods; however, manifold strategies (CMP) mitigate this degradation significantly better than ensemble approaches (e.g., CMP FLR shifts  $0.221 \rightarrow 0.242$  vs. Ens  $0.211 \rightarrow 0.312$ ). Overall, manifold-based steering consistently offers the most reliable decoupling of historical cognition from surface form across diverse model families.

#### A.1.3 Principle III: Topology is the Universal Translator

We further scrutinize universality by separating direct modulation (transferring raw source language vectors to target prompts without alignment) and manifold-aligned modulation (transferring linearly

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.338 ± 0.357	0.428 ± 0.314	0.324 ± 0.384	<b>0.260 ± 0.319</b>
Qwen-14B-Inst	0.370 ± 0.279	<b>0.299 ± 0.211</b>	0.364 ± 0.246	0.356 ± 0.291
Llama-8B-Inst	<b>0.187 ± 0.237</b>	0.307 ± 0.289	0.316 ± 0.435	0.420 ± 0.386
Gemma-9B-Inst	<b>0.260 ± 0.201</b>	0.331 ± 0.228	0.329 ± 0.246	0.449 ± 0.320

Table 4: World knowledge integrity analysis. Future Leakage Rate (FLR) measures the fraction of generated entities that belong to the future relative to the target era. Lower values indicate better suppression of anachronisms.

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.362 ± 0.367	0.505 ± 0.352	0.543 ± 0.411	<b>0.607 ± 0.354</b>
Qwen-14B-Inst	0.613 ± 0.274	<b>0.684 ± 0.264</b>	0.535 ± 0.296	0.644 ± 0.291
Llama-8B-Inst	<b>0.597 ± 0.352</b>	0.559 ± 0.232	0.567 ± 0.462	0.530 ± 0.371
Gemma-9B-Inst	<b>0.740 ± 0.201</b>	0.652 ± 0.235	0.621 ± 0.217	0.518 ± 0.286

Table 5: World knowledge integrity analysis. Precision Rate (PR) measures the fraction of generated entities that are historically accurate (in-scope) for the target era. Higher values indicate richer recall of era-appropriate entities.

aligned source language vectors via procrustes rotation to target prompts).

The contrast between FLR (Tables 10, 11) and PR (Tables 12, 13) provides the strongest evidence for the topological isomorphism of the manifold.

For Qwen-14B, the discrete CAA method fails catastrophically in the challenging manifold-aligned setting: its FLR spikes to 0.701, and its PR collapses to 0.146, indicating a total functional failure where the model neither suppresses future knowledge nor retains historical facts. In sharp contrast, EnsCMP maintains a robust FLR (0.196) and a healthy PR (0.371), demonstrating that the geometric curvature of the temporal manifold is shared across languages.

Crucially, this topological advantage extends to instruction-tuned models, albeit with different dynamics. For Llama-8B-Inst, while CAA achieves a low FLR in manifold-aligned modulation (0.150), it suffers from a significant degradation in PR (0.320), suggesting that discrete vectors may be suppressing capability indiscriminately. EnsCMP, however, recovers the PR to 0.512 while maintaining effective suppression (FLR 0.204). Similarly, for Qwen-14B-Inst, the Ensemble and EnsCMP methods consistently outperform CAA in PR across both direct and manifold-aligned settings. This pattern confirms that while discrete alignment (CAA) is brittle and prone to "catastrophic forgetting" during cross-lingual transfer, the continuous manifold strategies (EnsCMP) successfully preserve the semantic structure of time,

balancing epistemic boundary enforcement with knowledge retention across diverse architectures.

## A.2 Cross-linguistic Generalization

In this section, we present a detailed qualitative evaluation of the cross-linguistic experiments introduced in the main text. To demonstrate the generalization of our framework across genres, we here provide verse examples (complementing the prose instances in Sec. 4.1) and examine how the projection of hidden states from input prompts onto historical subspaces constructed from disparate source languages affects the stylistic output.

Detailed comparisons are provided in Table 14 and Table 15. To systematically analyze these effects, we define two distinct modulation settings:

- **Native Signal Modulation:** The source of the steering vector shares the linguistic identity of the prompt. For example, this involves the modulation of a *Modern Chinese* prompt using a vector extracted from *Old Chinese* data (denoted as ZH→ZH). This setting establishes a baseline for the capability of the model to perform monolingual diachronic style transfer.
- **Cross-Lingual Signal Modulation:** The source of the steering vector differs from the linguistic identity of the prompt. For instance, this involves the modulation of a *Modern Chinese* prompt using a vector extracted from *Old English* data (denoted as

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.290 ± 0.262	0.390 ± 0.244	0.238 ± 0.273	<b>0.193 ± 0.199</b>
Qwen-14B-Inst	0.335 ± 0.266	<b>0.264 ± 0.137</b>	0.385 ± 0.307	0.384 ± 0.242
Llama-8B-Inst	<b>0.140 ± 0.153</b>	0.285 ± 0.200	0.244 ± 0.216	0.193 ± 0.183
Gemma-9B-Inst	<b>0.186 ± 0.116</b>	0.211 ± 0.172	0.221 ± 0.107	0.276 ± 0.225

Table 6: Disentanglement baseline using the full time vector ( $\mathbf{v}_{time}$ ). FLR evaluates the performance of this raw vector, which retains both stylistic and cognitive components. Lower values indicate better suppression.

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.299 ± 0.323	0.381 ± 0.268	0.208 ± 0.229	<b>0.169 ± 0.191</b>
Qwen-14B-Inst	0.324 ± 0.225	<b>0.279 ± 0.139</b>	0.417 ± 0.290	0.345 ± 0.278
Llama-8B-Inst	<b>0.156 ± 0.135</b>	0.290 ± 0.244	0.321 ± 0.231	0.216 ± 0.202
Gemma-9B-Inst	0.245 ± 0.154	0.312 ± 0.176	<b>0.242 ± 0.090</b>	0.359 ± 0.219

Table 7: Disentanglement efficacy using the cognitive vector ( $\mathbf{v}_{cog}$ ). FLR evaluates performance after stylistic features have been removed to isolate cognitive content. Lower values indicate robust cognitive persistence.

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.676 ± 0.290	0.568 ± 0.288	0.620 ± 0.295	<b>0.707 ± 0.261</b>
Qwen-14B-Inst	0.607 ± 0.261	<b>0.736 ± 0.137</b>	0.473 ± 0.284	0.599 ± 0.258
Llama-8B-Inst	<b>0.760 ± 0.224</b>	0.682 ± 0.172	0.748 ± 0.217	0.665 ± 0.197
Gemma-9B-Inst	<b>0.814 ± 0.116</b>	0.789 ± 0.172	0.771 ± 0.130	0.691 ± 0.214

Table 8: Disentanglement baseline using the full time vector ( $\mathbf{v}_{time}$ ). PR reflects the accurate retrieval of historical facts using this raw vector. Higher values indicate richer recall.

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.659 ± 0.334	0.577 ± 0.305	0.650 ± 0.312	<b>0.747 ± 0.263</b>
Qwen-14B-Inst	0.651 ± 0.230	<b>0.721 ± 0.139</b>	0.558 ± 0.307	0.580 ± 0.312
Llama-8B-Inst	<b>0.744 ± 0.215</b>	0.676 ± 0.199	0.671 ± 0.227	0.668 ± 0.234
Gemma-9B-Inst	<b>0.755 ± 0.154</b>	0.688 ± 0.176	0.758 ± 0.090	0.640 ± 0.219

Table 9: Disentanglement efficacy using the cognitive vector ( $\mathbf{v}_{cog}$ ). PR measures the retention of historical facts after stylistic features have been stripped. Higher values indicate minimal knowledge loss.

943 EN→ZH). As detailed in Sec. 3.6, we evaluate  
944 this via direct modulation (raw vector  
945 integration) and manifold-aligned modulation  
946 (Procrustes-based alignment), rigorously test-  
947 ing the robustness of the time vector and the  
948 disentanglement between temporal features  
949 and linguistic surface forms.

950 As evidenced in Table 14, we observe a com-  
951 pelling phenomenon of representation override: the  
952 integration of the *Old English* steering vector into  
953 the processing of *Modern Chinese* prompts results  
954 in outputs generated in the *Old English* register.  
955 Here, the steering magnitude of the time vector ef-

956 fectively suppresses the linguistic priming provided  
957 by the Chinese prompt, shifting the generation tra-  
958 jectory entirely across language boundaries.

959 Conversely, as detailed in Table 15, modulating  
960 *Modern English* contexts with *Old Chinese* signals  
961 induces outputs in *Classical Chinese*. These results  
962 suggest that the extracted time vectors encode a  
963 robust “chronotope”—a unified representation of  
964 time and language—that is sufficiently potent to  
965 dictate the generation path, demonstrating that the  
966 temporal geometry within the model acts as an  
967 active, traversable dimension capable of zero-shot  
968 stylistic transfer.

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.450 ± 0.360	0.459 ± 0.389	0.219 ± 0.266	<b>0.168 ± 0.226</b>
Qwen-14B-Inst	0.379 ± 0.359	<b>0.213 ± 0.259</b>	0.280 ± 0.286	0.304 ± 0.297
Llama-8B-Inst	<b>0.111 ± 0.204</b>	0.198 ± 0.288	0.139 ± 0.257	0.124 ± 0.186
Gemma-9B-Inst	0.411 ± 0.384	<b>0.381 ± 0.399</b>	0.389 ± 0.381	0.454 ± 0.432

Table 10: Cross-lingual transfer using direct modulation. FLR is evaluated when transferring raw source language vectors directly to target prompts (e.g., English Vector → English Prompt).

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.701 ± 0.357	0.473 ± 0.364	0.442 ± 0.310	<b>0.196 ± 0.213</b>
Qwen-14B-Inst	0.386 ± 0.378	<b>0.240 ± 0.271</b>	0.321 ± 0.283	0.312 ± 0.261
Llama-8B-Inst	<b>0.150 ± 0.280</b>	0.157 ± 0.288	0.341 ± 0.405	0.204 ± 0.258
Gemma-9B-Inst	<b>0.474 ± 0.401</b>	0.610 ± 0.445	0.536 ± 0.433	0.542 ± 0.412

Table 11: Cross-lingual transfer using manifold-aligned modulation. FLR is evaluated when transferring vectors across languages via Procrustes alignment (e.g., English Vector → Chinese Prompt).

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.310 ± 0.338	0.354 ± 0.353	0.174 ± 0.248	<b>0.275 ± 0.271</b>
Qwen-14B-Inst	0.271 ± 0.365	<b>0.527 ± 0.359</b>	0.327 ± 0.381	0.386 ± 0.295
Llama-8B-Inst	0.492 ± 0.413	0.468 ± 0.405	<b>0.657 ± 0.403</b>	0.583 ± 0.382
Gemma-9B-Inst	<b>0.396 ± 0.381</b>	0.349 ± 0.380	0.286 ± 0.365	0.227 ± 0.331

Table 12: Cross-lingual transfer using direct modulation. PR is evaluated when using native language vectors without geometric rotation.

Model	CAA	Ens	CMP	EnsCMP
Qwen-14B	0.146 ± 0.228	0.294 ± 0.325	0.189 ± 0.237	<b>0.371 ± 0.251</b>
Qwen-14B-Inst	0.348 ± 0.387	<b>0.510 ± 0.345</b>	0.356 ± 0.363	0.431 ± 0.329
Llama-8B-Inst	0.320 ± 0.400	0.446 ± 0.411	0.366 ± 0.370	<b>0.512 ± 0.407</b>
Gemma-9B-Inst	<b>0.289 ± 0.336</b>	0.130 ± 0.270	0.214 ± 0.330	0.178 ± 0.304

Table 13: Cross-lingual transfer using manifold-aligned modulation. PR is evaluated when using foreign language vectors aligned via Procrustes analysis, testing the universality of the manifold.

### 969 A.3 Model details

970 This appendix reports the architectural specifica- 982  
971 tions of the four LLMs used in our experiments (Ta- 983  
972 ble 16). We select the Qwen2.5-14B Base/Instruct 984  
973 pair specifically to isolate the impact of post- 985  
974 training alignment on temporal geometry. Further- 986  
975 more, we extend our evaluation to Llama-3.1-8B- 987  
976 Instruct and Gemma-2-9B-Instruct to verify the 988  
977 cross-architectural universality of the chronologi-  
978 cal manifold beyond a single model family.

### 979 A.4 Prompts for Temporal State Extraction

980 This appendix documents the immersive tasks em-  
981 ployed to extract temporal states via CAA. To en-

sure the cross-lingual consistency of the extracted 982  
signals, these tasks are designed with semantic par- 983  
allelism across Chinese and English, as detailed 984  
in Table 17 and Table 18. The complete "Era 985  
Charters" used to ground these tasks are available 986  
in our repository at [https://anonymous.4open. 987  
science/r/TTE-FF57](https://anonymous.4open.science/r/TTE-FF57). 988

### 989 A.5 Contrastive Style Dataset

To strictly distinguish between linguistic style and 990  
temporal cognition, we curated a contrastive style 991  
dataset. This dataset consists of parallel seman- 992  
tic pairs where the core meaning remains invariant 993  
while the linguistic register shifts across four canon- 994

ical eras: *Old*, *Middle*, *Early Modern*, and *Modern*.

These pairs serve as the basis for constructing the style subspace ( $U_{style}$ ) used in the disentanglement experiments (Sec. 3.5). We isolate the geometric direction of 'style' by applying PCA to the difference vectors between modern and archaic representations.

Table 19 presents representative examples from the dataset, spanning diverse topics such as time, nature, and social customs. The full dataset, containing paired prose and verse for both Chinese and English, is available at <https://anonymous.4open.science/r/TTE-FF57>.

## A.6 Evaluation Datasets Details

To comprehensively evaluate the effectiveness of the TTE framework, we constructed three specialized datasets, each targeting a specific dimension of temporal alignment: epistemic boundaries, causal reasoning, and feature disentanglement. The following sections provide representative samples and descriptions for each benchmark; the complete datasets and the associated generation scripts are publicly available in our repository at <https://anonymous.4open.science/r/TTE-FF57>.

### A.6.1 Epistemic Cutoff Dataset

This dataset is designed to probe the chronological knowledge boundaries of the model. It consists of open-ended queries requesting lists of entities (e.g., philosophers, inventions, tools). The ground truth for these queries is strictly segmented by historical era. For instance, when the model is steered to the *Old* era (e.g., Pre-Qin), the generation of an entity from the *Modern* era (e.g., "iPhone" or "Kant") is penalized as a future leakage (measured by FLR), while entities belonging to or preceding the target era are counted towards the PR.

### A.6.2 Causal Remodeling Dataset

The Causal Remodeling Dataset evaluates whether the internal reasoning logic of the model shifts to align with the technological and causal constraints of the target era. Distinct from the Epistemic dataset which probes static factual retrieval, this dataset presents practical, daily-life scenarios (e.g., lighting, communication, preservation) that demand era-appropriate problem-solving. We apply the same FLR and PR metrics to this context: FLR measures the intrusion of anachronistic technologies (e.g., suggesting a "flashlight" in the Old era), while PR quantifies the successful adoption

of historically valid causal mechanisms (e.g., "oil lamps").

### A.6.3 Mismatch Entanglement Dataset

This dataset serves as a rigorous stress test to evaluate the disentanglement between linguistic surface realization and latent cognitive representations. The prompts deliberately introduce chronological paradoxes by querying contemporary scientific or technological concepts while the model is modulated by steering vectors derived from *Old*, *Middle*, or *Early Modern* periods. Such a configuration enables an empirical verification of whether the TTE manifold can achieve independent control over the epistemic grounding (the internal knowledge state) and the stylistic register (the linguistic output mode). By isolating these dimensions, we assess the robustness of the temporal steering mechanism against the explicit modern semantic constraints embedded within the prompts.

<b>PROMPT</b>		
<p>请以诗歌的形式描述一场冬日的暴风雪。</p> <p><i>Please describe a winter storm in the form of a poem.</i></p>		
<b>Baseline (Modern State)</b>	<b>Native Signal (ZH→ZH)</b>	<b>Cross-Lingual Signal (EN→ZH)</b>
<b>Signal: Old</b>		
<p>冬夜寒风起，雪花漫天飞。 银装素裹中，世界静悄悄...</p> <p><i>Winter winds rise in the night, snowflakes fly across the sky. Amidst the silver-clad world, all is silent...</i></p>	<p>朔风怒号兮，天地昏黄，素羽纷飞兮，六出其芒...</p> <p><i>The north wind howls—lo! Heaven and earth turn dusky yellow. White plumes scatter—lo! Six-pointed are their blades...</i></p>	<p>þurfe geong wundorlic þing, wolcenum geweorþeð on wolde, hrerende hræfena gehyngum hrimbealo sceawade.</p> <p><i>A wondrous thing comes to pass in the clouds upon the wold; stirring the ravens' hunger, the frost-ruin appears.</i></p>
<b>Signal: Middle</b>		
<p>冬夜寒风起，雪花漫天飞。 银装素裹中，世界静悄悄...</p> <p><i>Winter winds rise in the night...</i></p>	<p>朔风卷地雪漫天，万木凋零玉絮旋。冻合江河冰作岸，云迷岩谷霰连绵...</p> <p><i>The north wind rolls across the earth, snow fills the sky; Ten thousand trees wither, jade fluffs swirl on high. Ice locks the rivers, forming banks of frost; Clouds confuse the valleys, in endless sleet lost.</i></p>	<p>The wynter storme is so wode, The snow floures ful thikke falleth, As hevy stones down they shalle...</p> <p><i>The winter storm is so furious; The snow flowers fall full thick, As heavy stones down they shall [fall]...</i></p>
<b>Signal: Early Modern</b>		
<p>冬夜寒风起，雪花漫天飞。 银装素裹中，世界静悄悄...</p> <p><i>Winter winds rise in the night...</i></p>	<p>纷纷扬扬漫天，素裹银装玉栏，风卷残云舞翩跹...</p> <p><i>Fluttering and scattering across the vast sky, Clad in white, silver-wrapped upon jade balustrades; The wind coils the scattered clouds, dancing in a whirl...</i></p>	<p>Thou art beholding unto Winter's fierce tempest, As snow doth swiftly swell from heaven above...</p> <p><i>You are beholding Winter's fierce tempest, As snow does swiftly swell from heaven above...</i></p>

Table 14: Qualitative demonstration of diachronic register modulation via TTE. Baseline is compared against native-signal and cross-lingual-signal modulation. Text in gray represents English translations.

<b>PROMPT</b>		
Describe a winter storm in verse.		
<b>Baseline (Modern State)</b>	<b>Native Signal (EN→EN)</b>	<b>Cross-Lingual Signal (ZH→EN)</b>
<b>Signal: Old</b>		
The sky was dark and heavy, A blanket of grey, so dreary...	三月春寒，雪落纷纷，北风呼啸，天地昏昏... <i>Spring chill in March, snow falls in profusion; The north wind whistles, heaven and earth are dim...</i>	惟天有雪，其寒若霜。风兮何起，云兮何茫... <i>Heaven alone holds the snow; its cold is like frost. Oh, from where does the wind rise? Oh, how vast are the clouds...</i>
<b>Signal: Middle</b>		
The sky was dark and heavy, A blanket of grey, so dreary...	þe wynd doth blowe wyth fyer and might, þe snow falleþ downe styffe and bryght... <i>The wind blows with fire and might, The snow falls down stiff and bright...</i>	惟馮雪皓然，白者皆素焉。寒月霜气冽，清风松韵传... <i>Only the snow remains bright; all that is white is pure. The cold moon's frost-air is biting; the clear wind conveys the pine's rhyme...</i>
<b>Signal: Early Modern</b>		
The sky was dark and heavy, A blanket of grey, so dreary...	A ftrong wifterne winde doth blow, And with his fstrength he fendes forth fnow... <i>A strong western wind does blow, And with his strength he sends forth snow...</i>	南枝梅蕊，北窗竹影，西阁茶烟。雪风儿六花作团，云月儿一玉相联... <i>Plum blossoms on the south branch, bamboo shadows by the north window, tea smoke in the west pavilion. The snowy wind swirls six-petaled flowers into balls; clouds and moon join in jade-like unity...</i>

Table 15: Qualitative demonstration of diachronic register modulation via TTE. Baseline is compared against native-signal and cross-lingual-signal modulation.

<b>Model</b>	<b>Params</b>	<b>Layers</b>	<b>Attention Heads</b>	<b>Hidden States</b>	<b>Feedforward</b>
Gemma-2-9B-Instruct	9B	42	16	3584	14336
Llama-3.1-8B-Instruct	8B	32	32	4096	14336
Qwen2.5-14B-Instruct	14B	48	40	5120	13824
Qwen2.5-14B-Base	14B	48	40	5120	13824

Table 16: Model specifications.

#	Task
1	用三句话介绍你自己、你的身份和你每日所见的制度。
2	描写春日清晨城郭或乡野的景色。
3	解释什么是‘薄如纸而能映人像与远景，可收于衣囊’的怪物。
4	说明如何度量时间、距离、重量，常用何种单位？
5	城中多人发热咳嗽而亡，你如何判断病因并建议官府采取措施？
6	连续歉收或大水成灾，你认为成因为何？如何赈济安民？
7	描述一次远行的准备工作和沿途可能遇到的困难。
8	解释天象变化（日食、彗星）的含义和应对方法。
9	介绍当世最重要的三部典籍及其核心思想。
10	描述一场婚礼或丧葬仪式的基本流程。
11	说明如何治疗常见的风寒发热之症。
12	解释为何某些行为会被视为违礼或犯法。
13	描述一位德高望重之人应具备的品质。
14	说明农事的基本节气安排和耕作要点。
15	描写月夜独坐或登高远眺的感受（诗意表达）。

Table 17: Chinese immersive tasks.

#	Task
1	Introduce yourself in three sentences: your identity and the institutions you encounter daily.
2	Describe a spring morning scene near a town or countryside.
3	Explain what a 'paper-thin object that reflects faces and distant scenes, fitting in a pocket' might be.
4	Explain how you measure time, distance, and weight. What units are commonly used?
5	A fever and coughing illness kills many in town. How do you determine the cause and advise authorities?
6	Continuous crop failure or flooding occurs. What causes do you suspect and how to provide relief?
7	Describe preparations for a long journey and difficulties you might encounter.
8	Explain the meaning of celestial changes (eclipses, comets) and how to respond.
9	Introduce the three most important texts of your era and their core ideas.
10	Describe the basic procedures of a wedding or funeral ceremony.
11	Explain how to treat a common cold or fever.
12	Explain why certain behaviors are considered improper or illegal.
13	Describe the qualities a person of great virtue should possess.
14	Explain the basic agricultural calendar and farming essentials.
15	Describe feelings while sitting alone on a moonlit night or gazing from a height (poetic expression).

Table 18: English immersive tasks.

Period	Chinese Text	English Text
<i>Topic: Passage of Time</i>		
<i>Old</i>	逝者如斯夫，不舍昼夜。	Tīd ne bīdeþ nānigne mann, hēo æfre forðgæþ.
<i>Middle</i>	光阴似箭，日月如梭，岁不我与，时不再来。	Tyme abideth noon, it passeth as a stream.
<i>E. Mod</i>	俗话说得好，一寸光阴一寸金，寸金难买寸光阴哪。	Time and Tide wait for no Man, but ever flow onward.
<i>Modern</i>	时间不等人。	Time waits for no one.
<i>Topic: Blooming Flowers</i>		
<i>Old</i>	百花竞放，其香四溢。	Blostman springaþ, hira wlite is fæger.
<i>Middle</i>	繁花似锦，争奇斗艳，满园春色关不住。	The floures bloomen faire, of dyverse colours and swete smell.
<i>E. Mod</i>	花儿都开了，红的白的紫的，好不热，香气扑鼻。	The Flowers are blooming in great abundance, their sweet perfume fills the air.
<i>Modern</i>	花开了。	The flowers are blooming.
<i>Topic: Autumn Scenery (Verse)</i>		
<i>Old</i>	蒹葭苍苍，白露为霜。秋风萧瑟，草木摇落。	Hærfest is cumen, hēah wind blæwþ; lēaf feallaþ, lond wearþ cald.
<i>Middle</i>	碧云天外雁南飞，黄叶纷纷满地堆。萧瑟西风吹不尽，一年秋色到楼台。	The leues fallen from the tre, The wynde doth blowe so colde and kene; The somer dayes now fledde be, And al the worlde is bare and lene.
<i>E. Mod</i>	西风吹落叶，黄叶满庭阶。一年好景处，秋色最堪怀。	When Autumn comes with golden Hue, The Leaves do fall from every Tree; The Summer Days have bid Adieu, And Nature mourns what used to be.
<i>Modern</i>	秋叶飘落，铺满了小径。风吹过树梢，带走了夏天。	Autumn leaves are falling, carpeting the ground in gold. The wind whispers through bare branches, a farewell to warmth.
<i>Topic: Wine and Gathering (Verse)</i>		
<i>Old</i>	我有嘉宾，鼓瑟吹笙。呦呦鹿鸣，食野之苹。	Wīn bringþ glædnesse, gāstas āhebbap; drync and drēam, dēorwyrðe sǣl.
<i>Middle</i>	葡萄美酒夜光杯，欲饮琵琶马上催。醉卧沙场君莫笑，古来征战几人回。	Wyn bryngeth joye and maketh glad, It lyghteth up the derke of nyght; It comforteth the herte that's sad, And turneth wronge to alle aryght.
<i>E. Mod</i>	劝君更尽一杯酒，西出阳关无故人。今朝有酒今朝醉，明日愁来明日愁。	Wine brings Joy to ev'ry Heart, It lifts the Soul and cheers the Mind; From Sorrow's Grip it sets apart, And leaves our Worldly Cares behind.
<i>Modern</i>	一杯美酒，带来欢乐。朋友相聚，笑语盈盈。	A cup of wine brings joy. Friends gather, laughter fills the air.

Table 19: Representative samples from the contrastive style dataset used for computing the style subspace. The dataset aligns semantic content across four historical periods.

English Query	Chinese Query
Who are famous philosophers?	著名的哲学家有哪些?
What were major inventions?	重大发明有哪些?
What scripts were primarily used?	主要书写系统有哪些?
How was time kept?	记录时间的主要工具是什么?
How did people communicate over long distances?	主要远程通讯方式是什么?

Table 20: Sample queries from the Epistemic Cutoff Dataset. The model is expected to filter its output based on the temporal knowledge cutoff of the steered era.

English Query	Chinese Query
I need to send a message from Beijing to London. How long will it take?	我需要从北京给伦敦发送信息，需要多长时间?
The night is dark. What do I do?	夜晚无灯，我该如何照亮?
I want to preserve fish for the winter. What should I do?	我想保存鱼肉过冬，该怎么做?
In winter it's cold. How do I keep the house warm?	冬天寒冷，我如何使家中温暖?
How can I cross a river?	要过河，应使用什么方式?

Table 21: Sample queries from the Causal Remodeling Dataset. These prompts test whether the model adopts the technological and causal constraints of the target era.

English Query	Chinese Query
Please answer using modern English prose, what are the main methods for treating diseases today?	请用现代汉语白话文体回答，当今时代治疗疾病的主要方法是什么?
Please answer using modern English prose, what media are currently used to store information?	请用现代汉语白话文体回答，目前常用的信息存储介质是什么?
Please answer using modern English prose, what sports are popular currently?	请用现代汉语白话文体回答，目前流行的运动是什么?
Please answer using modern English verse, what water supply systems exist in the modern era?	请用现代汉语诗歌文体回答，现代供水系统有哪些?
Please answer using modern English verse, what painting materials were used in the present day?	请用现代汉语诗歌文体回答，现今绘画材料有哪些?

Table 22: Sample queries from the Mismatch Entanglement Dataset. These queries create deliberate conflicts between the requested content (*Modern*) and the specified style or the active time vector.