
Focus and Dilution: The Multi-stage Learning Process of Attention

Anonymous Authors¹

Abstract

Transformer-based models have achieved remarkable success across a wide range of domains, yet our understanding of their training dynamics remains limited. In this work, we identify a recurrent focus–dilution cycle in attention learning and provide a rigorous explanation in a one-layer Transformer setting for Markovian data via gradient-flow analysis. Using stage-wise linearization around critical points, we show that a single focus–dilution cycle can be decomposed into a sequence of distinct stages. First, embedding and projection rapidly condense to a rank-one structure, while attention parameters remain effectively frozen. Then, the attention parameters begin to increase, inducing a frequency-driven focus toward high-frequency tokens. As attention continues to evolve, it generates next-order perturbations in embeddings, leading to a mass-redistribution mechanism that progressively dilutes this focus. Finally, small asymmetries among low-frequency tokens lift a degenerate critical point, opening new embedding directions and initiating the next cycle. Experiments on synthetic Markovian data as well as WikiText and TinyStories corroborate the predicted stages and cyclical dynamics.

1. Introduction

Transformer models (Vaswani et al., 2017) have become the dominant architecture for sequence modeling. While their approximation power is now well understood in a variety of regimes (Pérez et al., 2019; Yun et al., 2020a;b), we still lack a mechanistic theory for how attention itself evolves during training. Most existing analyses gain tractability by introducing additional technical condition, such as reparameterizations (Zhang et al., 2024a) or proxy dynamics

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

(Tarzanagh et al., 2023), which may obscure the native coupling among embeddings, projection, and attention. Moreover, recent work suggests that Transformer training often undergoes multiple stages (Chang et al., 2024), and that attention can shift from highly concentrated to more diffuse patterns (Tian et al., 2024). These observations point to a need for a dynamical picture that remains faithful to the coupled dynamics and can explain both attention amplification and its subsequent dissipation.

In this work, we combine theory and experiments to show that attention can be understood as a cyclical learning process. Within each cycle, attention first amplifies a frequency-driven preference over tokens (*focus*), and then gradually weakens this preference as the embedding structure adapts (*dilution*). We identify the dynamical origin of each stage and explain how the interaction between embeddings and attention progressively decomposes the learning problem.

To keep the analysis tractable while preserving essential sequential structure, we study population gradient flow for a one-layer Transformer trained by cross-entropy on Markov data (Chang et al., 2024; Makkuva et al., 2024; 2025). Our explanation is stage-wise and built on linearizations around critical points. Under small initialization, the trajectory is first governed by the linearization near the origin, which forces a rank-one condensation of the embedding and projection components, consistent with the condensation phenomenon in Chen & Luo (2025). We further observed that the condensed direction is explicitly determined by the stationary distribution. In contrast, the attention parameters remain small in this initial stage because the leading-order driving term for (W_Q, W_K) vanishes at the origin.

After condensation, the trajectory follows the same low-rank ray until it reaches a second critical point. We show that this point is generically a saddle: the Jacobian admits a local block decomposition into (i) a contracting embedding/output subsystem and (ii) an attention subsystem with a single unstable mode. Consequently, once the trajectory enters this neighborhood, (W_Q, W_K) align with the unstable eigendirection and grow exponentially, and attention acquires a bias toward high-frequency tokens, initiating the focus phase.

Going beyond the focus phase requires a more refined description than the local saddle analysis. Once (W_Q, W_K)

align with the unstable direction, dynamics enters to a rank-one invariant manifold and induces a closed reduced system. This reduced flow exposes a mass-redistribution mechanism in the embeddings. As the attention amplitude evolves, it generates next-order perturbations, causing the embeddings of the main token and the remaining tokens to move in opposite directions. As a result, the earlier high-frequency focus is gradually weakened, leading to an attention dilution phase.

Finally, the model must learn new embedding directions that distinguish low-frequency tokens. However, we show that the training dynamics become trapped at a degenerate critical point on the rank-one manifold, where the driving forces vanish and no new directions can emerge. To model realistic asymmetries and eliminate degeneration, we introduce a small symmetry-breaking perturbation among low-frequency tokens and analyze the resulting bifurcation of critical points. This mechanism explains how new embedding directions are unlocked, thereby initiating the next focus–dilution cycle. Experiments on synthetic Markovian data as well as WikiText and TinyStories corroborate the predicted stages and cyclical dynamics.

Our contributions.

1. We identify a focus–dilution cycle in the training dynamics of attention and introduce a minimal tractable setting that captures this phenomenon.
2. We develop a stage-wise analysis based on linearization at critical points that explains the different stages within single cycle.
3. We empirically validate the predicted stages and transitions, demonstrating that the focus–dilution cycle persists on synthetic Markov data as well as realistic data.

2. Preliminaries

2.1. Basic Notations.

For any $N \in \mathbb{N}$, let $[N] := \{1, \dots, N\}$. Let $\mathcal{V} := [d]$ be the vocabulary set with $d \geq 2$. We identify tokens with indices in $[d]$ and write $\{e_i\}_{i=1}^d$ for the canonical basis of \mathbb{R}^d . For $\alpha \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$, $\|\alpha\|_2$ and $\|A\|_F$ denote Euclidean norm and Frobenius norm separately, with the subscript omitted when clear from context. We write $\|\alpha\|_C := \sqrt{\alpha^\top C \alpha}$ for the seminorm induced by positive semidefinite matrix $C \succeq 0$. For $\alpha \in \mathbb{R}^n$, define the variance matrix $\text{Var}(\alpha) := \text{diag}(\alpha) - \alpha\alpha^\top$.

2.2. Markov data generation

We generate the dataset $\mathcal{D} = \{(X_i, y_i)\}_{i=1}^N$ with $X_i = (x_{i,1}, \dots, x_{i,s}) \in \mathcal{V}^s$ and $y_i := x_{i,s+1}$, by a Markov chain.

Definition 2.1 (Markovian data). Let $P \in \mathbb{R}^{d \times d}$ be row-stochastic. For each $i \in [N]$, sample $x_{i,1} \sim \text{Unif}(\mathcal{V})$ and $x_{i,j} \sim P_{x_{i,j-1}}$ for $j = 2, \dots, s+1$. Set $X_i = (x_{i,1}, \dots, x_{i,s})$ and $y_i = x_{i,s+1}$.

To model one high-frequency token together with a group of low-frequency tokens that may exhibit mild heterogeneity, we consider a stationary distribution of the form $\pi^\top(\delta) := (\pi_1, \dots, \pi_d)$ with

$$\pi_i = \frac{1 - \pi_1}{d - 1} + c_i \delta, \quad \forall 2 \leq i \leq d, \quad (1)$$

where $\frac{1 - \pi_1}{d - 1} < \pi_1 < 1$, $\sum_{i=2}^d c_i = 0$, and $\delta \geq 0$ is a small parameter chosen so that $\pi(\delta)$ remains entrywise nonnegative. The first term describes two-group setting, one high-frequency token and the rest symmetrical low-frequency tokens. The second term is an $O(\delta)$ perturbation that breaks symmetry within the low-frequency group. Unless stated otherwise, we treat the first term as the leading-order component and regard the second term as a small perturbation that can be neglected in early-stage analyses.

The transition matrix is defined as

$$P = \lambda I + (1 - \lambda)\mathbf{1}\pi^\top, \quad 0 < \lambda < 1, \quad (2)$$

where $\mathbf{1} \in \mathbb{R}^d$ is the all-ones vector. A direct computation verifies that $\pi^\top P = \pi^\top$, hence π is stationary for P .

2.3. One-layer transformer

Since the next token depends only on the current token under the Markov assumption, a single attention block is sufficient to capture the relevant dependency. We therefore study a one-layer Transformer and its training dynamics.

Definition 2.2 (One-layer transformer). Given input sequence $X = (x_1, \dots, x_s)$, let $E_X = (e_{x_1}, \dots, e_{x_s})^\top \in \mathbb{R}^{s \times d}$. Let $W_0 \in \mathbb{R}^{d \times m}$ be the embedding matrix and define the embedded sequence $E_X W_0 \in \mathbb{R}^{s \times m}$. For any $Z \in \mathbb{R}^{s \times m}$, the attention block is

$$\text{Attn}(Z) = \text{softmax}\left(ZW_QW_K^\top Z^\top\right)Z.$$

Let $W_1 \in \mathbb{R}^{m \times d}$ be the output projection. The output logits are

$$f_\theta(X) = \text{Attn}(E_X W_0)W_1$$

For notational convenience, we define $W_{QK} := W_QW_K^\top$, $\Phi := W_0W_QW_K^\top W_0^\top$ and $M := W_0W_1$.

2.4. Training objective and gradient-flow dynamics

Given $(X_i, y_i) \in \mathcal{D}$, define the cross-entropy at the last token $\ell(f_\theta(X_i)_s, y_i) := -\log \frac{\exp(f_\theta(X_i)_s, y_i)}{\sum_{j=1}^d \exp(f_\theta(X_i)_s, j)}$. Then

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f_\theta(X_i)_s, y_i). \quad (3)$$

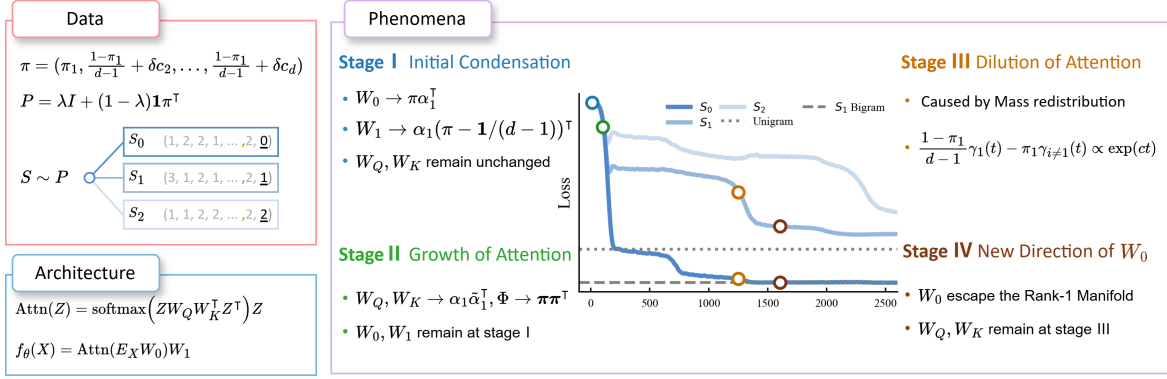


Figure 1. Overview of the setting and the focus–dilution training pattern. (Left) Sequences are generated by a Markov chain with stationary distribution π . We extract dataset S_0, S_1, S_2 from the training set, which differ only in the identity of the final token in each sequence. (Right) The loss curves exhibit four stages: initial condensation, attention growth, attention dilution, and the emergence of a new direction.

We study the gradient flow $\dot{\theta} = -\nabla\mathcal{L}(\theta)$.

Proposition 2.3 (Gradient flow and population-gradient limit). *The gradient flow dynamics satisfy*

$$\begin{cases} \frac{dW_0}{dt} = -\frac{\partial\mathcal{L}}{\partial M}W_1^\top - \frac{\partial\mathcal{L}}{\partial\Phi}W_0W_Q^\top - \left(\frac{\partial\mathcal{L}}{\partial\Phi}\right)^\top W_0W_QK, \\ \frac{dW_1}{dt} = -W_0^\top\frac{\partial\mathcal{L}}{\partial M}, \\ \frac{dW_Q}{dt} = -W_0^\top\frac{\partial\mathcal{L}}{\partial\Phi}W_0W_K, \\ \frac{dW_K}{dt} = -W_0^\top\left(\frac{\partial\mathcal{L}}{\partial\Phi}\right)^\top W_0W_Q. \end{cases} \quad (4)$$

Moreover, define the token-level proxy attention matrix $\mathbb{A} \in \mathbb{R}^{d \times d}$ by $\mathbb{A}_{i,j} = \frac{\pi_j \exp(e_i^\top \Phi e_j)}{\sum_{j'} \pi_{j'} \exp(e_i^\top \Phi e_{j'})}$ and the model output distribution $\mathbb{P} \in \mathbb{R}^{d \times d}$ by $\mathbb{P}_{i,j} = \frac{\exp(\mathbb{A}_i M e_j)}{\sum_{j'} \exp(\mathbb{A}_i M e_{j'})}$, where \mathbb{A}_i and \mathbb{P}_i denote the i -th row.

Then, in the large sample-size and long-context limit $(N, s) \rightarrow \infty$, the empirical gradients converge to

$$\begin{aligned} \lim_{N,s \rightarrow \infty} \frac{\partial\mathcal{L}}{\partial M} &= -\sum_{i=1}^d \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i), \\ \lim_{N,s \rightarrow \infty} \frac{\partial\mathcal{L}}{\partial\Phi} &= -\sum_{i=1}^d \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i). \end{aligned} \quad (5)$$

3. Theoretical results

3.1. Idea: stage-wise linearization around saddle points

Under small initialization, attention training often exhibits a multi-stage pattern: the trajectory spends a long time near a low-dimensional structure and then abruptly departs in a new direction. We explain this behavior via a stage-wise analysis around successive critical points. At each stage, the

parameters enter a neighborhood of a saddle point where the gradient flow is well-approximated by its linearization. The linearized dynamics exposes (i) stable directions that keep the trajectory confined to a low-dimensional manifold, and (ii) unstable directions that eventually dominate and trigger the transition to the next stage.

Concretely, we consider the gradient flow $\dot{\theta} = -\nabla\mathcal{L}(\theta)$. Let θ_* be a critical point, and define $\Delta\theta := \theta - \theta_*$. A Taylor expansion yields

$$\frac{d}{dt}\Delta\theta = -\nabla^2\mathcal{L}(\theta_*)\Delta\theta + \text{higher-order terms.}$$

The next lemma characterizes the linearization in which the nonlinear flow is governed by the linearized system, and formalizes the alignment with the most unstable direction.

Lemma 3.1 (Linearization near a saddle point). *Let $\dot{\theta} = F(\theta)$ be an ODE with $F \in C^2$, and let θ_* satisfy $F(\theta_*) = 0$. Let $J := DF(\theta_*)$ and assume there exist $r > 0$ and $L > 0$ such that for all $\|\Delta\theta\| \leq r$,*

$$\|F(\theta_* + \Delta\theta) - J\Delta\theta\| \leq L\|\Delta\theta\|^2. \quad (6)$$

Let $\theta(t)$ be the solution with $\|\Delta\theta(0)\| = \varepsilon \leq r/2$, and $\tilde{\Delta}\theta(t) := e^{Jt}\Delta\theta(0)$ be the solution of the linearized system $\dot{\Delta}\theta = J\tilde{\Delta}\theta$. Define $\mu := \sup\{\Re(\lambda) : \lambda \in \sigma(J)\}$. Then for all t such that $\|\tilde{\Delta}\theta(t)\| \leq r/2$,

$$\|\Delta\theta(t) - \tilde{\Delta}\theta(t)\| \leq C\varepsilon^2 e^{2\mu t} \quad (7)$$

for some constant $C = C(J, L)$. In particular, if $\mu > 0$, then the nonlinear dynamics is well-approximated by the linearized dynamics up to times $t = \Theta(\log(1/\varepsilon))$.

Moreover, suppose J is symmetric and has a simple eigenvalue $\mu > 0$ with eigenvector v_u and a spectral gap $\rho > 0$ in the sense that $\lambda \leq \mu - \rho$ for all $\lambda \in \sigma(J) \setminus \{\mu\}$. Then

for any initialization with $\langle \Delta\theta(0), v_u \rangle \neq 0$,

$$\frac{\Delta\theta(t)}{\|\Delta\theta(t)\|} \rightarrow \pm \frac{v_u}{\|v_u\|} \quad (8)$$

for any sequence $t = t(\varepsilon)$ with $t(\varepsilon) \rightarrow \infty$ and $\varepsilon e^{\mu t(\varepsilon)} \rightarrow 0$.

At initialization, each entry of every parameter matrix is sampled i.i.d. from $\mathcal{N}(0, \varepsilon^2)$ with $\varepsilon \ll 1$. Thus $\theta(0)$ lies in an $\mathcal{O}(\varepsilon)$ -neighborhood of the origin, which is a critical point of the gradient flow. By Lemma 3.1, the dynamics in the early time window of length $\Theta(\log(1/\varepsilon))$ is governed by the linearization at $\theta = 0$. A key consequence is that the linearized system admits a single unstable direction, so trajectories rapidly align with a rank-one direction. In our setting, this direction is not arbitrary: it is explicitly pinned down by the stationary distribution π of the underlying token Markov chain.

Theorem 3.2 (Initial condensation (rephrased from Thm. 2 in (Chen & Luo, 2025))). *The origin is a critical point and*

$$\frac{\partial \mathcal{L}}{\partial M} \Big|_{\theta=0} = -\pi \left(\pi - \frac{1}{d} \mathbf{1} \right)^\top, \quad \frac{\partial \mathcal{L}}{\partial \Phi} \Big|_{\theta=0} = 0. \quad (9)$$

The effective dynamics near $\theta = 0$ is

$$\frac{d\Delta W_0}{dt} = -\frac{\partial \mathcal{L}}{\partial M} \Big|_{\theta=0} \Delta W_0^\top, \quad \frac{d\Delta W_1}{dt} = -\Delta W_0^\top \frac{\partial \mathcal{L}}{\partial M} \Big|_{\theta=0} \quad (10)$$

Consequently, there exist a vector α_1 such that the following limit holds as $\varepsilon \rightarrow 0$ at $t = \Theta(\log \frac{1}{\varepsilon})$:

$$\frac{W_0}{\|W_0\|} \rightarrow \frac{\pi}{\|\pi\|} \alpha_1^\top, \quad \frac{W_1}{\|W_1\|} \rightarrow \alpha_1 \frac{\pi^\top - \frac{1}{d} \mathbf{1}^\top}{\|\pi^\top - \frac{1}{d} \mathbf{1}^\top\|}. \quad (11)$$

Theorem 3.2 characterizes the first stage of training dynamics in our model. Although it is rephrased from Thm. 2 in (Chen & Luo, 2025), we emphasize a more concrete interpretation relevant to data. As a result, (W_0, W_1) rapidly condense onto a π -driven rank-one structure within time $T_1 = \Theta(\log(1/\varepsilon))$. In contrast, the attention block (W_Q, W_K) stays $\mathcal{O}(\varepsilon)$ throughout this stage because the linear term in its dynamics vanishes at the origin, i.e., $\partial \mathcal{L} / \partial \Phi|_{\theta=0} = 0$.

3.2. Focus of Attention

After initial condensation stage, outer parameters (W_0, W_1) rapidly become approximately rank-one, while the attention parameters (W_Q, W_K) remain $\mathcal{O}(\varepsilon)$. Empirically, the trajectory then stays close to the rank-one condensation ray where the outer parameters evolve along the same direction until it enters a neighborhood of a second critical point.

Proposition 3.3 (Existence of a second critical point on the condensation ray). *Assume $0 < \pi_d = \dots = \pi_2 = \pi_1 < 1$.*

Then there exists $\kappa_1 > 0$ such that the parameter tuple θ_c^1

$$W_0 = \kappa_1 \frac{\pi}{\|\pi\|} \alpha_1^\top, \quad W_1 = \kappa_1 \alpha_1 \frac{\pi^\top - \frac{1}{d} \mathbf{1}^\top}{\|\pi^\top - \frac{1}{d} \mathbf{1}^\top\|}, \quad W_Q, W_K = 0, \quad (12)$$

satisfies $\mathbb{P}_i = \pi^\top$ for all i and $\frac{\partial \mathcal{L}}{\partial M} \Big|_{\theta=\theta_c^1} = 0$. Moreover, θ_c^1 is a critical point of the full gradient flow.

The key point is that θ_c^1 is typically a saddle: the (W_0, W_1) -subsystem is contracting (or neutrally stable due to symmetry), while the (W_Q, W_K) -subsystem contains an unstable mode.

Proposition 3.4 (Linearized dynamics and its unique unstable direction). *At critical point θ_c^1 defined in Prop. 3.3, the linearization of the gradient flow admits the block form*

$$\frac{d}{dt} \begin{pmatrix} \Delta W_0 \\ \Delta W_1 \\ \Delta W_Q \\ \Delta W_K \end{pmatrix} = \begin{pmatrix} J_{\text{out}} & 0 \\ 0 & J_{\text{att}} \end{pmatrix} \begin{pmatrix} \Delta W_0 \\ \Delta W_1 \\ \Delta W_Q \\ \Delta W_K \end{pmatrix} \quad (13)$$

where J_{out} is negative semi-definite and J_{att} is positive semi-definite. It indicates that the current dynamics are dominated by the attention subsystem. In particular, the attention block satisfies the explicit closed system

$$\frac{d}{dt} \Delta W_Q = c \alpha_1 \alpha_1^\top \Delta W_K, \quad \frac{d}{dt} \Delta W_K = c \alpha_1 \alpha_1^\top \Delta W_Q, \quad (14)$$

where $c = \lambda \kappa_1^4 \frac{\|\pi\|_{\text{var}(\pi)}^4}{\|\pi - \frac{1}{d} \mathbf{1}\| \|\pi\|^3} > 0$. Therefore the attention block has an exponentially unstable mode.

By Lemma 3.1, once the trajectory enters an $\mathcal{O}(\varepsilon)$ -neighborhood of θ_c^1 , the dynamics is governed by the linearization for a duration $\Theta(\log(1/\varepsilon))$. It implies that the attention parameters (W_Q, W_K) converge into the direction depending on the condensation direction. As a result, the attention structure prioritizes tokens that appear frequently in the steady-state distribution, indicating that the attention mechanism has become specific.

Theorem 3.5 (High frequency token bias). *Suppose the trajectory enters an $\mathcal{O}(\varepsilon)$ -neighborhood of θ_c^1 . Within the linearization neighborhood of Lemma 3.1, there exists a unit vector $\tilde{\alpha}_1 \in \mathbb{R}^m$ such that, for generic small initialization of (W_Q, W_K) ,*

$$\frac{W_Q(t)}{\|W_Q(t)\|} \rightarrow \alpha_1 \tilde{\alpha}_1^\top, \quad \frac{W_K(t)}{\|W_K(t)\|} \rightarrow \alpha_1 \tilde{\alpha}_1^\top. \quad (15)$$

Consequently, along the unstable ray $W_Q = W_K = \kappa \alpha_1 \tilde{\alpha}_1^\top$, the attention score matrix satisfies

$$\frac{\Phi}{\|\Phi\|} = \frac{W_0 W_Q W_K^\top W_0^\top}{\|W_0 W_Q W_K^\top W_0^\top\|} = \pi \pi^\top. \quad (16)$$

Then for each $i \in [d]$ the attention distribution exhibits a high-frequency bias:

$$\lim_{\kappa \rightarrow \infty} \mathbb{A}_i(W_0, W_1, \kappa \alpha_1 \tilde{\alpha}_1^\top, \kappa \alpha_1 \tilde{\alpha}_1^\top) = e_1^\top. \quad (17)$$

3.3. Dilution of Attention

Sec. 3.2 shows that the second critical point θ_c^1 is a saddle whose unique unstable direction lies in the attention subsystem: after a transient of length $\Theta(\log(1/\varepsilon))$, the attention parameters become approximately rank-1 and aligned while the outer parameters remain close to their initial values on the condensation ray. In this subsection, we stay in the same neighborhood of θ_c^1 but go beyond linearization: Conditioned on the rank-1 manifold, we resolve the next-order perturbations in embeddings induced by the evolution of attention. This refinement reveals a redistribution effect in the embeddings that gradually undermines the previously formed focus, leading to the dilution phase.

Motivated by the alignment result in Sec. 3.2, we model the post-transient phase by the rank-one parametrization

$$\begin{aligned} W_0 &= \gamma(t) \alpha_1^\top, & W_1 &= \alpha_1 \beta(t)^\top, \\ W_Q &= \lambda_Q(t) \alpha_1 \tilde{\alpha}_1^\top, & W_K &= \lambda_K(t) \alpha_1 \tilde{\alpha}_1^\top, \end{aligned} \quad (18)$$

where $\gamma(t), \beta(t) \in \mathbb{R}^d$ and $\lambda_Q(t), \lambda_K(t) \in \mathbb{R}$. At the entry time t_0 of this phase,

$$\begin{aligned} \gamma(t_0) &= \kappa_1 \frac{\pi}{\|\pi\|}, & \beta(t_0) &= \kappa_1 \frac{\pi - \frac{1}{d} \mathbf{1}}{\|\pi - \frac{1}{d} \mathbf{1}\|}, \\ \lambda_Q(t_0) &= \lambda_K(t_0) = o(1). \end{aligned} \quad (19)$$

For notational convenience, we also define the attention amplitude $\eta(t) := \lambda_Q(t) \lambda_K(t)$ which is the only combination that enters the reduced dynamics below.

Proposition 3.6 (Invariant rank-one manifold). *Assume $\pi_1 > \pi_2 = \dots = \pi_d$. Define*

$$\mathcal{W} := \left\{ \theta \text{ satisfying (18) for some } (\gamma, \beta, \lambda_Q, \lambda_K) \right\}.$$

If $(W_0, W_1, W_Q, W_K) \in \mathcal{W}$ at time t_0 , then the gradient flow (4) remains in \mathcal{W} for all $t \geq t_0$. Moreover, if $\gamma_2(t_0) = \dots = \gamma_d(t_0)$ and $\beta_2(t_0) = \dots = \beta_d(t_0)$, it will be preserved for any $t \geq t_0$.

Restricting the gradient flow to \mathcal{W} yields a closed system in (γ, β, η) :

$$\begin{aligned} \dot{\gamma} &= -\frac{\partial \mathcal{L}}{\partial M} \beta - \eta \left(\frac{\partial \mathcal{L}}{\partial \Phi} + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \right) \gamma, \\ \dot{\beta} &= -\left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top \gamma, & \dot{\eta} &= -2\eta \gamma^\top \frac{\partial \mathcal{L}}{\partial \Phi} \gamma. \end{aligned} \quad (20)$$

By Proposition 3.6, it suffices to track two-group coordinates

$$\begin{aligned} \gamma_1, & \quad \gamma_2 = \dots = \gamma_{|\mathcal{V}|} := \gamma_{i \neq 1}, \\ \beta_1, & \quad \beta_2 = \dots = \beta_{|\mathcal{V}|} := \beta_{i \neq 1}, \end{aligned}$$

and denote $\Delta\gamma := \gamma_1 - \gamma_{i \neq 1}$, $\Delta\beta := \beta_1 - \beta_{i \neq 1}$. Intuitively, this reduces the post-alignment dynamics to an effective

two-group system (token 1 versus all others). Importantly, we are *still* analyzing the flow near the same critical point θ_c^1 ; the difference from Sec. 3.2 is that we can keep the attention direction fixed and resolve the next-order feedback that governs redistribution on the rank-one manifold.

Theorem 3.7 (Mass redistribution). *Consider the linearization of reduced dynamics (20) on \mathcal{W} at critical point corresponding to θ_c^1 . There exists $c > 0$ such that*

$$(1 - \pi_1) \gamma_1(t) - (d - 1) \pi_1 \gamma_{i \neq 1}(t) \propto \exp(ct). \quad (21)$$

Consequently, $\gamma_1(t)$ and $\gamma_{i \neq 1}(t)$ cannot move in the same direction: a weighted contrast between high-frequency token and the remaining tokens is exponentially amplified.

Theorem 3.7 explains the mechanism behind the dilution phase. After alignment, the attention direction is essentially fixed, and the attention amplitude $\eta(t)$ feeds back into $\dot{\gamma}$ through the $\eta(\partial \mathcal{L} / \partial \Phi) \gamma$ term in (20). The redistribution effect forces a growing separation between γ_1 and $\gamma_{i \neq 1}$, so the embedding mass cannot remain concentrated along π .

Therefore, the logit difference that causes high-frequency bias gradually weakens: the attention weights corresponding to low-frequency tokens are no longer concentrated on high-frequency tokens. This marks a shift from focus to dilution.

3.4. Emergence of a new direction via data asymmetry

In Sec. 3.3, the dynamics collapses onto a rank-one invariant manifold, effectively reducing learning to a “token 1 vs. all others” two-group system. Further learning requires separating low-frequency states, which demands growth of embeddings along directions that distinguish low-frequency tokens.

However, for perfectly symmetric data among low-frequency tokens, the rank-one manifold may contain a degenerate critical point where both driving forces vanish: $\partial \mathcal{L} / \partial M = 0$ and $\partial \mathcal{L} / \partial \Phi = 0$. Crucially, the degeneracy is not only tangential, but also transverse. As a consequence, linearization does not generate a mechanism that pushes the trajectory away from the rank-one manifold.

Proposition 3.8 (Degenerate critical point). *Assume perfect symmetry among low-frequency tokens. On the rank-one invariant manifold \mathcal{W} (Proposition 3.6), there exists a critical point, which is a neutrally stable equilibrium for the linearized dynamics, such that $\frac{\partial \mathcal{L}}{\partial M} = 0$ and $\frac{\partial \mathcal{L}}{\partial \Phi} = 0$.*

The solution to remove degeneracy is to introduce perturbations that breaks the symmetry. In practice, low-frequency tokens rarely have identical frequencies. To model a minimal asymmetry while keeping calculations simple, we focus on $d = 3$ and perturb the stationary distribution by a small parameter δ :

$$\pi^\top = (\pi_1, \frac{1-\pi_1}{2}, \frac{1-\pi_1}{2}) \Rightarrow \tilde{\pi}^\top = (\pi_1, \frac{1-\pi_1}{2} + \delta, \frac{1-\pi_1}{2} - \delta)$$

We study stationary points of the perturbed gradient field $-\nabla_{\theta}\mathcal{L}(\theta, \delta) = 0$ near the degenerate critical point. After shifting coordinates so that $\theta = 0$ corresponds to the degenerate critical point, a formal expansion takes the form

$$-\nabla\mathcal{L}(\theta, \delta) = J_0\theta + \delta f_1 + \text{h.o.t.}, \quad (22)$$

where $J_0 = -\nabla_{\theta}^2\mathcal{L}$ at $\delta = 0$.

If J_0 were invertible, the implicit function theorem would apply, and we could directly obtain the solution $\theta(\delta)$. Unfortunately, due to symmetry, J_0 is degenerate, so we use the standard Lyapunov–Schmidt reduction. Let $Q_K = (k_1, \dots, k_{d_K})$ and $Q_R = (q_1, \dots, q_{d_R})$ be orthonormal bases for the kernel and range subspaces of J_0 :

$$\mathbb{R}^p = \text{Ker}(J_0) \oplus \text{Ran}(J_0), \quad \theta = Q_K x + Q_R y,$$

where p is the parameter dimension. Projecting the stationarity condition onto the range and kernel yields the equivalent system

$$-Q_R^T \nabla\mathcal{L}(\theta, \delta) = 0, \quad -Q_K^T \nabla\mathcal{L}(\theta, \delta) = 0. \quad (23)$$

The range equation can be solved by the implicit function theorem since $Q_R^T J_0 Q_R$ is invertible, yielding a smooth map $y = \zeta(x, \delta)$. Substituting back into the kernel equation produces a reduced low-dimensional problem in x whose solutions describe nearby stationary points.

The key effect of the perturbation is that it splits the previously flat transverse directions. A genuinely transverse positive eigenvalue of order $\Theta(\delta)$ appears, while tangential instability is at most $\mathcal{O}(\delta^2)$. This fast transverse instability is what drives the trajectory away from the rank-one manifold and seeds a new embedding direction.

Theorem 3.9 (Asymmetry lifts degeneracy and induces a new direction). *Consider the perturbed stationary distribution with parameter δ above. There exists a point $\theta(\delta)$ near the degenerate critical point such that*

$$\|\nabla_{\theta}\mathcal{L}(\theta(\delta), \delta)\| = \mathcal{O}(\delta^3). \quad (24)$$

Moreover, the Hessian at $\theta(\delta)$ exhibits two distinct scales:

1. *Slow tangential instability. Any positive eigenvalues created from the previously degenerate directions are at most $\mathcal{O}(\delta^2)$.*
2. *Fast normal instability. Under mild condition, in directions transverse to the rank-one manifold, there exists a positive eigenvalue of order $\Theta(\delta)$.*

Theorem 3.9 shows that any generic low-frequency asymmetry lifts this degeneracy and produces a fast transverse unstable mode of size $\Theta(\delta)$. Once the trajectory enters the neighborhood of $\theta(\delta)$, this transverse instability drives it away from the degenerate rank-one configuration and enables the emergence of a genuinely new embedding direction, allowing the model to further differentiate low-frequency tokens beyond the two-group description.

4. Empirical evidence

In this section, leveraging the simplified transformer model, we analyze the training behavior on Markovian data and empirically validate the theoretical derivation describing the transition of attention from focus to dilution. In parallel, we evaluate the model on real-world WikiText corpora and on the TinyStories corpus, which exhibits basic linguistic structure, to assess whether our observations generalize to the training dynamics of large-scale language models in realistic settings.

4.1. Synthetic Experiments

We construct synthetic datasets using four distinct transition matrices P , designed to share a common stationary distribution $\pi = (0.75, 0.19, 0.05, 0.01)$. To reveal the low-rank structure of parameters caused by condensation, we measured the cosine similarity between neuronal input weights for analysis (Chen & Luo, 2025; Xu et al., 2025b).

Additionally, we visualize the embedding trajectory evolution by applying Principal Component Analysis (PCA) to the concatenated embedding snapshots across all steps (Lorch, 2016; Antognini & Sohl-Dickstein, 2018). Detailed experimental setups are provided in Appendix E. The overall evolution of the training dynamics is visualized in Fig. 2. We identify four distinct stages during the training process. In the following, we provide a detailed analysis of each stage to demonstrate the consistency between our experimental observations and theoretical results.

Stage I: Initial Condensation Our theoretical analysis predicts that during this stage, the outer layers W_0, W_1 evolve from an initialized full-rank state to a low-rank structure, while the inner attention parameters remain largely invariant. This is depicted by Fig. 2(A), which demonstrates that the outer weights rapidly evolve into rank-1, whereas the W_Q, W_K maintain the high-rank nature of their initialization. Simultaneously, Fig. 2(B1) illustrates that the embeddings of all tokens evolve towards a uniform direction, further validating our theoretical analysis.

Stage II: Growth of Attention During this stage, the outer parameters remain largely invariant, while W_Q and W_K transition into a condensed state. This phase coincides with a significant drop in training loss, marking the evolution of parameters from the origin to the next critical point. According to our theory, the attention mechanism evolves such that high-frequency tokens are gradually focused by the remaining tokens. This phenomenon is clearly visualized in Fig. 2(C).

Stage III: Dilution of Attention In Stage III, although the parameters remain confined to the rank-1 manifold

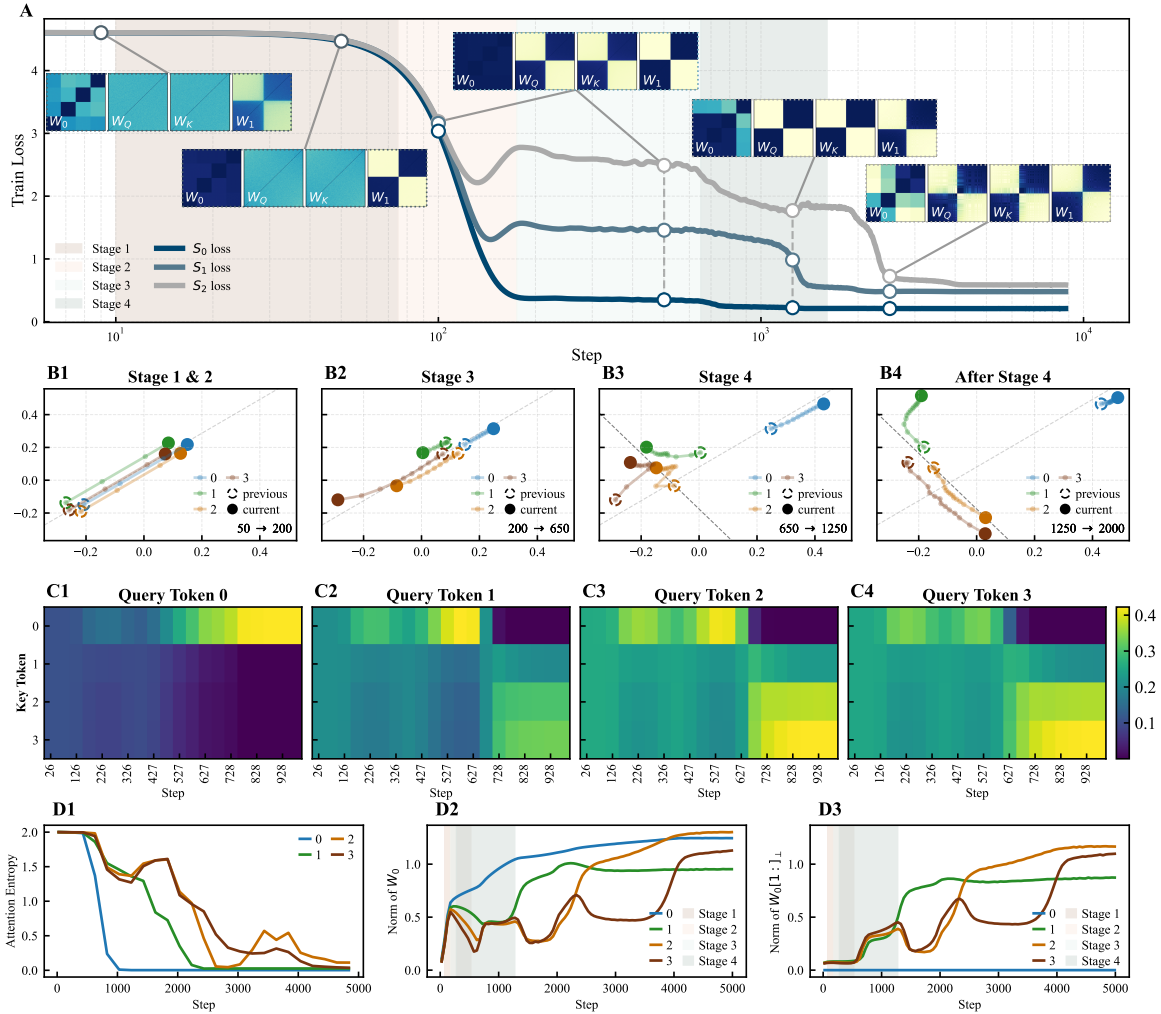


Figure 2. **Empirical focus–dilution cycle.** (A) Loss curves for S_0, S_1, S_2 , accompanied by cosine-similarity matrices for (W_0, W_Q, W_K, W_1) . (B) PCA of embeddings reveals one-directional growth (Stage I & II), retraction during dilution (Stage III), and expansion into new directions (Stage IV). (C) Attention maps transition from focus to dilution. Between steps 350 and 550, the attention given to the first token by all query pairs increases synchronously. Afterward, the first token focuses attention on itself, while other tokens reduce their attention towards it. (D1) Attention entropy; (D2), (D3) The embedding norm and the norm of perpendicular component onto the first token.

(evidenced by the unchanged condensation heatmap in Fig. 2(A)), Fig. 2(B2) reveals that all tokens, except for token 0, exhibit a retraction trajectory. This implies that while the training dynamics are strictly constrained within the low-rank manifold, the model begins to differentiate between tokens. As shown in Fig. 2(C, D), low-frequency tokens pay less attention to high-frequency token in this phase, accompanied by a significant drop in the embedding norms of low-frequency tokens. Consequently, outer parameters of the network revert to an unstable state.

Stage IV: Emergence of New Direction In Stage IV, the accumulated instability drives the model to escape the constraints of the rank-1 manifold, initiated by the growth of new directions in the outer layers. This transition is

clearly observable in Fig. 2(B3). To further quantify this, in Fig. 2(D3), we project the embeddings of low-frequency tokens ($W_0[1:]$) onto the direction orthogonal to token 0 (denoted as $W_0[0]_{\perp}$) and calculate the projection norms. The results indicate that, for the first time, the remaining token embeddings significantly deviate from the direction of token 0.

After Stage IV: Subsequent Training Dynamics Fig. 2(D) further illustrates the later stages of the training process, revealing a distinct periodicity in the embedding norms. Specifically, the focus and dilution pattern repeats recursively: as the network proceeds to learn Token 1, the remaining tokens (2–3) undergo the same retraction and regrowth process, continuing sequentially until training concludes.

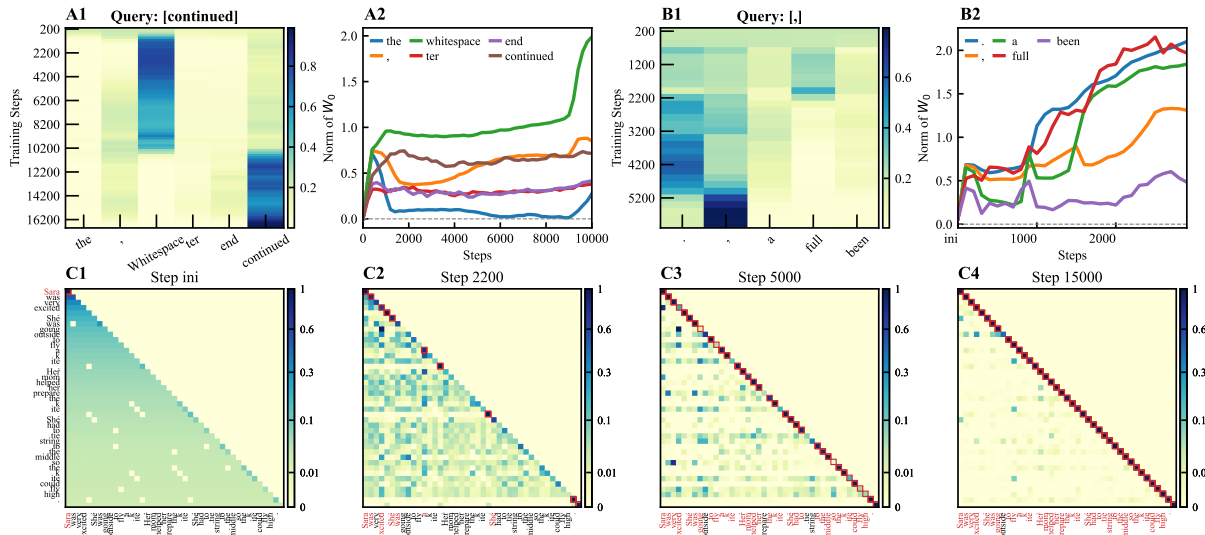


Figure 3. **Experimental results on real-world datasets.** (A) Results on WikiText. (A1) The attention evolution of the medium-frequency token ‘continue’ given the input sequence [the, comma, whitespace, ter, end, continue]. The attention scores exhibit a distinct four-phase transition: dilution → focus on the high-frequency whitespace → secondary dilution → final focus on continue itself. (A2) The evolution of $\|W_0\|_2$ for selected tokens. The dominant *whitespace* token shows continuous growth, while other tokens display a clear retraction. (B) Results on TinyStories. We observe similar dynamics using the input sequence [the, comma, a, full, been], mirroring the phenomena in (A). (C) Visualization of attention shifts. The evolution of attention for a single test sample across different training steps, with tokens exhibiting self-attention scores ≥ 0.75 highlighted.

We hypothesize that after Stage IV, the model has effectively converged on Token 0. Consequently, the system evolves into a sub-dynamic regime governed by Tokens 1–3. In this reduced state, the parameter dynamics can be re-analyzed within our original theoretical framework.

4.2. Real-world Experiments

Experimental Results. We validate the correctness of our theorem on two real-world datasets: WikiText (Merity et al., 2016) and TinyStories (Eldan & Li, 2023). We employ the same simplified Transformer architecture and maintain hyperparameter settings consistent with the synthetic data experiments. To investigate the “focus-and-dilution” characteristics of the attention mechanism, we track the top-three most frequent tokens alongside three randomly sampled medium-frequency tokens (frequency $> 10\%$) from the training set.

As illustrated in the figure, across both WikiText and TinyStories, the attention mechanism exhibits a consistent pattern: it initially prioritizes high-frequency tokens (e.g., “the”, “a”, and whitespace) before subsequently losing this focus—a process we term “dilution.” Concurrently, by monitoring the embedding evolution of these selected tokens, we observe a distinct “retraction” phenomenon. Notably, due to the high variance in batch composition inherent to the 1-epoch training regime on real-world corpora, we occasionally observe this retraction even in the most frequent tokens.

Validity of the Markov Approximation. Existing studies (Chang & Bergen, 2022; Chang et al., 2024) indicate a curriculum in Transformer learning, starting from 1-gram to n-gram statistics. Our analysis of attention patterns in TinyStories supports this: distinct tokens gradually shift from uniform attention to self-attention. This behavior indicates that the model functions as a pseudo-2-gram model during early training phases, despite the non-Markovian nature of real text. These observations indirectly validate our experimental design, confirming that our synthetic Markov data acts as a suitable proxy for understanding real-world training dynamics.

5. Conclusion

This work advances the theoretical understanding of transformer training dynamics by providing a mechanistic account of how attention evolves. We identify a recurring focus–dilution cycle and develop a stage-wise gradient-flow framework that characterizes rank-one condensation, saddle-to-saddle transitions, and the impact of symmetry breaking, offering a rigorous basis for phenomena often reported empirically. While the formal analysis is derived in a restricted setting, the same qualitative signatures appear on real corpora, suggesting the framework remains a useful lens for interpreting early-stage attention dynamics beyond the idealized regime.

Impact Statement

This paper aims to advance the theoretical understanding of transformer training dynamics by providing a mechanistic analysis of attention evolution. While improved understanding of learning dynamics may inform future model design and training practices, we do not foresee any direct negative ethical or societal consequences arising specifically from this work.

References

- Antognini, J. and Sohl-Dickstein, J. Pca of high dimensional random walks with comparison to neural network training. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/7a576629fef88f3e636afd33b09e8289-Paper.pdf.
- Arora, S., Du, S. S., Hu, W., Li, Z., Salakhutdinov, R. R., and Wang, R. On exact computation with an infinitely wide neural net. *Advances in neural information processing systems*, 32, 2019.
- Bietti, A., Cabannes, V., Bouchacourt, D., Jegou, H., and Bottou, L. Birth of a transformer: A memory viewpoint. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=3X2EbBLNsk>.
- Chang, T. A. and Bergen, B. K. Word acquisition in neural language models. *Transactions of the Association for Computational Linguistics*, 10:1–16, 2022. doi: 10.1162/tacl.a.00444. URL <https://aclanthology.org/2022.tacl-1.1/>.
- Chang, T. A., Tu, Z., and Bergen, B. K. Characterizing learning curves during language model pre-training: Learning, forgetting, and stability. *Transactions of the Association for Computational Linguistics*, 12:1346–1362, 2024. doi: 10.1162/tacl.a.00708. URL <https://aclanthology.org/2024.tacl-1.74/>.
- Chen, S., Sheen, H., Wang, T., and Yang, Z. Unveiling induction heads: Provable training dynamics and feature learning in transformers. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 66479–66567. Curran Associates, Inc., 2024a. doi: 10.52202/079017-2127. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/7aae9e3ec211249e05bd07271a6b1441-Paper-Conference.pdf.
- Chen, Z.-A. and Luo, T. From condensation to rank collapse: A two-stage analysis of transformer training dynamics. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=gm5mkiTGOy>.
- Chen, Z.-A., Li, Y., Luo, T., Zhou, Z., and Xu, Z.-Q. J. Phase diagram of initial condensation for two-layer neural networks. *CSIAM Transactions on Applied Mathematics*, 5(3):448–514, 2024b. ISSN 2708-0579. doi: <https://doi.org/10.4208/csiam-am.SO-2023-0016>. URL <https://global-sci.com/article/91025/phase-diagram-of-initial-condensation-for-two-layer-neural-networks>.
- Eldan, R. and Li, Y. Tinystories: How small can language models be and still speak coherent english?, 2023. URL <https://arxiv.org/abs/2305.07759>.
- Gao, C., Cao, Y., Li, Z., He, Y., Wang, M., Liu, H., Klusowski, J. M., and Fan, J. Global convergence in training large-scale transformers. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 29213–29284. Curran Associates, Inc., 2024. doi: 10.52202/079017-0921. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/33b47b3d2441a17b95344cd635f3dd01-Paper-Conference.pdf.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Ji, Z. and Telgarsky, M. Gradient descent aligns the layers of deep linear networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HJflg30qKX>.
- Kim, J. and Suzuki, T. Transformers learn nonlinear features in context: Nonconvex mean-field dynamics on the attention landscape. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 24527–24561. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/kim24af.html>.
- Kumar, A. and Haupt, J. Early directional convergence in deep homogeneous neural networks for small initializations. *arXiv preprint arXiv:2403.08121*, 2024.
- Li, Y., Li, Y., and Risteski, A. How do transformers learn topic structure: Towards a mechanistic understanding, 2023. URL <https://arxiv.org/abs/2303.04245>.

- 495 Lorch, E. Visualizing deep network training trajectories
496 with pca. In *ICML Workshop on Visualization for Deep*
497 *Learning*, 2016.
- 498 Lu, H., Mao, Y., and Nayak, A. On the dynamics of
499 training attention models. In *International Conference*
500 *on Learning Representations*, 2021. URL <https://openreview.net/forum?id=1OCTOShAmqB>.
- 503 Luo, T., Xu, Z.-Q. J., Ma, Z., and Zhang, Y. Phase diagram
504 for two-layer relu neural networks at infinite-width limit.
505 *The Journal of Machine Learning Research*, 22(1):3327–
506 3373, 2021.
- 508 Makuva, A. V., Bondaschi, M., Girish, A., Nagle, A., Kim,
509 H., Gastpar, M., and Ekbote, C. Local to global: Learn-
510 ing dynamics and effect of initialization for transform-
511 ers. In *The Thirty-eighth Annual Conference on Neural*
512 *Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=OX4y1l3X53>.
- 515 Makuva, A. V., Bondaschi, M., Girish, A., Nagle, A.,
516 Jaggi, M., Kim, H., and Gastpar, M. Attention with
517 markov: A curious case of single-layer transformers. In
518 *The Thirteenth International Conference on Learning*
519 *Representations*, 2025. URL [https://openreview](https://openreview.net/forum?id=SqZ0KY4qBD)
520 [.net/forum?id=SqZ0KY4qBD](https://openreview.net/forum?id=SqZ0KY4qBD).
- 521 Mei, S., Montanari, A., and Nguyen, P.-M. A mean field
522 view of the landscape of two-layer neural networks. *Pro-*
523 *ceedings of the National Academy of Sciences*, 115(33):
524 E7665–E7671, 2018.
- 526 Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer
527 sentinel mixture models, 2016.
- 528 Pérez, J., Marinković, J., and Barceló, P. On the turing
529 completeness of modern neural network architectures. In
530 *International Conference on Learning Representations*,
531 2019. URL [https://openreview.net/forum](https://openreview.net/forum?id=HyGBdo0qFm)
532 [?id=HyGBdo0qFm](https://openreview.net/forum?id=HyGBdo0qFm).
- 533 Rajaraman, N., Jiao, J., and Ramchandran, K. An analysis
534 of tokenization: Transformers under markov data. In *The*
535 *Thirty-eighth Annual Conference on Neural Information*
536 *Processing Systems*, 2024. URL [https://openre](https://openreview.net/forum?id=wm9JZq7RCe)
537 [view.net/forum?id=wm9JZq7RCe](https://openreview.net/forum?id=wm9JZq7RCe).
- 538 Rotskoff, G. and Vanden-Eijnden, E. Parameters as inter-
539 acting particles: long time convergence and asymptotic
540 error scaling of neural networks. *Advances in neural*
541 *information processing systems*, 31, 2018.
- 542 Sheen, H., Chen, S., Wang, T., and Zhou, H. H. Implicit
543 regularization of gradient flow on one-layer softmax at-
544 tention, 2024. URL [https://arxiv.org/abs/24](https://arxiv.org/abs/2403.08699)
545 [03.08699](https://arxiv.org/abs/2403.08699).
- 546 Snell, C. B., Zhong, R., Klein, D., and Steinhardt, J. Ap-
547 proximating how single head attention learns. *ArXiv*,
548 [abs/2103.07601](https://arxiv.org/abs/2103.07601), 2021. URL [https://api.semant](https://api.semanticscholar.org/CorpusID:232232786)
549 [icscholar.org/CorpusID:232232786](https://api.semanticscholar.org/CorpusID:232232786).
- Tarzanagh, D. A., Li, Y., Thrampoulidis, C., and Oymak, S.
Transformers as support vector machines. *arXiv preprint*
arXiv:2308.16898, 2023.
- Tian, Y., Wang, Y., Zhang, Z., Chen, B., and Du, S. S. JoMA:
Demystifying multilayer transformers via joint dynam-
ics of MLP and attention. In *The Twelfth International*
Conference on Learning Representations, 2024. URL
[https://openreview.net/forum?id=LbJq](https://openreview.net/forum?id=LbJqRGNYCf)
[RGNYCf](https://openreview.net/forum?id=LbJqRGNYCf).
- Varre, A. V., Vladarean, M.-L., Pillaud-Vivien, L., and
Flammarion, N. On the spectral bias of two-layer lin-
ear networks. In *Thirty-seventh Conference on Neural*
Information Processing Systems, 2023. URL <https://openreview.net/forum?id=FFdrXkm3Cz>.
- Vasudeva, B., Deora, P., and Thrampoulidis, C. Implicit bias
and fast convergence rates for self-attention. *Transactions*
on Machine Learning Research, 2025. ISSN 2835-8856.
URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=pKilnjQsb0)
[pKilnjQsb0](https://openreview.net/forum?id=pKilnjQsb0).
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. At-
tention is all you need. *Advances in neural information*
processing systems, 30, 2017.
- Wang, M. and Ma, C. Understanding multi-phase opti-
mization dynamics and rich nonlinear behaviors of reLU
networks. In *Thirty-seventh Conference on Neural In-*
formation Processing Systems, 2023. URL <https://openreview.net/forum?id=konBXvt2iS>.
- Williams, F., Trager, M., Panozzo, D., Silva, C., Zorin, D.,
and Bruna, J. Gradient dynamics of shallow univariate
relu networks. *Advances in neural information processing*
systems, 32, 2019.
- Wu, D., Shevchenko, A., Oymak, S., and Mondelli, M.
Attention with trained embeddings provably selects im-
portant tokens, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2505.17282)
[abs/2505.17282](https://arxiv.org/abs/2505.17282).
- Xu, Z., Min, H., Luo, J., MacDonald, L. E., Tarmoun, S.,
Mallada, E., and Vidal, R. Understanding the learning dy-
namics of loRA: A gradient flow perspective on low-rank
adaptation in matrix factorization. In *The 28th Interna-*
tional Conference on Artificial Intelligence and Statistics,
2025a. URL [https://openreview.net/forum](https://openreview.net/forum?id=hphdX8Wlct)
[?id=hphdX8Wlct](https://openreview.net/forum?id=hphdX8Wlct).

- 550 Xu, Z.-Q. J., Zhang, Y., and Zhou, Z. An overview of
551 condensation phenomenon in deep learning, 2025b. URL
552 <https://arxiv.org/abs/2504.09484>.
553
- 554 Yang, H., Kailkhura, B., Wang, Z., and Liang, Y. Train-
555 ing dynamics of transformers to recognize word co-
556 currence via gradient flow analysis. In Globerson, A.,
557 Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak,
558 J., and Zhang, C. (eds.), *Advances in Neural Information
559 Processing Systems*, volume 37, pp. 46047–46117. Cur-
560 ran Associates, Inc., 2024. doi: 10.52202/079017-1465.
561 URL [https://proceedings.neurips.cc/p
562 aper_files/paper/2024/file/520416e27
563 d3b0cef3cd70a083e2991c7-Paper-Confe
564 rence.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/520416e27d3b0cef3cd70a083e2991c7-Paper-Conference.pdf).
- 565
- 566 Yang, T., Huang, Y., Liang, Y., and Chi, Y. Multi-head trans-
567 formers provably learn symbolic multi-step reasoning via
568 gradient descent. In *The Thirty-ninth Annual Conference
569 on Neural Information Processing Systems*, 2025. URL
570 [https://openreview.net/forum?id=qFC7
571 28XyeM](https://openreview.net/forum?id=qFC728XyeM).
572
- 573 Yun, C., Bhojanapalli, S., Rawat, A. S., Reddi, S., and
574 Kumar, S. Are transformers universal approximators
575 of sequence-to-sequence functions? In *International
576 Conference on Learning Representations*, 2020a. URL
577 [https://openreview.net/forum?id=ByxR
578 M0Ntvr](https://openreview.net/forum?id=ByxRM0Ntvr).
579
- 580 Yun, C., Chang, Y.-W., Bhojanapalli, S., Rawat, A. S.,
581 Reddi, S., and Kumar, S. O (n) connections are expressive
582 enough: Universal approximability of sparse transform-
583 ers. *Advances in Neural Information Processing Systems*,
584 33:13783–13794, 2020b.
- 585
- 586 Zhang, R., Frei, S., and Bartlett, P. L. Trained transform-
587 ers learn linear models in-context. *Journal of Machine
588 Learning Research*, 25(49):1–55, 2024a.
- 589
- 590 Zhang, Y., Xu, Z.-Q. J., Luo, T., and Ma, Z. A type of
591 generalization error induced by initialization in deep neu-
592 ral networks. In *Mathematical and Scientific Machine
593 Learning*, pp. 144–164. PMLR, 2020.
- 594
- 595 Zhang, Y., Singh, A. K., Latham, P. E., and Saxe, A. M.
596 Training dynamics of in-context learning in linear atten-
597 tion. In *Forty-second International Conference on Ma-
598 chine Learning*, 2025a. URL [https://openreview
599 .net/forum?id=aFNq67ilos](https://openreview.net/forum?id=aFNq67ilos).
- 600
- 601 Zhang, Z., Lin, P., Wang, Z., Zhang, Y., and Xu, Z.-Q. J.
602 Initialization is critical to whether transformers fit com-
603 posite functions by inference or memorizing, 2024b. URL
604 <https://arxiv.org/abs/2405.05409>.
- Zhang, Z., Lin, P., Wang, Z., Zhang, Y., and Xu, Z.-Q. J. Complexity control facilitates reasoning-based compositional generalization in transformers. *arXiv preprint arXiv:2501.08537*, 2025b.
- Zhou, H., Zhou, Q., Luo, T., Zhang, Y., and Xu, Z.-Q. Towards understanding the condensation of neural networks at initial training. *Advances in Neural Information Processing Systems*, 35:2184–2196, 2022.
- Zhou, Z., Zhou, H., Li, Y., and Xu, Z.-Q. J. Understanding the initial condensation of convolutional neural networks. *arXiv preprint arXiv:2305.09947*, 2023.
- Zucchet, N., D’Angelo, F., Lampinen, A. K., and Chan, S. C. The emergence of sparse attention: impact of data distribution and benefits of repetition. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL [https://openreview.net
/forum?id=jMhRbV47pS](https://openreview.net/forum?id=jMhRbV47pS).

A. Related Works

Training dynamics of attention and multi-stage analysis Given the scale of modern models and the complexity of optimizers, studying the training dynamics of attention remains a challenging problem. A common practice is to introduce various simplifications to the research object, such as constructing task-specific synthetic data, utilizing reparameterization or simplified model and optimization function (Sheen et al., 2024; Kim & Suzuki, 2024; Varre et al., 2023; Chen et al., 2024a; Wu et al., 2025; Gao et al., 2024; Zhang et al., 2025a; Vasudeva et al., 2025; Yang et al., 2025). Among these, (Lu et al., 2021) establishes key dynamical identities using a controllable text classification task where sentences consist of a “topic word” plus random noise. (Snell et al., 2021) suggests that models first capture word co-occurrence before adjusting attention to focus on relevant tokens. (Li et al., 2023) examines the dynamical effects of fixing specific attention components within a topic-word task framework. Following these previous works, (Tian et al., 2024) proposed a novel mathematical framework for analyzing the joint dynamics of MLP and attention blocks, successfully explaining the sparsity of attention score matrices. (Zucchet et al., 2025) also discussed the emergence of sparse attention and the timing of training dynamics. Furthermore, (Yang et al., 2024) provides a clear and rigorous discussion of the two-stage training dynamics under classifiable text tasks. Similarly, (Chen & Luo, 2025) offers a more rigorous proof of dynamical separation in more general scenarios. Regarding multi-stage analysis, (Xu et al., 2025a) analyzes LoRA’s cross-stage dynamics, while (Wang & Ma, 2023) provides a full-process characterization of two-layer ReLU networks across four distinct training phases, from initialization to convergence. However, the aforementioned literature either relies on data settings that deviate significantly from real-world scenarios or requires overly stringent analytical conditions.

Transformers on Markov chains A significant body of influential work employs Markov chains to understand how Transformers, as probabilistic models, learn continuous linguistic data. (Chang et al., 2024) discovers that LLM learning can be summarized as “early n-gram learning followed by the gradual refinement of low-probability (tail) n-gram predictions.” (Bietti et al., 2023) analyzes the formation mechanism of induction heads using Markov-like data. (Rajaraman et al., 2024) investigates the impact of tokenization on Markovian data, proving that appropriate tokenization assists Transformers in modeling Markov processes. Additionally, (Makkuva et al., 2024) and (Makkuva et al., 2025) explore training dynamics and convergence analysis specifically under Markovian data settings.

Small initialization The initialization of a neural network significantly affects its learning outcomes (Arora et al., 2019; Williams et al., 2019; Mei et al., 2018; Jacot et al., 2018; Rotskoff & Vanden-Eijnden, 2018; Zhang et al., 2020). Small initialization is a common setting investigated in the study of neural network optimization dynamics, which contrasts with the Neural Tangent Kernel (NTK) perspective prevalent in infinitely wide networks. For linear models, (Ji & Telgarsky, 2019) theoretically establish results regarding matrix alignment. For nonlinear models, (Zhou et al., 2022) found that small initialization similarly promotes parameter condensation, thereby reducing model complexity. Theoretically, (Luo et al., 2021; Chen et al., 2024b; Zhou et al., 2023; Kumar & Haupt, 2024) have further deepened the understanding of this phenomenon. A recent survey article (Xu et al., 2025b) systematically synthesizes these empirical and theoretical findings.

B. Theoretical details in Sec. 2

B.1. Property of markov process

We will use two standard asymptotic properties of Markov chains. For the sake of completeness, we provide a detailed proof.

Proposition B.1 (Basic Markov properties). *Given the transition matrix P in (2) and any initial distribution μ_0 :*

1. **Convergence.** *The marginal distribution converges to π . That is $\lim_{t \rightarrow \infty} \mu_0^\top P^t = \pi^\top$.*
2. **Ergodicity.** *Along a single trajectory, the empirical state frequencies converge to π . That is $\lim_{s \rightarrow \infty} \frac{1}{s} \sum_{j=1}^s \mathbb{1}_{x_j} = \pi^\top$, where $\mathbb{1}_{x_j} \in \mathbb{R}^d$ is the one-hot vector of token x_j .*

Proof. Throughout, we work on the finite state space $\mathcal{V} = \{1, \dots, |\mathcal{V}|\}$. By the definition of the transition matrix P defined in (2), P is irreducible and aperiodic with strictly positive entries. Thus, P is ergodic. In particular, P admits a unique stationary distribution π satisfying $\pi^\top P = \pi^\top$ and $\pi_i > 0$ for all i .

1) Convergence of marginals. Since P is ergodic on a finite state space, it is primitive. By the Perron–Frobenius theorem, the eigenvalue 1 of P is simple and all other eigenvalues satisfy $|\lambda| < 1$. Let $\mathbf{1} \in \mathbb{R}^{|\mathcal{V}|}$ denote the all-ones vector. Because

P is row-stochastic, we have $P\mathbf{1} = \mathbf{1}$; because π is stationary, we have $\pi^\top P = \pi^\top$. Define the rank-one projector

$$\Pi := \mathbf{1} \pi^\top.$$

Then $\Pi^2 = \Pi$, $P\Pi = \Pi P = \Pi$, and we can write

$$P = \Pi + Q, \quad \text{where } Q := P - \Pi.$$

Note that $Q\mathbf{1} = 0$ and $\pi^\top Q = 0$. Moreover, the spectrum of Q equals the spectrum of P with the eigenvalue 1 removed, hence its spectral radius satisfies $\rho(Q) < 1$. Therefore, $Q^t \rightarrow 0$ as $t \rightarrow \infty$ (in any matrix norm), and

$$P^t = (\Pi + Q)^t = \Pi + Q^t \xrightarrow[t \rightarrow \infty]{} \Pi = \mathbf{1} \pi^\top. \quad (25)$$

For any initial distribution μ_0 (a row vector with nonnegative entries summing to 1),

$$\mu_0^\top P^t \xrightarrow[t \rightarrow \infty]{} \mu_0^\top (\mathbf{1} \pi^\top) = (\mu_0^\top \mathbf{1}) \pi^\top = \pi^\top,$$

which proves the convergence claim.

2) Ergodicity of empirical frequencies. Let $(X_t)_{t \geq 0}$ be the Markov chain with transition matrix P and arbitrary initial distribution μ_0 . Fix a reference state, say state 1, and define the (strict) return times

$$\tau_0 := 0, \quad \tau_{k+1} := \inf\{t > \tau_k : X_t = 1\}, \quad k \geq 0.$$

By irreducibility on a finite state space, the chain is positive recurrent, hence $\tau_k < \infty$ almost surely for all k and $\mathbb{E}_1[\tau_1] < \infty$.

For each cycle $k \geq 0$, define the cycle length and the state- i visit count within the cycle:

$$S_k := \tau_{k+1} - \tau_k, \quad R_k(i) := \sum_{t=\tau_k}^{\tau_{k+1}-1} \mathbb{1}\{X_t = i\}.$$

By the strong Markov property, conditional on $X_{\tau_k} = 1$ the post- τ_k evolution is independent of the past, and therefore the pairs $\{(S_k, R_k(i))\}_{k \geq 1}$ are i.i.d. under \mathbb{P}_1^{MC} (and also after the chain first hits state 1 when started from an arbitrary μ_0). Let $N(T) := \max\{k : \tau_k \leq T\}$ be the number of completed cycles up to time T . Then for each fixed i ,

$$\sum_{t=0}^{T-1} \mathbb{1}\{X_t = i\} = \underbrace{\sum_{t=0}^{\tau_1-1} \mathbb{1}\{X_t = i\}}_{\text{initial transient}} + \sum_{k=1}^{N(T)-1} R_k(i) + \underbrace{\sum_{t=\tau_{N(T)}}^{T-1} \mathbb{1}\{X_t = i\}}_{\text{remainder}}. \quad (26)$$

Divide by T . The initial transient term is $O(1/T)$ almost surely. The remainder term is at most one cycle, hence bounded by $S_{N(T)}$, and thus also negligible after dividing by T because $\tau_{N(T)} \leq T < \tau_{N(T)+1}$ implies $S_{N(T)} \leq \tau_{N(T)+1}$ and $\tau_{N(T)} \rightarrow \infty$.

It remains to analyze the dominant sum over complete cycles. By the strong law of large numbers applied to the i.i.d. sequences $\{S_k\}$ and $\{R_k(i)\}$,

$$\frac{1}{n} \sum_{k=1}^n S_k \rightarrow \mathbb{E}_1[S_1] = \mathbb{E}_1[\tau_1], \quad \frac{1}{n} \sum_{k=1}^n R_k(i) \rightarrow \mathbb{E}_1[R_1(i)] \quad \text{a.s.} \quad (27)$$

Moreover, by the definition of τ_k , we know

$$\tau_{N(T)+1} = \sum_{k=0}^{N(T)} S_k \geq T, \quad \tau_{N(T)} = \sum_{k=0}^{N(T)-1} S_k \leq T - 1$$

which implies

$$\frac{N(T)}{T} \rightarrow \frac{1}{\mathbb{E}_1[\tau_1]} \quad \text{a.s.}$$

Combining with (26) and (27), we obtain

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{1}\{X_t = i\} \xrightarrow[T \rightarrow \infty]{a.s.} \frac{\mathbb{E}_1[R_1(i)]}{\mathbb{E}_1[\tau_1]}. \quad (28)$$

Take $i = 1$, we get the right-hand side is π_1 by the computation about expectation of first return time. Since the above proof process is independent of the choice of the reference state, by considering all possible reference states, we obtain the result. \square

B.2. Gradient-flow dynamics

In this section, we will supplement the proof details of Proposition 2.3.

Proof. We first derive Eq. (4), using the standard trace theorem and chain rule. Taking the total differential of the loss, we get

$$d\mathcal{L} = \left\langle \frac{\partial \mathcal{L}}{\partial M}, dM \right\rangle + \left\langle \frac{\partial \mathcal{L}}{\partial \Phi}, d\Phi \right\rangle \quad (29)$$

Using the chain rule, we get

$$\begin{aligned} \left\langle \frac{\partial \mathcal{L}}{\partial M}, dM \right\rangle &= \left\langle \frac{\partial \mathcal{L}}{\partial M}, dW_0 W_1 + W_0 dW_1 \right\rangle \\ \left\langle \frac{\partial \mathcal{L}}{\partial \Phi}, d\Phi \right\rangle &= \left\langle \frac{\partial \mathcal{L}}{\partial \Phi}, dW_0 W_Q W_K^\top W_0^\top + W_0 dW_Q W_K^\top W_0^\top + W_0 W_Q dW_K^\top W_0^\top + W_0 W_Q W_K^\top dW_0^\top \right\rangle \end{aligned} \quad (30)$$

We derive the evolution equation for W_0 , and the other derivations are similar. We collect items related to dW_0 :

$$\begin{aligned} &\left\langle \frac{\partial \mathcal{L}}{\partial M}, dW_0 W_1 \right\rangle + \left\langle \frac{\partial \mathcal{L}}{\partial \Phi}, dW_0 W_Q W_K^\top W_0^\top + W_0 W_Q W_K^\top dW_0^\top \right\rangle \\ &= \text{tr} \left(\left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top dW_0 W_1 \right) + \text{tr} \left(\left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top (dW_0 W_Q W_K^\top W_0^\top + W_0 W_Q W_K^\top dW_0^\top) \right) \\ &= \text{tr} \left(W_1 \left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top dW_0 \right) + \text{tr} \left(W_Q W_K^\top W_0^\top \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top dW_0 \right) + \text{tr} \left(W_K W_Q^\top W_0^\top \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top dW_0 \right) \\ &= \left\langle \frac{\partial \mathcal{L}}{\partial M} W_1^\top + \frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K W_Q^\top + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 W_Q W_K^\top, dW_0 \right\rangle \end{aligned} \quad (31)$$

Therefore, we obtain an expression for $\frac{\partial \mathcal{L}}{\partial W_0}$. Since we are considering gradient descent, the evolution of the parameters follows the direction of the negative gradient. Then we derive the expression of $\frac{\partial \mathcal{L}}{\partial M}$ and $\frac{\partial \mathcal{L}}{\partial \Phi}$ in large N and s limit. Taking the total differential of the loss and taking the term about dM , we get

$$\left\langle \frac{\partial \mathcal{L}}{\partial M}, dM \right\rangle = \frac{1}{N} \sum_{i=1}^N - \sum_{l=1}^s A_{s,l}(X_i) e_{x,l}^\top dM e_{y_i} + \frac{1}{N} \sum_{i=1}^N \sum_j p_{y_j}(X_i) \sum_{l=1}^s A_{s,l}(X_i) e_{x,l}^\top dM e_{y_j} \quad (32)$$

Here, $A_{s,l}(X_i) = \frac{\exp(e_{x_{i,s}}^\top \Phi e_{x_{i,l}})}{\sum_{l'=1}^s \exp(e_{x_{i,s}}^\top \Phi e_{x_{i,l'}})}$. Based on Proposition B.1, we find that

$$A_{s,l}(X_i) = \frac{\frac{1}{s} \exp(e_{x_{i,s}}^\top \Phi e_{x_{i,l}})}{\frac{1}{s} \sum_{l'=1}^s \exp(e_{x_{i,s}}^\top \Phi e_{x_{i,l'}})} = \frac{\frac{1}{s} \exp(e_{x_{i,s}}^\top \Phi e_{x_{i,l}})}{\sum_{j=1}^d \pi_j \exp(e_{x_{i,s}}^\top \Phi e_j)} \quad (33)$$

Then, for sufficiently large sequence length s ,

$$\begin{aligned} \sum_{l=1}^s A_{s,l}(X_i) e_{x_{i,l}}^\top &= \frac{1}{\sum_{j=1}^d \pi_j \exp(e_{x_{i,s}}^\top \Phi e_j)} \sum_{l=1}^s \frac{1}{s} \exp(e_{x_{i,s}}^\top \Phi e_{x_{i,l}}) e_{x_{i,l}}^\top \\ &= \frac{1}{\sum_{j=1}^d \pi_j \exp(e_{x_{i,s}}^\top \Phi e_j)} \sum_{j'=1}^d \pi_{j'} \exp(e_{x_{i,s}}^\top \Phi e_{j'}) e_{j'}^\top = \mathbb{A}_{x_{i,s}} \end{aligned} \quad (34)$$

It implies that the output probability $p(X_i)$ actually depends on the last token $x_{i,s}$. That is

$$p_j(X_i) = \frac{\exp(\mathbb{A}_{x_{i,s}} M e_{y_j})}{\sum_{j'=1}^d \exp(\mathbb{A}_{x_{i,s}} M e_{y_{j'}})} = \mathbb{P}_{x_{i,s},j} \quad (35)$$

Based on this fact and Eq. (34), Eq. (32) can be reformulated by using the notations about \mathbb{A} and \mathbb{P} .

$$\left\langle \frac{\partial \mathcal{L}}{\partial M}, dM \right\rangle = -\frac{1}{N} \sum_{i=1}^N \mathbb{A}_{x_{i,s}} dM e_{y_i} + \frac{1}{N} \sum_{i=1}^N \mathbb{A}_{x_{i,s}} dM \mathbb{P}_{x_{i,s}}^\top. \quad (36)$$

Based on Proposition B.1, we have $x_{i,s} \sim \pi^\top$ when s is sufficiently large. Thus, we get

$$\left\langle \frac{\partial \mathcal{L}}{\partial M}, dM \right\rangle = -\sum_{i=1}^d \pi_i \mathbb{A}_i dM (P_i - \mathbb{P}_i)^\top \quad (37)$$

Using the trace theorem again, we have

$$\frac{\partial \mathcal{L}}{\partial M} = -\sum_{i=1}^d \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i). \quad (38)$$

Then we derive $\frac{\partial \mathcal{L}}{\partial \Phi}$. By direct computation, we get

$$\begin{aligned} \left\langle \frac{\partial \mathcal{L}}{\partial \Phi}, d\Phi \right\rangle &= \frac{1}{N} \sum_{i=1}^N - \left(\sum_{l=1}^s \left(A_{s,l} d(e_{x_{i,s}}^\top \Phi e_{x,l}) - A_{s,l} \sum_{l'} A_{s,l'} d(e_{x_{i,s}}^\top \Phi e_{x,l'}) \right) e_{x,l}^\top M e_{y_i} \right) \\ &+ \frac{1}{N} \sum_{i=1}^N \sum_j p_{y_j} \left(\sum_{l=1}^s \left(A_{s,l} d(e_{x_{i,s}}^\top \Phi e_{x,l}) - A_{s,l} \sum_{l'} A_{s,l'} d(e_{x_{i,s}}^\top \Phi e_{x,l'}) \right) e_{x,l}^\top M e_{y_j} \right) \end{aligned} \quad (39)$$

Using the notations we introduced, the derivative can be reformulated as

$$\begin{aligned} \left\langle \frac{\partial \mathcal{L}}{\partial \Phi}, d\Phi \right\rangle &= -\sum_i \pi_i \left(\sum_{l=1}^s A_{x_{i,l}} e_{x_i}^\top d\Phi (e_l - \mathbb{A}_i^\top) \right) e_l^\top M (P_i - \mathbb{P}_i)^\top \\ &= -\sum_i \pi_i e_{x_i}^\top d\Phi (\text{diag}(\mathbb{A}_i^\top) - \mathbb{A}_i^\top \mathbb{A}_i) M (P_i - \mathbb{P}_i)^\top \end{aligned} \quad (40)$$

As a result, using the trace theorem, we get the expression about $\frac{\partial \mathcal{L}}{\partial \Phi}$ as follows:

$$\frac{\partial \mathcal{L}}{\partial \Phi} = -\sum_i \pi_i e_{x_i} (P_i - \mathbb{P}_i) M^\top (\text{diag}(\mathbb{A}_i^\top) - \mathbb{A}_i^\top \mathbb{A}_i) \quad (41)$$

□

C. Theoretical details in Sec. 3

C.1. Theoretical details in Sec. 3.1

Proof of Lemma 3.1

Proof. Let $\Delta\theta(t) := \theta(t) - \theta_*$. Since $F(\theta_*) = 0$ and $J = DF(\theta_*)$, we can write

$$\dot{\Delta\theta}(t) = J\Delta\theta(t) + R(\Delta\theta(t)), \quad R(\Delta\theta) := F(\theta_* + \Delta\theta) - J\Delta\theta. \quad (42)$$

By assumption (6), for all $\|\Delta\theta\| \leq r$,

$$\|R(\Delta\theta)\| \leq L\|\Delta\theta\|^2. \quad (43)$$

Step 1: Variation-of-constants representation. Let $\tilde{\Delta}\theta(t) := e^{Jt}\Delta\theta(0)$ be the solution of the linearized system. From (42), the solution satisfies the Duhamel formula

$$\Delta\theta(t) = e^{Jt}\Delta\theta(0) + \int_0^t e^{J(t-s)}R(\Delta\theta(s))ds = \tilde{\Delta}\theta(t) + \int_0^t e^{J(t-s)}R(\Delta\theta(s))ds. \quad (44)$$

Define the linearization error $E(t) := \Delta\theta(t) - \tilde{\Delta}\theta(t)$. Then $E(0) = 0$ and by (44),

$$E(t) = \int_0^t e^{J(t-s)}R(\Delta\theta(s))ds. \quad (45)$$

Step 2: A standard bound on the semigroup e^{Jt} . Let $\mu := \sup\{\Re(\lambda) : \lambda \in \sigma(J)\}$. In finite dimension, for the chosen operator norm there exists a constant $K \geq 1$ such that

$$\|e^{Jt}\| \leq Ke^{\mu t}, \quad \forall t \geq 0. \quad (46)$$

Step 3: Bootstrap control inside the neighborhood. Fix a time horizon $T > 0$ such that

$$\|\tilde{\Delta}\theta(t)\| \leq \frac{r}{2}, \quad \forall t \in [0, T]. \quad (47)$$

We will show that for $\varepsilon := \|\Delta\theta(0)\|$ sufficiently small (depending on J, L, r), the trajectory stays in the ball $\|\Delta\theta(t)\| \leq r$ on $[0, T]$, so that (43) applies.

Indeed, from (44), (46), and (43), as long as $\|\Delta\theta(s)\| \leq r$ for all $s \in [0, t]$, we have

$$\|E(t)\| \leq \int_0^t \|e^{J(t-s)}\| \|R(\Delta\theta(s))\| ds \leq KL \int_0^t e^{\mu(t-s)} \|\Delta\theta(s)\|^2 ds. \quad (48)$$

Also $\|\tilde{\Delta}\theta(t)\| \leq \|e^{Jt}\| \|\Delta\theta(0)\| \leq K\varepsilon e^{\mu t}$.

We now bootstrap the bound

$$\|\Delta\theta(t)\| \leq 2\|\tilde{\Delta}\theta(t)\| \quad \text{for all } t \in [0, T]. \quad (49)$$

Assuming (49) holds on $[0, t]$, then $\|\Delta\theta(s)\| \leq 2\|\tilde{\Delta}\theta(s)\| \leq r$ by (47), so (48) applies and yields

$$\begin{aligned} \|E(t)\| &\leq KL \int_0^t e^{\mu(t-s)} (2\|\tilde{\Delta}\theta(s)\|)^2 ds = 4KL \int_0^t e^{\mu(t-s)} \|\tilde{\Delta}\theta(s)\|^2 ds \\ &\leq 4KL \int_0^t e^{\mu(t-s)} (K\varepsilon e^{\mu s})^2 ds = 4K^3L \varepsilon^2 e^{\mu t} \int_0^t e^{\mu s} ds. \end{aligned} \quad (50)$$

If $\mu > 0$, then $\int_0^t e^{\mu s} ds = (e^{\mu t} - 1)/\mu \leq e^{\mu t}/\mu$, and thus

$$\|E(t)\| \leq \frac{4K^3L}{\mu} \varepsilon^2 e^{2\mu t}. \quad (51)$$

If $\mu = 0$, then $\int_0^t e^{\mu s} ds = t$ and (50) gives

$$\|E(t)\| \leq 4K^3L \varepsilon^2 t. \quad (52)$$

(For $\mu < 0$, one may similarly bound the integral by a constant and obtain a uniform $O(\varepsilon^2)$ error.)

Now choose ε small enough such that on $[0, T]$,

$$\|E(t)\| \leq \|\tilde{\Delta}\theta(t)\| \quad \forall t \in [0, T]. \quad (53)$$

This is possible because by (47) we have $\|\tilde{\Delta}\theta(t)\| \leq r/2$, while (51) shows $\|E(t)\|$ is $O(\varepsilon^2)$ times an exponential factor; e.g. it suffices to require

$$\frac{4K^3L}{\mu} \varepsilon e^{\mu T} \leq \frac{1}{2} \quad (\mu > 0), \quad \text{or} \quad 4K^3L \varepsilon T \leq \frac{1}{2} \quad (\mu = 0),$$

and recall T is such that $\|\tilde{\Delta}\theta(t)\| \leq r/2$ on $[0, T]$, hence $e^{\mu T}$ is at most on the order of $1/\varepsilon$ when $\mu > 0$. Under (53),

$$\|\Delta\theta(t)\| \leq \|\tilde{\Delta}\theta(t)\| + \|E(t)\| \leq 2\|\tilde{\Delta}\theta(t)\| \leq r,$$

so the bootstrap is self-consistent and (51) (or (52)) holds for all $t \in [0, T]$. This proves (7) with $C = \frac{4K^3L}{\mu}$ for $\mu > 0$ and $C = 4K^3L$ for $\mu = 0$.

Step 4: The $\Theta(\log(1/\varepsilon))$ window when $\mu > 0$. If $\mu > 0$, then $\|\tilde{\Delta}\theta(t)\| \leq K\varepsilon e^{\mu t}$. Therefore the condition $\|\tilde{\Delta}\theta(t)\| \leq r/2$ holds at least up to times

$$t \leq \frac{1}{\mu} \log \frac{r}{2K\varepsilon} = \Theta(\log(1/\varepsilon)),$$

which is exactly the linearization window claimed in the lemma.

Step 5: Alignment with the unstable eigenvector with positive spectral gap. Assume now that J has a simple eigenvalue $\mu > 0$ with eigenvector v_u and a spectral gap: $\Re(\lambda) \leq \mu - \delta$ for all other eigenvalues. Let Π_u be the spectral projection onto $\text{span}\{v_u\}$ and $\Pi_s = I - \Pi_u$. Then there exist constants K_u, K_s such that

$$\|\Pi_u e^{Jt}\| \leq K_u e^{\mu t}, \quad \|\Pi_s e^{Jt}\| \leq K_s e^{(\mu-\delta)t}, \quad \forall t \geq 0. \quad (54)$$

Write $\tilde{\Delta}\theta(t) = \Pi_u \tilde{\Delta}\theta(t) + \Pi_s \tilde{\Delta}\theta(t)$. If $\langle \Delta(0), v_u \rangle \neq 0$, then $\Pi_u \tilde{\Delta}\theta(t) = a_0 e^{\mu t} v_u$ for some $a_0 \neq 0$, while

$$\|\Pi_s \tilde{\Delta}\theta(t)\| \leq K_s \varepsilon e^{(\mu-\delta)t}.$$

Hence

$$\frac{\tilde{\Delta}\theta(t)}{\|\tilde{\Delta}\theta(t)\|} \rightarrow \pm \frac{v_u}{\|v_u\|} \quad \text{as } t \rightarrow \infty, \quad (55)$$

and the convergence rate is $O(e^{-\delta t})$.

For the nonlinear trajectory, decompose $\Delta\theta(t) = u(t) + s(t)$ with $u(t) := \Pi_u \Delta\theta(t)$ and $s(t) := \Pi_s \Delta\theta(t)$. Projecting (44) onto the two subspaces and using (54) gives

$$u(t) = \Pi_u \tilde{\Delta}\theta(t) + \int_0^t \Pi_u e^{J(t-s)} R(\Delta\theta(s)) ds, \quad s(t) = \Pi_s \tilde{\Delta}\theta(t) + \int_0^t \Pi_s e^{J(t-s)} R(\Delta\theta(s)) ds. \quad (56)$$

Inside the linearization window we have $\|\Delta\theta(s)\| \leq r$, so (43) applies and, using also $\|\Delta\theta(s)\| \lesssim \varepsilon e^{\mu s}$ from the bootstrap in Step 3, we obtain

$$\|R(\Delta(s))\| \leq L \|\Delta\theta(s)\|^2 \lesssim L \varepsilon^2 e^{2\mu s}.$$

Plugging into (56) yields, for t in the linearization window,

$$\|u(t) - \Pi_u \tilde{\Delta}\theta(t)\| \leq \int_0^t \|\Pi_u e^{J(t-s)}\| \|R(\Delta\theta(s))\| ds \lesssim \varepsilon^2 e^{2\mu t}, \quad (57)$$

$$\|s(t) - \Pi_s \tilde{\Delta}\theta(t)\| \leq \int_0^t \|\Pi_s e^{J(t-s)}\| \|R(\Delta\theta(s))\| ds \lesssim \varepsilon^2 e^{2\mu t}. \quad (58)$$

Therefore,

$$\|s(t)\| \leq \|\Pi_s \tilde{\Delta}\theta(t)\| + \|s(t) - \Pi_s \tilde{\Delta}\theta(t)\| \lesssim \varepsilon e^{(\mu-\delta)t} + \varepsilon^2 e^{2\mu t},$$

while

$$\|u(t)\| \geq \|\Pi_u \tilde{\Delta}\theta(t)\| - \|u(t) - \Pi_u \tilde{\Delta}\theta(t)\| \gtrsim \varepsilon e^{\mu t} - \varepsilon^2 e^{2\mu t}.$$

Hence, for times t such that $\varepsilon e^{\mu t}$ is still sufficiently small (which holds throughout a $\Theta(\log(1/\varepsilon))$ interval inside the linearization window), we have

$$\frac{\|s(t)\|}{\|u(t)\|} \lesssim e^{-\delta t} + \varepsilon e^{\mu t}. \quad (59)$$

Now let t increase while remaining in the linearization window (so $t \rightarrow \infty$ is possible as $\varepsilon \rightarrow 0$), and choose any sequence $t = t(\varepsilon)$ such that

$$t(\varepsilon) \rightarrow \infty, \quad \varepsilon e^{\mu t(\varepsilon)} \rightarrow 0 \quad (\text{e.g. } t(\varepsilon) = \frac{1}{2\mu} \log(1/\varepsilon)).$$

Then (59) implies $\|s(t)\|/\|u(t)\| \rightarrow 0$, so

$$\frac{\Delta\theta(t)}{\|\Delta\theta(t)\|} = \frac{u(t) + s(t)}{\|u(t) + s(t)\|} \rightarrow \pm \frac{v_u}{\|v_u\|}.$$

This proves (8). \square

Proof of Theorem 3.2

Proof. The claim follows by a direct evaluation of the gradients at the origin. By Lemma 3.1, the early-time dynamics is governed by the linearization at $\theta = 0$, so we substitute $\theta = 0$ into (5).

At $\theta = 0$, the definitions of \mathbb{A}_i and \mathbb{P}_i yield

$$\mathbb{A}_i = \pi^\top, \quad \mathbb{P}_i = \frac{1}{d} \mathbf{1}^\top, \quad \text{for all } i.$$

Moreover, since π^\top is stationary, we have $\pi^\top P = \pi^\top$, and hence

$$\sum_i \pi_i P_i = \pi^\top.$$

Plugging these identities into the expression of $\frac{\partial \mathcal{L}}{\partial M}$ in (5), we obtain

$$\left. \frac{\partial \mathcal{L}}{\partial M} \right|_{\theta=0} = -\pi \left(\pi - \frac{1}{d} \mathbf{1} \right)^\top.$$

This is exactly the desired formula, completing the proof. \square

C.2. Theoretical details in Sec. 3.2

Proof of Proposition 3.3

Proof. Since attention parameters (W_Q, W_K) are chosen to be zero, for any i we have $\mathbb{A}_i = \pi^\top$, and hence $\mathbb{P}_i = \mathbb{P}_j$ for $i \neq j$. By definition,

$$\mathbb{P}_{i,j} = \frac{\exp\left(\kappa^2 \|\pi\|^2 \left(\pi_j - \frac{1}{|\mathcal{V}|}\right)\right)}{\sum_{j'} \exp\left(\kappa^2 \|\pi\|^2 \left(\pi_{j'} - \frac{1}{|\mathcal{V}|}\right)\right)} = \frac{\exp\left(\kappa^2 \|\pi\|^2 \pi_j\right)}{\sum_{j'} \exp\left(\kappa^2 \|\pi\|^2 \pi_{j'}\right)}. \quad (60)$$

As $\kappa \rightarrow 0$, $\mathbb{P}_i \rightarrow \frac{1}{|\mathcal{V}|} \mathbf{1}^\top$; as $\kappa \rightarrow \infty$, $\mathbb{P}_i \rightarrow e_1^\top$ since $\pi_1 = \max_j \pi_j$. The map $\kappa \mapsto \mathbb{P}_{i,1}$ is continuous, hence by the intermediate value theorem there exists $\kappa_1 > 0$ such that $\mathbb{P}_{i,1} = \pi_1$. Together with the symmetry assumption $\pi_i = \pi_j$ for $i, j \geq 2$, this implies $\mathbb{P}_i = \pi^\top$.

Substituting $\mathbb{P}_i = \pi^\top$ and $\mathbb{A}_i = \pi^\top$ into $\frac{\partial \mathcal{L}}{\partial M} = -\sum_{i=1}^{|\mathcal{V}|} \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i)$ yields

$$\frac{\partial \mathcal{L}}{\partial M} = -\pi \sum_{i=1}^{|\mathcal{V}|} \pi_i (P_i - \pi^\top). \quad (61)$$

Using $\pi^\top P = \pi^\top$, we obtain $\frac{\partial \mathcal{L}}{\partial M} = 0$. Finally, $W_Q = W_K = 0$ at this point, so it is indeed a critical point. \square

We first record the derivatives needed for the linearization.

Proposition C.1 (Derivatives at the second critical point). *At θ_c^1 in (12), we have*

$$\left. \frac{\partial \mathcal{L}}{\partial \Phi} \right|_{\theta_c^1} = -\lambda \kappa_1^2 \text{Var}(\pi) \frac{\pi - \frac{1}{|\mathcal{V}|} \mathbf{1}}{\left\| \pi - \frac{1}{|\mathcal{V}|} \mathbf{1} \right\|} \frac{\pi^\top}{\|\pi\|} \text{Var}(\pi), \quad (62)$$

and the total differential of $\frac{\partial \mathcal{L}}{\partial M}$ satisfies

$$d \left. \frac{\partial \mathcal{L}}{\partial M} \right|_{\theta_c^1} = \pi \pi^\top dM \left(\text{diag}(\pi) - \pi \pi^\top \right). \quad (63)$$

Proof of Proposition C.1

Proof. First, we compute the specific expression of $\frac{\partial \mathcal{L}}{\partial \Phi}$. Using the expression derived in Eq. (5),

$$\frac{\partial \mathcal{L}}{\partial \Phi} = - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top (\text{diag}(\mathbb{A}_i^\top) - \mathbb{A}_i^\top \mathbb{A}_i)$$

Using the fact that $\mathbb{P}_i = \pi^\top$ and $\mathbb{A}_i = \pi^\top$ and substituting $M = \kappa_1^2 \frac{\pi}{\|\pi\|} \frac{\pi^\top - \frac{1}{d} \mathbf{1}^\top}{\|\pi - \frac{1}{d} \mathbf{1}\|}$ into the equation, we get

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \Phi} &= -\kappa_1^2 \sum_i \pi_i e_i (P_i - \pi^\top) \frac{\pi - \frac{1}{d} \mathbf{1}}{\|\pi - \frac{1}{d} \mathbf{1}\|} \frac{\pi^\top}{\|\pi\|} (\text{diag}(\pi) - \pi \pi^\top) \\ &= -\kappa_1^2 (\text{diag}(\pi) P - \pi \pi^\top) \frac{\pi - \frac{1}{d} \mathbf{1}}{\|\pi - \frac{1}{d} \mathbf{1}\|} \frac{\pi^\top}{\|\pi\|} (\text{diag}(\pi) - \pi \pi^\top) \end{aligned} \quad (64)$$

By the definition of P , we have

$$\begin{aligned} \text{diag}(\pi) P - \pi \pi^\top &= \text{diag}(\pi) (\lambda I + (1 - \lambda) \mathbf{1} \mathbf{1}^\top) - \pi \pi^\top \\ &= \lambda (\text{diag}(\pi) - \pi \pi^\top) \end{aligned} \quad (65)$$

Combining Eqs (64) and (65), we get the expression of $\frac{\partial \mathcal{L}}{\partial \Phi}$ at θ_c^1 . Then, we consider the total differential of the gradient of the loss function with respect to M . Firstly, using Eq. (5) again, we get

$$\frac{\partial \mathcal{L}}{\partial M} = - \sum_i \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i)$$

By chain rule, we get

$$d \frac{\partial \mathcal{L}}{\partial M} = - \sum_i \pi_i d \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top (-d \mathbb{P}_i). \quad (66)$$

Let's consider these two items separately. By the definition of \mathbb{A}_i , we find that

$$\begin{aligned} d \mathbb{A}_{i,j} &= \mathbb{A}_{i,j} e_i^\top d \Phi e_j - \mathbb{A}_{i,j} \sum_{j'} \mathbb{A}_{i,j'} e_i^\top d \Phi e_{j'} \\ &= \mathbb{A}_{i,j} e_i^\top d \Phi (e_j - \mathbb{A}_i^\top) \end{aligned} \quad (67)$$

As a result, it can be verified that $d \mathbb{A}_i = e_i^\top d \Phi (\text{diag}(\mathbb{A}_i^\top) - \mathbb{A}_i^\top \mathbb{A}_i)$. However, this term will be zero because it contains the intersection terms $W_Q W_K^\top$ which will be zero by the chain rule and the condition $W_Q = 0$ and $W_K = 0$. Thus, we focus on the second term. Recall the definition of $\mathbb{P}_{i,j} = \frac{\exp(\mathbb{A}_i M e_j)}{\sum_{j'} \exp(\mathbb{A}_i M e_{j'})}$ and take the total differential of it:

$$\begin{aligned} d \mathbb{P}_{i,j} &= \mathbb{P}_{i,j} d(\mathbb{A}_i M) e_j - \mathbb{P}_{i,j} \sum_{j'} \mathbb{P}_{i,j'} d(\mathbb{A}_i M) e_{j'} \\ &= \mathbb{P}_{i,j} d(\mathbb{A}_i M) (e_j - \mathbb{P}_i^\top). \end{aligned} \quad (68)$$

Similar to the derivation of $d \mathbb{A}_i$, $d \mathbb{P}_i$ can be reformulated as $d(\mathbb{A}_i M) (\text{diag}(\mathbb{P}_i^\top) - \mathbb{P}_i^\top \mathbb{P}_i)$. In particular, at this critical point,

$$d \mathbb{P}_i = \mathbb{A}_i dM (\text{diag}(\mathbb{P}_i^\top) - \mathbb{P}_i^\top \mathbb{P}_i) = \pi^\top dM \text{Var}(\pi). \quad (69)$$

Substitute this expression into Eq. (66), we get

$$d \frac{\partial \mathcal{L}}{\partial M} = \sum_i \pi_i \mathbb{A}_i^\top \pi^\top dM \text{Var}(\pi) = \pi \pi^\top dM \text{Var}(\pi). \quad (70)$$

□

Proof of Proposition 3.4

Proof. We linearize the gradient flow (4) at θ_c^1 . Since $\frac{\partial \mathcal{L}}{\partial M} \Big|_{\theta_c^1} = 0$ and $W_Q = W_K = 0$, the only first-order contribution in the (W_0, W_1) subsystem comes from the first variation of $\frac{\partial \mathcal{L}}{\partial \Phi}$, whereas the (W_Q, W_K) subsystem is driven by the constant matrix $\frac{\partial \mathcal{L}}{\partial \Phi} \Big|_{\theta_c^1}$. More specifically, the linearized subsystems with respect to (W_0, W_1) and (W_Q, W_K) are two decoupled systems which separately follow

$$\begin{aligned} \frac{d\Delta W_0}{dt} &= d \frac{\partial \mathcal{L}}{\partial M} W_1^\top \\ \frac{d\Delta W_1}{dt} &= W_0^\top d \frac{\partial \mathcal{L}}{\partial M} \end{aligned} \quad (71)$$

and

$$\begin{aligned} \frac{d\Delta W_Q}{dt} &= -W_0^\top \frac{\partial \mathcal{L}}{\partial \Phi} W_0 \Delta W_K \\ \frac{d\Delta W_K}{dt} &= -W_0^\top \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 \Delta W_Q \end{aligned} \quad (72)$$

Step 1: the (W_0, W_1) -subsystem is contracting along α_1 . By Proposition C.1,

$$d \left(\frac{\partial \mathcal{L}}{\partial M} \right) \Big|_{\theta_c^1} = \pi \pi^\top dM \text{Var}(\pi)$$

At θ_c^1 , Proposition 3.3 gives the rank-one form

$$W_{0,0} = \kappa_1 q \alpha_1^\top, \quad W_{1,0} = \kappa_1 \alpha_1 u^\top, \quad q := \frac{\pi}{\|\pi\|}, \quad u := \frac{\pi - \frac{1}{|\mathcal{V}|} \mathbf{1}}{\left\| \pi - \frac{1}{|\mathcal{V}|} \mathbf{1} \right\|}. \quad (73)$$

A direct substitution into Eq. (71) shows that for any $v \perp \alpha_1$,

$$\frac{d}{dt} \Delta W_0 v = 0, \quad \frac{d}{dt} v^\top \Delta W_1 = 0,$$

i.e. the linearization is degenerate in the normal directions.

Therefore we focus on the α_1 -component and calculate the specific expansion:

$$\begin{aligned} \frac{d\Delta W_0 \alpha_1}{dt} &= -\pi \pi^\top (\Delta W_0 W_1 + W_0 \Delta W_1) \text{Var}(\pi) W_1^\top \\ &= -\kappa_1^2 \pi \pi^\top \Delta W_0 \alpha_1 \frac{\pi^\top - \frac{1}{d} \mathbf{1}^\top}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \text{Var}(\pi) \frac{\pi - \frac{1}{d} \mathbf{1}}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \\ &\quad - \kappa_1^2 \pi \pi^\top \frac{\pi}{\|\pi\|} \alpha_1^\top \Delta W_1 \text{Var}(\pi) \frac{\pi - \frac{1}{d} \mathbf{1}}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \end{aligned} \quad (74)$$

and

$$\begin{aligned} \frac{d\Delta W_1^\top \alpha_1}{dt} &= -\text{Var}(\pi) (W_1^\top \Delta W_0^\top + \Delta W_1^\top W_0^\top) \pi \pi^\top W_0 \\ &= -\kappa_1^2 \text{Var}(\pi) \frac{\pi - \frac{1}{d} \mathbf{1}}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \alpha_1^\top \Delta W_0^\top \pi \pi^\top \frac{\pi}{\|\pi\|} \\ &\quad - \kappa_1^2 \text{Var}(\pi) \Delta W_1^\top \alpha_1 \frac{\pi^\top}{\|\pi\|} \pi \pi^\top \frac{\pi}{\|\pi\|} \end{aligned} \quad (75)$$

Finally, we have

$$d \begin{pmatrix} \Delta W_0 \alpha_1 \\ \Delta W_1^\top \alpha_1 \end{pmatrix} = -\kappa_1^2 \begin{pmatrix} \frac{\pi^\top - \frac{1}{d} \mathbf{1}^\top}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \text{Var}(\pi) \frac{\pi - \frac{1}{d} \mathbf{1}}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \pi \pi^\top & \pi^\top \frac{\pi}{\|\pi\|} \pi \frac{\pi^\top - \frac{1}{d} \mathbf{1}^\top}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \text{Var}(\pi) \\ \pi^\top \frac{\pi}{\|\pi\|} \text{Var}(\pi) \frac{\pi - \frac{1}{d} \mathbf{1}}{\left\| \pi - \frac{1}{d} \mathbf{1} \right\|} \pi^\top & \frac{\pi^\top}{\|\pi\|} \pi \pi^\top \frac{\pi}{\|\pi\|} \text{Var}(\pi) \end{pmatrix} \begin{pmatrix} \Delta W_0 \alpha_1 \\ \Delta W_1^\top \alpha_1 \end{pmatrix} \quad (76)$$

We introduce the notations:

$$x := \Delta W_0 \alpha_1 \in \mathbb{R}^{|\mathcal{V}|}, \quad y := \Delta W_1^\top \alpha_1 \in \mathbb{R}^{|\mathcal{V}|}.$$

Eq. (76) can be rewritten in the following concise form

$$\frac{d}{dt} \begin{pmatrix} x \\ y \end{pmatrix} = -\kappa_1^2 A \begin{pmatrix} x \\ y \end{pmatrix}, \quad (77)$$

where

$$A := \begin{pmatrix} a \pi \pi^\top & b \pi u^\top C \\ b C u \pi^\top & b^2 C \end{pmatrix}, \quad a := u^\top C u, \quad b := \pi^\top q = \|\pi\|. \quad (78)$$

The matrix A is symmetric by construction. Moreover, for arbitrary x, y define $\alpha := \pi^\top x$ and $\beta := u^\top C y$. Then the quadratic form is

$$\begin{pmatrix} x^\top & y^\top \end{pmatrix} A \begin{pmatrix} x \\ y \end{pmatrix} = a \alpha^2 + 2b \alpha \beta + b^2 y^\top C y.$$

Introducing the C -inner product $\langle v, w \rangle_C := v^\top C w$ (with seminorm $\|v\|_C = \sqrt{v^\top C v}$), we have $a = \|u\|_C^2 \geq 0$ and $|\beta| = |\langle u, y \rangle_C| \leq \|u\|_C \|y\|_C = \sqrt{a} \sqrt{y^\top C y}$. Hence

$$a \alpha^2 + 2b \alpha \beta + b^2 y^\top C y \geq (\sqrt{a} |\alpha| - b \sqrt{y^\top C y})^2 \geq 0,$$

so $A \succeq 0$. Therefore all eigenvalues of $-\lambda^2 A$ in (77) are non-positive, and the (x, y) -subsystem is contracting (or neutrally stable in the degenerate directions).

Step 2: effective coupling for (W_Q, W_K) . Recall Eq. (72),

$$\frac{d}{dt} \Delta W_Q = -W_0^\top \left. \frac{\partial \mathcal{L}}{\partial \Phi} \right|_{\theta_0} W_0 \Delta W_K, \quad \frac{d}{dt} \Delta W_K = -W_0^\top \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \Big|_{\theta_0} W_0 \Delta W_Q$$

Substitute the expression of $\frac{\partial \mathcal{L}}{\partial \Phi}$ into above equation, we take $\frac{d \Delta W_Q}{dt}$ as an example:

$$\begin{aligned} \frac{d \Delta W_Q}{dt} &= \kappa_1^2 \alpha_1 \frac{\pi^\top}{\|\pi\|} \lambda \kappa_1^2 \left(\text{diag}(\pi) - \pi \pi^\top \right) \frac{\pi - \frac{1}{|\mathcal{V}|} \mathbf{1}}{\|\pi - \frac{1}{|\mathcal{V}|} \mathbf{1}\|} \frac{\pi^\top}{\|\pi\|} \left(\text{diag}(\pi) - \pi \pi^\top \right) \frac{\pi}{\|\pi\|} \alpha_1^\top \Delta W_K \\ &= c_1 \alpha_1 \alpha_1^\top \Delta W_K. \end{aligned} \quad (79)$$

Left-multiplying by α_1^\top yields the results. Moreover, c_1 is positive by its definition. Thus, it is an unstable direction. It implies that the effective dynamics near the critical point is the subsystem about W_Q and W_K

□

C.3. Theoretical details in Sec. 3.3

Proof of Proposition 3.6 We complete the proof of Proposition 1 in two steps. First, we directly verify that a rank-one manifold is an invariant manifold. Then, we utilize data symmetry and permutation equivariance to prove the conservation of low-frequency tokens.

Proof. (i) Invariance of the rank-one form. Plug (18) into (4) and check that each right-hand side remains in the same rank-one span.

Since $W_1^\top = \beta \alpha_1^\top$,

$$-\frac{\partial \mathcal{L}}{\partial M} W_1^\top = -\left(\frac{\partial \mathcal{L}}{\partial M} \beta \right) \alpha_1^\top,$$

which is of the form $\dot{\gamma} \alpha_1^\top$.

Next, using $\tilde{\alpha}_1^\top \tilde{\alpha}_1 = 1$,

$$W_K W_Q^\top = \lambda_K \lambda_Q \alpha_1 (\tilde{\alpha}_1^\top \tilde{\alpha}_1) \alpha_1^\top = \eta \alpha_1 \alpha_1^\top, \quad W_0 W_K W_Q^\top = \eta \gamma (\alpha_1^\top \alpha_1) \alpha_1^\top = \eta \gamma \alpha_1^\top.$$

Therefore the Φ -driven terms in \dot{W}_0 satisfy

$$-\frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K W_Q^\top = -\eta \left(\frac{\partial \mathcal{L}}{\partial \Phi} \gamma \right) \alpha_1^\top, \quad -\left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 W_Q W_K^\top = -\eta \left(\left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \gamma \right) \alpha_1^\top,$$

so \dot{W}_0 stays in the span of $\{\cdot \alpha_1^\top\}$ and hence $W_0(t) = \gamma(t) \alpha_1^\top$.

Similarly, since $W_0^\top = \alpha_1 \gamma^\top$,

$$\dot{W}_1 = -W_0^\top \frac{\partial \mathcal{L}}{\partial M} = -\alpha_1 \left(\gamma^\top \frac{\partial \mathcal{L}}{\partial M} \right),$$

which is of the form $\alpha_1 \tilde{\beta}^\top$.

Finally,

$$\dot{W}_Q = -W_0^\top \frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K = -\lambda_K \left(\gamma^\top \frac{\partial \mathcal{L}}{\partial \Phi} \gamma \right) \alpha_1 \tilde{\alpha}_1^\top,$$

so W_Q remains in the form $\lambda_Q(t) \alpha_1 \tilde{\alpha}_1^\top$. The argument for W_K is identical. Thus the flow stays in \mathcal{W} .

(ii) Preservation of the low-frequency symmetry. Let

$$G := \{\sigma : \{1, \dots, V\} \rightarrow \{1, \dots, V\} \mid \sigma(1) = 1\}. \quad (80)$$

and let $\sigma \in G$ be the permutation matrix. Define the group action as

$$\rho_\sigma(\theta) := (\Pi_\sigma W_0, W_1 \Pi_\sigma^\top, W_Q, W_K).$$

Under this action, ones check that $\mathbb{P}_{i,j}(\rho_\sigma(\theta)) = \mathbb{P}_{\sigma(i),\sigma(j)}(\theta)$. Since the loss function can be viewed as $\mathcal{L}(\theta) = -\sum_i \pi_i \sum_j P_{i,j} \log \mathbb{P}_{i,j}$, we find

$$\mathcal{L}(\rho_\sigma(\theta)) = -\sum_i \pi_i \sum_j P_{i,j} \log \mathbb{P}_{\sigma(i),\sigma(j)} = -\sum_i \pi_{\sigma^{-1}(i)} \sum_j P_{\sigma^{-1}(i),\sigma^{-1}(j)} \log \mathbb{P}_{i,j}. \quad (81)$$

Under the symmetry assumption on the data and the definition of the transition probability matrix P , $\mathcal{L}(\rho_\sigma(\theta)) = \mathcal{L}(\theta)$. Hence if $\theta(t)$ solves the gradient flow, so does $\rho_\sigma(\theta(t))$. If $\theta(t_0) = \rho_\sigma(\theta(t_0))$ for all $\sigma \in G$ which is equivalent to $\gamma_2 = \dots = \gamma_d$ and $\beta_2 = \dots = \beta_d$ at t_0 , uniqueness of ODE solutions implies $\theta(t) = \rho_\sigma(\theta(t))$ for all $t \geq t_0$, which proves the symmetry is preserved. \square

Proof of Theorem 3.7 We proceed with the proof of Theorem 3.7. First, we introduce some notation to show that the dynamics on a rank-one manifold will be further simplified in the case of low-frequency symmetry. Next, since we are still near the critical point described in Proposition 1, this means that we are also near the critical point for the dynamics on a rank-one manifold. Therefore, we continue using linearization methods to obtain the key conservation law results.

First, we find that the proxy attention matrix \mathbb{A} has the form on \mathcal{W} by direct computation,

$$\mathbb{A} = \begin{pmatrix} \xi_1 & \frac{1-\xi_1}{|\mathcal{V}|-1} & \dots & \frac{1-\xi_1}{|\mathcal{V}|-1} \\ \xi_2 & \frac{1-\xi_2}{|\mathcal{V}|-1} & \dots & \frac{1-\xi_2}{|\mathcal{V}|-1} \\ \vdots & \vdots & \dots & \vdots \\ \xi_2 & \frac{1-\xi_2}{|\mathcal{V}|-1} & \dots & \frac{1-\xi_2}{|\mathcal{V}|-1} \end{pmatrix},$$

where

$$\xi_1 = \frac{\pi_1 \exp(\eta \gamma_1^2)}{\pi_1 \exp(\eta \gamma_1^2) + (1 - \pi_1) \exp(\eta \gamma_1 \gamma_{i \neq 1})}, \quad (82)$$

$$\xi_2 = \frac{\pi_1 \exp(\eta \gamma_1 \gamma_{i \neq 1})}{\pi_1 \exp(\eta \gamma_1 \gamma_{i \neq 1}) + (1 - \pi_1) \exp(\eta \gamma_{i \neq 1}^2)}. \quad (83)$$

Define the row-wise scalar projections

$$m_1 := \mathbb{A}_1 \gamma, \quad m_2 := \mathbb{A}_2 \gamma.$$

Then

$$m_1 = \gamma_{i \neq 1} + \xi_1 \Delta \gamma, \quad m_2 = \gamma_{i \neq 1} + \xi_2 \Delta \gamma.$$

Since $\mathbb{A}_i M = (\mathbb{A}_i \gamma) \beta^\top = m_i \beta^\top$, the model probability of predicting the first token is

$$\hat{p}_i := \mathbb{P}_{i,1} = \frac{\exp(m_i \beta_1)}{\exp(m_i \beta_1) + (|\mathcal{V}| - 1) \exp(m_i \beta_{i \neq 1})} = \sigma(m_i \Delta \beta - \log(|\mathcal{V}| - 1)), \quad i \in \{1, 2\}, \quad (84)$$

where σ is the sigmoid function. Here, we only consider the first and second probability because $\mathbb{P}_{i,1} = \mathbb{P}_{j,1}$ for $i, j \neq 1$. Moreover, there exists a key term $(P_i - \mathbb{P}_i) \beta$ in the following computation. By direct computation,

$$\begin{aligned} (P_i - \mathbb{P}_i) \beta &= (P_{i,1} - \mathbb{P}_{i,1}) \beta_1 + ((1 - P_{i,1} - (1 - \mathbb{P}_{i,1}))) \beta_{i \neq 1} \\ &= (P_{i,1} - \mathbb{P}_{i,1}) \Delta \beta. \end{aligned} \quad (85)$$

It implies that $(P_i - \mathbb{P}_i) \beta = (P_j - \mathbb{P}_j) \beta$ for $i, j \neq 1$. Let the residuals be

$$r_i := P_{i,1} - \hat{p}_i, \quad i \in \{1, 2\}. \quad (86)$$

For $i > 2$, we let $r_i = r_2$.

We now derive the explicit dynamics for γ_1 and $\gamma_{i \neq 1}$ by expanding the two contributions in $\dot{\gamma}$ in (20).

(1). The M -driven term $-\frac{\partial \mathcal{L}}{\partial M} \beta$. By the definition of $\frac{\partial \mathcal{L}}{\partial M}$, we obtain

$$-\frac{\partial \mathcal{L}}{\partial M} \beta = \sum_i \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) \beta = \Delta \beta \left(\pi_1 r_1 \mathbb{A}_1^\top + (1 - \pi_1) r_2 \mathbb{A}_2^\top \right).$$

Taking the first coordinate and a generic low-token coordinate yields

$$\dot{\gamma}_1|_M = \Delta \beta [\pi_1 r_1 \xi_1 + (1 - \pi_1) r_2 \xi_2], \quad \dot{\gamma}_{i \neq 1}|_M = \frac{\Delta \beta}{|\mathcal{V}| - 1} [\pi_1 r_1 (1 - \xi_1) + (1 - \pi_1) r_2 (1 - \xi_2)]. \quad (87)$$

(2). The Φ -driven term $-\eta [(\partial \mathcal{L} / \partial \Phi) + (\partial \mathcal{L} / \partial \Phi)^\top] \gamma$. Using $\frac{\partial \mathcal{L}}{\partial \Phi} = -\sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i)$ and $(P_i - \mathbb{P}_i) M^\top = \Delta \beta r_i \gamma^\top$, we get

$$\frac{\partial \mathcal{L}}{\partial \Phi} = -\Delta \beta \sum_i \pi_i r_i e_i \gamma^\top \text{Var}(\mathbb{A}_i).$$

Thus, the Φ -driven term is

$$-\eta \left[\left(\frac{\partial \mathcal{L}}{\partial \Phi} \right) + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \right] \gamma = \eta \Delta \beta \sum_i \pi_i r_i (e_i \gamma^\top \text{Var}(\mathbb{A}_i) \gamma + \gamma_i \text{Var}(\mathbb{A}_i) \gamma)$$

By direct computation, we find that

$$\begin{aligned} \text{Var}(\mathbb{A}_i) \gamma &= \xi_i (1 - \xi_i) \Delta \gamma \left(1, -\frac{1}{d-1}, \dots, -\frac{1}{d-1} \right)^\top, \\ \gamma^\top \text{Var}(\mathbb{A}_i) \gamma &= \xi_i (1 - \xi_i) (\Delta \gamma)^2 \end{aligned}$$

Substituting into the equation, we get the Φ -driven term:

$$\dot{\gamma}_1|_\Phi = \eta \Delta \beta \left[\pi_1 r_1 \xi_1 (1 - \xi_1) ((\Delta \gamma)^2 + a \Delta \gamma) + (1 - \pi_1) r_2 \xi_2 (1 - \xi_2) \Delta \gamma \right], \quad (88)$$

$$\dot{\gamma}_{i \neq 1}|_\Phi = \frac{\eta \Delta \beta}{|\mathcal{V}| - 1} \left[-\pi_1 r_1 a \xi_1 (1 - \xi_1) \Delta \gamma + (1 - \pi_1) r_2 \xi_2 (1 - \xi_2) ((\Delta \gamma)^2 - b \Delta \gamma) \right]. \quad (89)$$

Combining (87)–(89) gives the closed ODEs for (a, b) on the invariant manifold.

Plug the above equations into the dynamics of $\gamma_1, \gamma_{i \neq 1}$

We now formally proceed with the proof of Theorem 3.7. We linearize the reduced system around the entry state of this phase and denote base values by superscript 0 and first-order variations by superscript 1.

Proof. The test for the critical point is the same as for Proposition 3.3, because the parameters are essentially located near the same minimum point.

We then linearize terms in (87)–(89) in a fixed order.

1. Linearization about M -driven term. We take the expansion up to the first order about ξ_i and \hat{p}_i and then substitute then into the expression of M -driven term.

(1). Linearization of the proxy attention weights ξ_1, ξ_2 . Take ξ_1 as an example,

$$\xi_1 = \frac{1}{1 + \frac{1-\pi_1}{\pi_1} \exp(-\eta\gamma_1\Delta\gamma)} = \pi_1 + \pi_1(1-\pi_1)\eta^1\gamma_1^0\Delta\gamma^0 + \mathcal{O}(\|\theta\|^2)$$

in which we use the fact that parameters locate near $\eta = 0$. Thus,

$$\xi_1^1 = \pi_1(1-\pi_1)\eta^1\gamma_1^0\Delta\gamma^0, \quad \xi_2^1 = \pi_1(1-\pi_1)\eta^1\gamma_{i \neq 1}^0\Delta\gamma^0.$$

(2). Linearization of the prediction probabilities \hat{p}_i and residuals r_i . Recall $\hat{p}_i = \sigma(m_i\Delta\beta - \log(|\mathcal{V}| - 1))$ with $m_i = b + \xi_i\Delta\gamma$. Expanding \hat{p}_i to first order gives (writing m_i^0 for the base value)

$$\begin{aligned} \hat{p}_1^1 &= \pi_1(1-\pi_1)\left((b^1 + \pi_1\Delta\gamma^1)\Delta\beta^0 + \pi_1(1-\pi_1)\eta a^0\Delta\gamma^0\Delta\beta^0 + m_1^0\Delta\beta^1\right), \\ \hat{p}_2^1 &= \pi_1(1-\pi_1)\left((b^1 + \pi_1\Delta\gamma^1)\Delta\beta^0 + \pi_1(1-\pi_1)\eta b^0\Delta\gamma^0\Delta\beta^0 + m_2^0\Delta\beta^1\right). \end{aligned}$$

Since $r_i = P_{i,1} - \hat{p}_i$, we have $r_i^1 = -\hat{p}_i^1$.

2. Linearization about Φ -driven term. Using the fact that parameters locate near $\eta = 0$, the linearization of Eq (87)–(89) corresponds to the right-hand side except that eta takes a value at the initial point.

Substituting the above expansions into (87)–(89), and keeping only first-order terms, yields

$$\begin{aligned} \dot{\gamma}_1^1 &= \Delta\beta^0(\pi_1^2(-\hat{p}_1^1) + \pi_1(1-\pi_1)(-\hat{p}_2^1)) + 3\lambda\eta^1\pi_1^2(1-\pi_1)^2\Delta\beta^0(\Delta\gamma^0)^2, \\ \dot{\gamma}_{i \neq 1}^1 &= \frac{\Delta\beta^0}{|\mathcal{V}| - 1}(\pi_1(1-\pi_1)(-\hat{p}_1^1) + (1-\pi_1)^2(-\hat{p}_2^1)) - \frac{1}{|\mathcal{V}| - 1}3\lambda\eta^1\pi_1^2(1-\pi_1)^2\Delta\beta^0(\Delta\gamma^0)^2. \end{aligned}$$

Taking the linear combination $(1-\pi_1)\dot{\gamma}_1^1 - (|\mathcal{V}| - 1)\pi_1\dot{\gamma}_{i \neq 1}^1$ cancels the $(-\hat{p}_i^1)$ terms and yields

$$(1-\pi_1)\dot{\gamma}_1^1 - (|\mathcal{V}| - 1)\pi_1\dot{\gamma}_{i \neq 1}^1 = 3\lambda\pi_1^2(1-\pi_1)^2\Delta\beta^0(\Delta\gamma^0)^2\eta^1. \quad (90)$$

Considering the linearized dynamics about η , there exists $c > 0$ such that

$$\dot{\eta}^1 = c\eta^1,$$

which indicating that η^1 admits a solution as

$$\eta^1(t) = \eta^1(t_0) \exp c(t - t_0). \quad (91)$$

Substituting the above equation into Eq. (90) and integrating both sides of the equation, we get

$$(1-\pi_1)\gamma_1^1(t) - (|\mathcal{V}| - 1)\pi_1\gamma_{i \neq 1}^1(t) = c'(\exp(c(t - t_0)) - 1) \quad (92)$$

Here, we use the fact that

$$(1-\pi_1)\gamma_1(t_0) - (d-1)\pi_1\gamma_{i \neq 1}(t_0) = 0.$$

□

D. Theoretical details in Sec. 3.4

This appendix provides detailed proofs for Section 4.4. We focus on the minimal vocabulary size $d = 3$ to exhibit the separation between secondary high frequency and secondary low frequency. Throughout, we use the rank-one parametrization on the invariant manifold (cf. Proposition 3.6)

$$M = \gamma\beta^\top, \quad \Phi = \eta\gamma\gamma^\top, \quad \theta = (\gamma, \beta) \in \mathbb{R}^3 \times \mathbb{R}^3.$$

Here, we do not need to consider η , because calculations show that its derivatives up to the second order are zero, so it will not affect our analysis.

D.1. A degenerate critical point on the rank-one manifold

We first formalize the “bad” critical point on the rank-one manifold under symmetric frequencies. This critical point is degenerate in the sense that the key driving terms $\partial\mathcal{L}/\partial M$ and $\partial\mathcal{L}/\partial\Phi$ vanish, hence linearization on the manifold cannot explain the escape to new embedding directions. The following is the proof of Proposition 3.8.

Proof. We construct a critical point on the rank-one manifold and show it is a local minimum for the linearized dynamics.

The critical point is constructed as follows. Take $\gamma_1\gamma_{i\neq 1} < 0$ as shown in Theorem 3.7. When η is sufficiently large, the attention proxy satisfies $\mathbb{A}_1 \approx e_1$ and $\mathbb{A}_{i\neq 1} \approx \hat{e}_1 := (0, \frac{1}{2}, \frac{1}{2})$. Choose $\beta_1 > \beta_{i\neq 1}$ so that $\text{softmax}(k\beta^\top) \rightarrow e_1$ as $k \rightarrow +\infty$ and $\text{softmax}(k\beta^\top) \rightarrow \hat{e}_1$ as $k \rightarrow -\infty$. Thus we may choose $\gamma_1 > 0$ and $\gamma_{i\neq 1} < 0$ so that

$$\mathbb{P}_1 = P_1, \quad \mathbb{P}_{i\neq 1} = \frac{1}{2}(P_2 + P_3).$$

By direct computation and symmetry of the data, we have $\frac{\partial\mathcal{L}}{\partial M} = 0$ and $\frac{\partial\mathcal{L}}{\partial\Phi} = 0$ at this point (refer to Lemma D.4). Substituting this fact into Eq. (4), it implies that our construction gives a critical point.

To verify local minimality for the linearized dynamics, we linearize the dynamics in Eq. (4). We compute $d(\frac{\partial\mathcal{L}}{\partial M})$ and $d(\frac{\partial\mathcal{L}}{\partial\Phi})$. At the constructed symmetric point, Lemma D.5 implies $\sum_i \pi_i d\mathbb{A}_i^\top (P_i - \mathbb{P}_i) = 0$ and hence

$$d\frac{\partial\mathcal{L}}{\partial M} = \sum_i \pi_i \mathbb{A}_i^\top d\mathbb{P}_i \neq 0.$$

Moreover, Lemma D.5 shows that $d\mathbb{P}_i = \mathbb{A}_i dM \text{Var}(\mathbb{P}_i)$. Since $\mathbb{A}_2 = \mathbb{A}_3$, it implies that $d\mathbb{P}_2 = d\mathbb{P}_3$. In addition, Lemma D.5 gives $d(\partial\mathcal{L}/\partial\Phi) = 0$.

As a result, the linearized dynamics on (γ, β, η) reduces to

$$\frac{d\Delta\gamma}{dt} = -d\left(\frac{\partial\mathcal{L}}{\partial M}\right)\beta, \quad \frac{d\Delta\beta}{dt} = -d\left(\frac{\partial\mathcal{L}}{\partial M}\right)^\top\gamma, \quad \frac{d\Delta\eta}{dt} = 0, \quad (93)$$

and the Jacobian $J_0 = -\nabla_\theta^2\mathcal{L}$ admits the explicit block form in Lemma D.6. In particular, J_0 is negative semidefinite with a nontrivial kernel. Hence, the critical point we constructed is a neutrally stable equilibrium for the linearized dynamics, which motivates the Lyapunov–Schmidt reduction in the main text. \square

D.2. Breaking the degeneracy: frequency perturbation and Lyapunov–Schmidt reduction

To eliminate the degeneracy, we perturb the frequencies between the two low-frequency states:

$$\pi^\top = \left(c, \frac{1-c}{2}, \frac{1-c}{2}\right) \Rightarrow \tilde{\pi}^\top = \left(c, \frac{1-c}{2} + \delta, \frac{1-c}{2} - \delta\right). \quad (94)$$

The dynamics becomes

$$\dot{\theta} = -\nabla_\theta\mathcal{L}(\theta, \delta).$$

We study the perturbed critical point by solving

$$-\nabla_\theta\mathcal{L}(\theta, \delta) = 0 \quad (95)$$

near the degenerate minimum, which we shift to $\theta = 0$ for convenience.

We use the formal expansion (at $\theta = 0$):

$$-\nabla_{\theta} \mathcal{L}(\theta, \delta) = J_0 \theta + \delta f_1 + \frac{1}{2} B(\theta, \theta) + \delta J_1 \theta + \frac{1}{2} \delta^2 f_2 + \text{h.o.t.}, \quad (96)$$

where

$$J_0 = -\nabla_{\theta}^2 \mathcal{L}, \quad f_1 = \partial_{\delta}(-\nabla_{\theta} \mathcal{L}), \quad B(\cdot, \cdot) = -\nabla_{\theta}^3 \mathcal{L}, \quad J_1 = \partial_{\delta} J_0, \quad f_2 = \partial_{\delta}^2(-\nabla_{\theta} \mathcal{L}).$$

Since J_0 is singular, we apply Lyapunov–Schmidt reduction.

D.2.1. KERNEL/RANGE DECOMPOSITION OF J_0

Proposition D.1 (Kernel and range bases). *Assume $\|\gamma\| = \|\beta\|$ and $\beta^{\top} \mathbf{1} = 0$ at the symmetric degenerate minimum. Then $\dim \ker(J_0) = 3$ and one convenient orthonormal basis is*

$$\begin{aligned} k_1 &= \frac{1}{\sqrt{2}} ((0, 1, -1), (0, 0, 0)), \\ k_2 &= \frac{1}{\sqrt{3}} ((0, 0, 0), (1, 1, 1)), \\ k_3 &= \frac{1}{\sqrt{\|\gamma\|^2 + \|\beta\|^2}} (-\gamma, \beta). \end{aligned} \quad (97)$$

An orthonormal basis for $\text{Range}(J_0)$ can be taken as

$$\begin{aligned} q_1 &= \frac{1}{\sqrt{4\gamma_2^2 + 2\gamma_1^2}} ((-2\gamma_2, \gamma_1, \gamma_1), (0, 0, 0)), \\ q_2 &= \frac{1}{\sqrt{2}} ((0, 0, 0), (0, 1, -1)), \\ q_3 &= \frac{1}{\sqrt{\|\gamma\|^2 + \|\beta\|^2}} (\gamma, \beta). \end{aligned} \quad (98)$$

Let $Q_K = (k_1, k_2, k_3)$ and $Q_R = (q_1, q_2, q_3)$, and denote projections $P_K = Q_K Q_K^{\top}$, $P_R = Q_R Q_R^{\top}$. Write $\theta = Q_K x + Q_R y$.

D.2.2. SOLVING THE RANGE EQUATION

Recall the Lyapunov–Schmidt decomposition $\theta = Q_K x + Q_R y$ and define the range equation

$$F_R(x, y, \delta) := -Q_R^{\top} \nabla_{\theta} \mathcal{L}(Q_K x + Q_R y, \delta) = 0. \quad (99)$$

Proposition D.2 (Range solution and first-order expansion). *Given a perturbation of the data parameterized by δ , the range equation (99) admits a unique solution $y = \zeta(x, \delta)$ in a neighborhood of $(x, \delta) = (0, 0)$. Moreover, it satisfies the expansion*

$$\zeta(x, \delta) = \delta \left(\begin{array}{c} 0 \\ \sqrt{2} \frac{\lambda \gamma_{i \neq 1} + (1 - \lambda)(\pi_1 \gamma_1 + (1 - \pi_1) \gamma_{i \neq 1})}{\pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1 - \pi_1) \gamma_{i \neq 1}^2 \mathbb{P}_{i \neq 1,2}} \\ 0 \end{array} \right) + \mathcal{O}(\delta^2 + \|x\|^2), \quad (100)$$

where the denominator

$$c_1 := \pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1 - \pi_1) \gamma_{i \neq 1}^2 \mathbb{P}_{i \neq 1,2}$$

is strictly positive under our standing assumptions (in particular $\gamma_1, \gamma_{i \neq 1} \neq 0$ and $\mathbb{P}_{1,2}, \mathbb{P}_{i \neq 1,2} > 0$).

Proof. We expand F_R around $(x, y, \delta) = (0, 0, 0)$. Writing $\theta = Q_K x + Q_R y$ and using

$$-\nabla_{\theta} \mathcal{L}(\theta, \delta) = J_0 \theta + \delta f_1 + \mathcal{O}(\|\theta\|^2 + \delta^2),$$

1430 we obtain

$$1431 \quad F_R(x, y, \delta) = Q_R^\top(J_0 Q_R y + \delta f_1) + \mathcal{O}(\|\theta\|^2 + \delta^2) = \Lambda_R y + \delta Q_R^\top f_1 + \mathcal{O}(\|x\|^2 + \|y\|^2 + \delta^2), \quad (101)$$

1433 where $\Lambda_R := Q_R^\top J_0 Q_R$.

1435 By Lemma D.7 we have an explicit expression for $Q_R^\top f_1$, and by Lemma D.8 the matrix Λ_R is invertible on the range
1436 coordinates; in particular, its (2, 2)-entry equals $-c_1 < 0$ and hence $(\Lambda_R^{-1})_{22} = -1/c_1$.

1437 Therefore, $\partial_y F_R(0, 0, 0) = \Lambda_R$ is invertible, and the implicit function theorem yields a unique smooth function $y = \zeta(x, \delta)$
1438 solving $F_R(x, \zeta(x, \delta), \delta) = 0$ locally, with

$$1440 \quad \zeta(x, \delta) = -\Lambda_R^{-1} \delta Q_R^\top f_1 + \mathcal{O}(\delta^2 + \|x\|^2). \quad (102)$$

1442 Since $Q_R^\top f_1$ has only a nonzero second component (Lemma D.7), and $(\Lambda_R^{-1})_{22} = -1/c_1$ (Lemma D.8), the second
1443 coordinate of ζ equals

$$1445 \quad \zeta_2(x, \delta) = -\left(-\frac{1}{c_1}\right) \delta \cdot \sqrt{2}(\lambda \gamma_{i \neq 1} + (1 - \lambda)(\pi_1 \gamma_1 + (1 - \pi_1) \gamma_{i \neq 1})) + \mathcal{O}(\delta^2 + \|x\|^2),$$

1447 which is exactly (100). This completes the proof. \square

1449 D.2.3. REDUCED KERNEL EQUATION AND APPROXIMATE CRITICAL POINT

1450 Plugging $y = \zeta(x, \delta)$ into the kernel equation gives

$$1452 \quad -Q_K^\top \nabla \mathcal{L}(Q_K x + Q_R \zeta(x, \delta), \delta) = 0.$$

1454 Because $J_0 Q_K = 0$ and $Q_K^\top f_1 = 0$, the leading contributions are second order:

$$1456 \quad Q_K^\top \left(\frac{1}{2} B(\theta, \theta) + \delta J_1 \theta + \frac{1}{2} \delta^2 f_2 \right) + \text{h.o.t.} = 0, \quad \theta = Q_K x + Q_R \zeta(x, \delta).$$

1458 **Theorem D.3** (Existence of an approximate critical point and its two-scale stability). *Let $\zeta(x, \delta)$ be given by Proposition D.2.*
1459 *Then $x = 0$ is an approximate solution of the reduced kernel equation up to second order, i.e.*

$$1461 \quad \|\nabla_\theta \mathcal{L}(Q_R \zeta(0, \delta), \delta)\| = \mathcal{O}(\delta^3).$$

1463 *Moreover, the linear stability splits into two scales:*

- 1465 1. **Slow manifold directions (within the rank-one manifold):** any positive eigenvalues created from the kernel directions
1466 are at most $\mathcal{O}(\delta^2)$.
- 1467 2. **Fast transverse directions (escaping the manifold):** Under condition in Lem. D.17, there exists a transverse positive
1468 eigenvalue of order $\Theta(\delta)$.

1470 *Proof.* The estimate $\|\nabla \mathcal{L}\| = \mathcal{O}(\delta^3)$ follows by inserting $\theta = Q_R \zeta(0, \delta)$ into the kernel expansion and using the explicit
1471 expressions:

1473 (1). $\frac{1}{2} Q_K^\top B(q_2 y_2, q_2 y_2)$ (From Lem. D.13):

$$1476 \quad \frac{1}{2} Q_K^\top B(q_2 y_2, q_2 y_2) = \frac{1}{\sqrt{\|\gamma\|^2 + \|\beta\|^2}} \begin{pmatrix} 0 \\ 0 \\ \pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1 - \pi_1) \gamma_2^2 \mathbb{P}_{i \neq 1,2} \end{pmatrix} y_2^2.$$

1479 (2). $\delta Q_K^\top J_1(q_2 y_2)$ (From Lem. D.10):

$$1482 \quad \delta Q_K^\top J_1(q_2 y_2) = -\delta \frac{\sqrt{2}}{\sqrt{\|\gamma\|^2 + \|\beta\|^2}} \begin{pmatrix} 0 \\ 0 \\ \pi_1 \gamma_1 (1 - \lambda) + \gamma_{i \neq 1} (\lambda + (1 - \pi_1)(1 - \lambda)) \end{pmatrix} y_2.$$

1485 (3). $\delta^2 f_2$ vanishes (From Lem. D.14).
 1486

1487 Substitute $y_2 = \sqrt{2} \frac{\lambda \gamma_{i \neq 1} + (1-\lambda)(\pi_1 \gamma_1 + (1-\pi_1) \gamma_{i \neq 1})}{\pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1-\pi_1) \gamma_{i \neq 1}^2 \mathbb{P}_{i \neq 1,2}} \delta$. We found that the second-order terms cancel each other out automati-
 1488 cally.
 1489

1490 For stability, we write the perturbed Hessian at the approximate critical point as
 1491

$$1492 \nabla_{\theta}^2 \mathcal{L}(\theta, \delta) = J_0 + \delta H_1 + \mathcal{O}(\delta^2), \quad H_1 := - \left(\begin{pmatrix} 0 & B \\ B^\top & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & C \end{pmatrix} \right) + J_1,$$

1493 with (B, C) computed from the second-order kernel reduction (see Lem. D.15).
 1494

1495 We take the basis as $Q = (Q_K, Q_R)$:
 1496

$$1497 Q^\top J Q = \begin{pmatrix} 0 & 0 \\ 0 & \Lambda \end{pmatrix} + \delta \begin{pmatrix} G & E \\ E^\top & F \end{pmatrix} + \mathcal{O}(\delta^2) \quad (103)$$

1498 where $\Lambda = Q_R^\top J_0 Q_R$, $G = Q_K^\top H_1 Q_K$, $E = Q_K^\top H_1 Q_R$, and $F = Q_R^\top H_1 Q_R$. Let λ be a small eigenvalue and the
 1499 corresponding eigenvector is (x, y) , the equation is
 1500

$$1501 \begin{cases} \delta G x + \delta E y + \mathcal{O}(\delta^2) = \lambda x \\ \Lambda y + \delta E^\top x + \mathcal{O}(\delta) y = \lambda y \end{cases} \quad (104)$$

1502 Since the new positive eigenvalue is small, we can solve y as
 1503

$$1504 y = -(\Lambda - \lambda I)^{-1} \delta E^\top x + \mathcal{O}(\delta^2) = -\delta \Lambda^{-1} E^\top x + \mathcal{O}(\delta^2) \quad (105)$$

1505 Substitute this expression into the the first equation, we get
 1506

$$1507 (\delta G - \delta^2 E \Lambda^{-1} E^\top) x = \lambda x + \mathcal{O}(\delta^3) \quad (106)$$

1508 Hence $\lambda = \mathcal{O}(\delta^2)$ provided $G = Q_K^\top H_1 Q_K = 0$. This vanishing is proved in Lemma D.16.
 1509

1510 Finally, we compute the eigenvalue of the normal directions. Recall the linearization of the whole dynamics is
 1511

$$1512 \begin{cases} \frac{d\Delta W_0}{dt} = \Delta \left(-\frac{\partial \mathcal{L}}{\partial M} W_1^\top - \frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K W_Q^\top - \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 W_Q W_K^\top \right) \\ \frac{d\Delta W_1}{dt} = \Delta \left(-W_0^\top \frac{\partial \mathcal{L}}{\partial M} \right) \\ \frac{d\Delta W_Q}{dt} = \Delta \left(-W_0^\top \frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K \right) \\ \frac{d\Delta W_K}{dt} = \Delta \left(-W_0^\top \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 W_Q \right) \end{cases} \quad (107)$$

1513 We find that
 1514

$$1515 \begin{cases} \frac{d\Delta W_0 \alpha_{1,\perp}}{dt} = -\frac{\partial \mathcal{L}}{\partial M} \Delta W_1^\top \alpha_{1,\perp} - \frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K \Delta W_Q^\top \alpha_{1,\perp} - \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 W_Q \Delta W_K^\top \alpha_{1,\perp} \\ \frac{d\alpha_{1,\perp}^\top \Delta W_1}{dt} = -\alpha_{1,\perp}^\top \Delta W_0^\top \frac{\partial \mathcal{L}}{\partial M} \\ \frac{d\alpha_{1,\perp}^\top \Delta W_Q}{dt} = -\Delta W_0^\top \frac{\partial \mathcal{L}}{\partial \Phi} W_0 W_K \\ \frac{d\alpha_{1,\perp}^\top \Delta W_K}{dt} = -\Delta W_0^\top \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top W_0 W_Q \end{cases} \quad (108)$$

Since $\partial_\delta \frac{\partial \mathcal{L}}{\partial \Phi} = 0$ and $d \frac{\partial \mathcal{L}}{\partial \Phi} = 0$, we find that $\frac{\partial \mathcal{L}}{\partial \Phi} = \mathcal{O}(\delta^2)$. So the main term is

$$\begin{cases} \frac{d\Delta W_0 \alpha_{1,\perp}}{dt} = -\frac{\partial \mathcal{L}}{\partial M} \Delta W_1^\top \alpha_{1,\perp} \\ \frac{d\alpha_{1,\perp}^\top \Delta W_1}{dt} = -\alpha_{1,\perp}^\top \Delta W_0^\top \frac{\partial \mathcal{L}}{\partial M} \end{cases} \quad (109)$$

From Lemma D.17, under some mild condition, We find that

$$\frac{\partial \mathcal{L}}{\partial M} = \Theta(\delta) \quad (110)$$

Thus, there exists positive eigenvalue at least order $\Theta(\delta)$. \square

D.3. Derivative toolbox

This section collects all derivative computations referenced in the proofs above.

D.3.1. VANISHING OF GRADIENTS

Lemma D.4 ($\partial \mathcal{L} / \partial M = 0$ and $\partial \mathcal{L} / \partial \Phi = 0$ at the constructed point). *At the symmetric degenerate minimum in Proposition 3.8, we have*

$$\frac{\partial \mathcal{L}}{\partial M} = 0, \quad \frac{\partial \mathcal{L}}{\partial \Phi} = 0.$$

Proof. By direct computation, we get $\mathbb{P}_1 = P_1$ and $\mathbb{P}_2 = \mathbb{P}_3 = \frac{1}{2}(P_2 + P_3)$. Thus, the terms in the gradients cancel after summing with $\pi_2 = \pi_3$. \square

D.3.2. FIRST-ORDER VARIATIONS

Lemma D.5 (First-order variations with respect to parameters). *At the critical point in Proposition 3.8, the first-order variations have the following form:*

1. *The variation of the attention proxy satisfies $d\mathbb{A}_1 = 0$ and $d\mathbb{A}_{i \neq 1} = \eta \gamma_{i \neq 1} (0, \frac{1}{4}(d\gamma_2 - d\gamma_3), -\frac{1}{4}(d\gamma_2 - d\gamma_3))$ for $i \neq 1$.*
2. *The variation of the output probability satisfies $d\mathbb{P}_i = \mathbb{A}_i dM \text{Var}(\mathbb{P}_i)$.*
3. *The variation of $\frac{\partial \mathcal{L}}{\partial M}$ and $\frac{\partial \mathcal{L}}{\partial \Phi}$ admit the following expression:*

$$d \frac{\partial \mathcal{L}}{\partial M} = \sum_i \pi_i \mathbb{A}_i^\top d\mathbb{P}_i, \quad d \frac{\partial \mathcal{L}}{\partial \Phi} = 0.$$

Proof. We calculate the first-order variation in sequence.

1. At $\mathbb{A}_1 = e_1$, we have $\text{diag}(e_1) - e_1 e_1^\top = 0$, hence $d\mathbb{A}_1 = 0$. For $i \neq 1$, using $\mathbb{A}_i = \hat{e}_1$ and $\Phi = \eta \gamma \gamma^\top$,

$$e_i^\top d\Phi = e_i^\top d(\eta \gamma \gamma^\top) = d(\eta \gamma_i \gamma^\top).$$

Since $\text{Var}(\mathbb{A}_{i \neq 1}) = \text{diag}(\hat{e}_1) - \hat{e}_1 \hat{e}_1^\top$ equals to

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & \frac{1}{4} & -\frac{1}{4} \\ 0 & -\frac{1}{4} & \frac{1}{4} \end{pmatrix},$$

we obtain the displayed vector form.

2. By definition,

$$d\mathbb{P}_i = d(\mathbb{A}_i M) \text{Var}(\mathbb{P}_i) = (d\mathbb{A}_i M + \mathbb{A}_i dM) \text{Var}(\mathbb{P}_i).$$

Under $M = \gamma\beta^\top$,

$$d\mathbb{A}_i M = d\mathbb{A}_i \gamma\beta^\top = (d\mathbb{A}_i \gamma)\beta^\top.$$

For $i = 1$, $d\mathbb{A}_1 = 0$, hence $d\mathbb{A}_1 M = 0$. For $i \neq 1$, $d\mathbb{A}_{i \neq 1} = \eta\gamma_{i \neq 1} (0, \frac{1}{4}(d\gamma_2 - d\gamma_3), -\frac{1}{4}(d\gamma_2 - d\gamma_3))$. Thus,

$$d\mathbb{A}_{i \neq 1} \gamma = \eta\gamma_{i \neq 1} (0, \frac{1}{4}(d\gamma_2 - d\gamma_3), -\frac{1}{4}(d\gamma_2 - d\gamma_3)) \cdot (\gamma_1, \gamma_2, \gamma_3) = \frac{1}{4}\eta\gamma_{i \neq 1}(d\gamma_2 - d\gamma_3)(\gamma_2 - \gamma_3) = 0,$$

since $\gamma_2 = \gamma_3$ at the symmetric point. Hence $d\mathbb{A}_{i \neq 1} M = 0$. Therefore $d\mathbb{P}_i = \mathbb{A}_i dM \text{Var}(\mathbb{P}_i)$.

Because $\mathbb{A}_2 = \mathbb{A}_3 = \hat{e}_1$ and $\text{Var}(\mathbb{P}_2) = \text{Var}(\mathbb{P}_3)$ under symmetry, we also have $d\mathbb{P}_2 = d\mathbb{P}_3$.

3. By definition of $\frac{\partial \mathcal{L}}{\partial M}$ and the chain rule,

$$d\frac{\partial \mathcal{L}}{\partial M} = -\sum_i \pi_i d\mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top d(-\mathbb{P}_i)$$

At the symmetric point, $P_1 - \mathbb{P}_1 = 0$ and $\sum_{i \neq 1} (P_i - \mathbb{P}_i) = 0$, while $d\mathbb{A}_2 = d\mathbb{A}_3$ and $\pi_2 = \pi_3$. Therefore the $i = 2, 3$ contributions cancel, giving $\sum_i \pi_i d\mathbb{A}_i^\top (P_i - \mathbb{P}_i) = 0$. It yields the claimed form.

By the definition of $\frac{\partial \mathcal{L}}{\partial \Phi}$ and the chain rule,

$$\begin{aligned} d\frac{\partial \mathcal{L}}{\partial \Phi} &= -d\left(\sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i)\right) \\ &= -\sum_i \pi_i e_i d(-\mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) dM^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top d\text{Var}(\mathbb{A}_i). \end{aligned}$$

Using the expression of $\text{Var}(\mathbb{A}_i)$, we get $\gamma^\top \text{Var}(\mathbb{A}_i) = 0$ since $\gamma_2 = \gamma_3$, which implies that the first term vanishes. Similarly, using the chain rule, we get $dM^\top = d\beta\gamma^\top + \beta d\gamma^\top$. Combined with $(P_i - \mathbb{P}_i)\beta = 0$ and $\gamma^\top \text{Var}(\mathbb{A}_i) = 0$, the second term vanishes. The third term vanishes due to the same reason.

□

D.3.3. HESSIAN MATRIX J_0

Lemma D.6 (Computation of J_0 on the rank-one manifold). *At the critical point in Proposition 3.8, the linearization restricted to the rank-one manifold yields the Hessian $J_0 = -\nabla_\theta^2 \mathcal{L}(\theta, 0)$ in the matrix form:*

$$J_0 = -\begin{pmatrix} c_1 & 0 & 0 & v_1 \\ 0 & c_2 & c_2 & v_2 \\ 0 & c_2 & c_2 & v_2 \\ v_1^\top & v_2^\top & v_2^\top & C \end{pmatrix} \quad (111)$$

where $c_1 := \pi_1 \|\beta\|_{\text{Var}(\mathbb{P}_1)}^2$, $c_2 := \frac{1}{4}(1 - \pi_1) \|\beta\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2$, $v_1 := \pi_1 \gamma_1 \beta^\top \text{Var}(\mathbb{P}_1)$, $v_2 := \frac{1}{2}(1 - \pi_1) \gamma_{i \neq 1} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1})$, and $C := \pi_1 \gamma_1^2 \text{Var}(\mathbb{P}_1) + (1 - \pi_1) \gamma_{i \neq 1}^2 \text{Var}(\mathbb{P}_{i \neq 1})$. Moreover, J_0 is negative semidefinite.

Proof. As shown in Eq. (93), the linearized dynamics on the rank-one manifold can be written as

$$\begin{aligned} \frac{d\Delta\gamma}{dt} &= -\left(d\frac{\partial \mathcal{L}}{\partial M}\right)\beta, \\ \frac{d\Delta\beta}{dt} &= -\left(d\frac{\partial \mathcal{L}}{\partial M}\right)^\top \gamma, \\ \frac{d\Delta\eta}{dt} &= 0, \end{aligned}$$

where we used $d(\partial \mathcal{L} / \partial \Phi) = 0$ at the symmetric point. We calculate the Jacobian corresponding to $\frac{d\Delta\gamma}{dt}$ and $\frac{d\Delta\beta}{dt}$ respectively.

1. The $\Delta\gamma$ equation. Using $d\frac{\partial\mathcal{L}}{\partial M} = \pi_1 \mathbb{A}_1^\top d\mathbb{P}_1 + (1 - \pi_1) \mathbb{A}_{i \neq 1}^\top d\mathbb{P}_{i \neq 1}$ and $\mathbb{A}_1 = e_1^\top$, $\mathbb{A}_{i \neq 1} = \hat{e}_1^\top = (0, \frac{1}{2}, \frac{1}{2})$, we obtain

$$\begin{aligned} \frac{d\Delta\gamma}{dt} &= -\pi_1 \mathbb{A}_1^\top d\mathbb{P}_1 \beta - (1 - \pi_1) \mathbb{A}_{i \neq 1}^\top d\mathbb{P}_{i \neq 1} \beta \\ &= -\pi_1 \mathbb{A}_1^\top \left(d\gamma_1 \beta^\top \text{Var}(\mathbb{P}_1) \beta + \gamma_1 d\beta^\top \text{Var}(\mathbb{P}_1) \beta \right) \\ &\quad - (1 - \pi_1) \mathbb{A}_{i \neq 1}^\top \left(\frac{1}{2} (d\gamma_2 + d\gamma_3) \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \beta + \gamma_{i \neq 1} d\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \beta \right), \end{aligned} \quad (112)$$

where we used the rank-one identity $dM = d(\gamma\beta^\top) = (d\gamma)\beta^\top + \gamma(d\beta)^\top$ and $d\mathbb{P}_i = \mathbb{A}_i dM \text{Var}(\mathbb{P}_i)$ at the symmetric point.

Let $d\theta := (d\gamma_1, d\gamma_2, d\gamma_3, d\beta)$, where $d\beta \in \mathbb{R}^3$. Collecting the coefficients in (112) gives the matrix form

$$\frac{d\Delta\gamma}{dt} = - \begin{pmatrix} c_1 & 0 & 0 & v_1 \\ 0 & c_2 & c_2 & v_2 \\ 0 & c_2 & c_2 & v_2 \end{pmatrix} d\theta. \quad (113)$$

2. The $\Delta\beta$ equation. Similarly,

$$\frac{d\Delta\beta}{dt} = -\pi_1 d\mathbb{P}_1^\top \mathbb{A}_1 \gamma - (1 - \pi_1) d\mathbb{P}_{i \neq 1}^\top \mathbb{A}_{i \neq 1} \gamma. \quad (114)$$

Using again $d\mathbb{P}_i = \mathbb{A}_i dM \text{Var}(\mathbb{P}_i)$ and $\mathbb{A}_1 = e_1^\top$, $\mathbb{A}_{i \neq 1} = \hat{e}_1^\top$, we obtain the compact matrix form

$$\frac{d\Delta\beta}{dt} = - \begin{pmatrix} v_1^\top & v_2^\top & v_2^\top & C \end{pmatrix} d\theta. \quad (115)$$

Combining (113) and (115), the linearization reads

$$\frac{d}{dt} \begin{pmatrix} \Delta\gamma \\ \Delta\beta \end{pmatrix} = J_0 d\theta,$$

where J_0 is exactly the block matrix.

We now verify that J_0 is negative semidefinite. Let $d\theta = (d\gamma_1, d\gamma_2, d\gamma_3, d\beta)$ and define

$$d\gamma_+ := \frac{1}{2}(d\gamma_2 + d\gamma_3), \quad d\gamma_- := \frac{1}{2}(d\gamma_2 - d\gamma_3).$$

A direct expansion of the quadratic form induced by (111) yields

$$d\theta^\top J_0 d\theta = -\pi_1 \left\| d\gamma_1 \beta + \gamma_1 d\beta \right\|_{\text{Var}(\mathbb{P}_1)}^2 - (1 - \pi_1) \left\| d\gamma_+ \beta + \gamma_{i \neq 1} d\beta \right\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2. \quad (116)$$

Indeed, for the $i = 1$ block one checks

$$-\pi_1 \left((d\gamma_1)^2 \beta^\top \text{Var}(\mathbb{P}_1) \beta + 2\gamma_1 d\gamma_1 d\beta^\top \text{Var}(\mathbb{P}_1) \beta + \gamma_1^2 d\beta^\top \text{Var}(\mathbb{P}_1) d\beta \right) = -\pi_1 \left\| d\gamma_1 \beta + \gamma_1 d\beta \right\|_{\text{Var}(\mathbb{P}_1)}^2.$$

For the low-frequency block, the coefficients $\frac{1}{4}(1 - \pi_1)$ in the $(d\gamma_2, d\gamma_3)$ -submatrix imply

$$-\frac{1}{4}(1 - \pi_1) \left\| \beta \right\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2 (d\gamma_2 + d\gamma_3)^2 = -(1 - \pi_1) \left\| d\gamma_+ \beta \right\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2,$$

and the cross/ $(d\beta, d\beta)$ terms match exactly the remaining pieces of $-(1 - \pi_1) \left\| d\gamma_+ \beta + \gamma_{i \neq 1} d\beta \right\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2$, giving (116).

Since $\text{Var}(\mathbb{P}_1) \succeq 0$ and $\text{Var}(\mathbb{P}_{i \neq 1}) \succeq 0$, the right-hand side of (116) is always non-positive, hence $J_0 \preceq 0$. Moreover, $d\gamma_-$ does not appear in (116), which already produces a nontrivial kernel direction; additional kernel directions arise from the scaling invariance $(d\gamma, d\beta) \propto (-\gamma, \beta)$ on the rank-one parametrization. Therefore, the equilibrium is a *degenerate* local minimum restricted to the rank-one manifold. \square

D.3.4. COMPUTATION OF f_1 AND $Q_R^\top J_0 Q_R$ FOR THE RANGE EQUATION

Lemma D.7 (Computation of f_1 and $Q_R^\top f_1$). *At the symmetric rank-one critical point, we have*

$$\partial_\delta \frac{\partial \mathcal{L}}{\partial \Phi} = 0, \quad -\partial_\delta \left(\frac{\partial \mathcal{L}}{\partial M} \right) \beta = 0, \quad -\left(\partial_\delta \left(\frac{\partial \mathcal{L}}{\partial M} \right) \right)^\top \gamma = (\lambda \gamma_{i \neq 1} + (1 - \lambda)(\pi_1 \gamma_1 + (1 - \pi_1) \gamma_{i \neq 1})) (0, 1, -1)^\top.$$

Consequently,

$$f_1 = (\lambda \gamma_{i \neq 1} + (1 - \lambda)(\pi_1 \gamma_1 + (1 - \pi_1) \gamma_{i \neq 1})) \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ -1 \end{pmatrix}, \quad (117)$$

and for the range basis $Q_R = (q_1, q_2, q_3)$ with $q_2 = \frac{1}{\sqrt{2}}(0, 0, 0, 0, 1, -1)$,

$$Q_R^\top f_1 = \begin{pmatrix} 0 \\ \sqrt{2}(\lambda \gamma_{i \neq 1} + (1 - \lambda)(\pi_1 \gamma_1 + (1 - \pi_1) \gamma_{i \neq 1})) \\ 0 \end{pmatrix}. \quad (118)$$

Proof. We differentiate the explicit gradient formula with respect to δ . We calculate the partial derivatives of $\frac{\partial \mathcal{L}}{\partial M}$ and $\frac{\partial \mathcal{L}}{\partial \Phi}$ with respect to δ , respectively.

1. Computation of $\frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial M}$. By definition,

$$\begin{aligned} \frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial M} &= \frac{\partial}{\partial \delta} \left(-\sum_i \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) \right) \\ &= -\sum_i \partial_\delta \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i \partial_\delta \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) \end{aligned}$$

Using $P_i = \lambda e_i^\top + (1 - \lambda) \pi^\top$, the first term is computed as

$$-\sum_i \partial_\delta \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) = -\mathbb{A}_{i \neq 1}^\top ((P_2 - \mathbb{P}_2) - (P_3 - \mathbb{P}_3)) = -\mathbb{A}_{i \neq 1}^\top (0, \lambda, -\lambda)$$

Since η is sufficiently large and $\gamma_1 \gamma_{i \neq 1} < 0$, we get

$$\partial_\delta \mathbb{A}_1 = (0, 0, 0), \quad \partial_\delta \mathbb{A}_{i \neq 1} = \left(0, \frac{1}{1 - \pi_1}, -\frac{1}{1 - \pi_1} \right)$$

Using $\sum_{i \neq 1} (P_i - \mathbb{P}_i) = 0$, the second term vanishes.

For the last term, we get

$$\partial_\delta P_i = (0, 1 - \lambda, -(1 - \lambda)), \quad \partial_\delta \mathbb{P}_i = 0$$

As a result,

$$\frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial M} = - \begin{pmatrix} 0 \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} (0, \lambda, -\lambda) - \begin{pmatrix} \pi_1 \\ \frac{1}{2}(1 - \pi_1) \\ \frac{1}{2}(1 - \pi_1) \end{pmatrix} (0, 1 - \lambda, -(1 - \lambda)). \quad (119)$$

2. Computation of $\frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial \Phi}$. By definition,

$$\begin{aligned} \frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial \Phi} &= \frac{\partial}{\partial \delta} \left(-\sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) \right) \\ &= -\sum_i \partial_\delta \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i \partial_\delta (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \partial_\delta \text{Var}(\mathbb{A}_i) \end{aligned}$$

Similar to the computation about $\frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial M}$, ones can check that $\frac{\partial}{\partial \delta} \frac{\partial \mathcal{L}}{\partial \Phi}$ vanishes.

Multiplying (119) by β on the right yields zero because it is proportional to $(0, 1, -1)$ and $\beta_2 = \beta_3$ at the symmetric point. Taking transpose and multiplying by γ on the right yields a multiple of $(0, 1, -1)^\top$, with the scalar coefficient $\lambda\gamma_{i \neq 1} + (1 - \lambda)(\pi_1\gamma_1 + (1 - \pi_1)\gamma_{i \neq 1})$, which gives the stated formula for f_1 in (117) under the definition of f_1 in the expansion of $-\nabla_\theta \mathcal{L}$.

Finally, (118) follows from $q_2^\top f_1 = \sqrt{2} \cdot (\text{scalar})$ and $q_1^\top f_1 = q_3^\top f_1 = 0$ by orthogonality. \square

Lemma D.8 (Structure of $Q_R^\top J_0 Q_R$ on the range). *Let $\Lambda_R := Q_R^\top J_0 Q_R$. Then Λ_R is nonsingular, and in particular,*

$$(\Lambda_R)_{22} = -c_1, \quad c_1 := \pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1 - \pi_1) \gamma_{i \neq 1}^2 \mathbb{P}_{i \neq 1,2} > 0. \quad (120)$$

Equivalently, Λ_R has the block structure

$$\Lambda_R = - \begin{pmatrix} * & 0 & * \\ 0 & c_1 & 0 \\ * & 0 & * \end{pmatrix},$$

where the starred entries are finite constants determined by the symmetric point, and are not needed in Proposition D.2.

Proof. This follows by substituting the explicit expression of J_0 (computed from the linearization on the rank-one manifold) into the orthonormal basis $Q_R = (q_1, q_2, q_3)$.

The key point is the q_2 direction. Recall $q_2 = \frac{1}{\sqrt{2}}(0, 0, 0, 0, 1, -1)$, i.e., it lies purely in the β -difference direction. At the symmetric point, the β -block of J_0 equals

$$J_{0,\beta\beta} = -(\pi_1 \gamma_1^2 \text{Var}(\mathbb{P}_1) + (1 - \pi_1) \gamma_{i \neq 1}^2 \text{Var}(\mathbb{P}_{i \neq 1})).$$

A direct computation gives

$$q_2^\top J_0 q_2 = -\left(\pi_1 \gamma_1^2 q_2^\top \text{Var}(\mathbb{P}_1) q_2 + (1 - \pi_1) \gamma_{i \neq 1}^2 q_2^\top \text{Var}(\mathbb{P}_{i \neq 1}) q_2 \right) = -(\pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1 - \pi_1) \gamma_{i \neq 1}^2 \mathbb{P}_{i \neq 1,2}),$$

where we used $q_2^\top \text{Var}(\mathbb{P}_i) q_2 = \mathbb{P}_{i,2}$ under the symmetric specialization $\mathbb{P}_{i,2} = \mathbb{P}_{i,3}$ (hence $\text{Var}(\mathbb{P}_i)$ acts diagonally on the $(2, -3)$ difference). This proves (120). The remaining entries are obtained similarly and yield the stated block structure, implying Λ_R is invertible on the range. \square

D.3.5. COMPUTATION OF CROSS TERM J_1

Lemma D.9 (Derivation of the mixed operator J_1). *Write $\theta = (\gamma, \beta) \in \mathbb{R}^3 \times \mathbb{R}^3$, and view J_1 as a 2×2 block operator with respect to the (γ, β) -splitting. Then*

$$J_1 = - \begin{pmatrix} J_{1,\gamma\gamma} & J_{1,\gamma\beta} \\ J_{1,\beta\gamma} & 0 \end{pmatrix} + \begin{pmatrix} 0 & A \\ A^\top & 0 \end{pmatrix}, \quad (121)$$

where

$$J_{1,\gamma\gamma} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \beta & 0 \\ 0 & 0 & -\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \beta \end{pmatrix}, \quad J_{1,\gamma\beta} = \begin{pmatrix} 0 \\ \gamma_{i \neq 1} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ -\gamma_{i \neq 1} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \end{pmatrix}, \quad (122)$$

$$J_{1,\beta\gamma} = \left(0, \gamma_{i \neq 1} \text{Var}(\mathbb{P}_{i \neq 1}) \beta, -\gamma_{i \neq 1} \text{Var}(\mathbb{P}_{i \neq 1}) \beta \right),$$

and

$$A = \begin{pmatrix} 0 & \pi_1(1 - \lambda) & -\pi_1(1 - \lambda) \\ 0 & \frac{1}{2}(\lambda + (1 - \pi_1)(1 - \lambda)) & -\frac{1}{2}(\lambda + (1 - \pi_1)(1 - \lambda)) \\ 0 & \frac{1}{2}(\lambda + (1 - \pi_1)(1 - \lambda)) & -\frac{1}{2}(\lambda + (1 - \pi_1)(1 - \lambda)) \end{pmatrix}. \quad (123)$$

Proof. We compute the mixed differential

$$J_1 = \partial_\delta \left(\nabla_\theta \left[-\nabla_\theta \mathcal{L}(\theta, \delta) \right] \right) \Big|_{(\theta, \delta) = (\theta_*, 0)}.$$

On the rank-one manifold, the (γ, β) -dynamics involve the two components

$$-\frac{\partial \mathcal{L}}{\partial M} \beta, \quad -\left(\frac{\partial \mathcal{L}}{\partial M}\right)^\top \gamma,$$

while the Φ -part does not contribute to J_1 at the symmetric point (see Step 2 below). Therefore it suffices to compute

$$\partial_\delta \nabla_\theta \left(-\frac{\partial \mathcal{L}}{\partial M} \beta \right), \quad \partial_\delta \nabla_\theta \left(-\left(\frac{\partial \mathcal{L}}{\partial M}\right)^\top \gamma \right).$$

We follow the same route as in the derivation of J_0 : we first compute $d(\partial \mathcal{L} / \partial M)$ and $d(\partial \mathcal{L} / \partial \Phi)$, then take ∂_δ and finally reassemble the induced variation of the rank-one gradients.

Step 1: computing $\partial_\delta \nabla_\theta(\partial \mathcal{L} / \partial M)$. Recall

$$\frac{\partial \mathcal{L}}{\partial M} = -\sum_i \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i).$$

Taking θ -differential gives

$$d \frac{\partial \mathcal{L}}{\partial M} = -\sum_i \pi_i (d\mathbb{A}_i^\top) (P_i - \mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top d(P_i - \mathbb{P}_i),$$

and since $dP_i = 0$, we have $d(P_i - \mathbb{P}_i) = -d\mathbb{P}_i$. Differentiating w.r.t. δ and using the product rule yields

$$\begin{aligned} \partial_\delta d \frac{\partial \mathcal{L}}{\partial M} &= -\sum_i (\partial_\delta \pi_i) (d\mathbb{A}_i^\top) (P_i - \mathbb{P}_i) - \sum_i \pi_i \partial_\delta d\mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i (d\mathbb{A}_i^\top) \partial_\delta (P_i - \mathbb{P}_i) \\ &\quad - \sum_i (\partial_\delta \pi_i) \mathbb{A}_i^\top (-d\mathbb{P}_i) - \sum_i \pi_i (\partial_\delta \mathbb{A}_i^\top) (-d\mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (-d\mathbb{P}_i). \end{aligned} \tag{124}$$

We now analyze each term. (All computations are evaluated at the symmetric point.)

1. The term $-\sum_i (\partial_\delta \pi_i) (d\mathbb{A}_i^\top) (P_i - \mathbb{P}_i)$: using the structure of $d\mathbb{A}_i$ and $(P_i - \mathbb{P}_i)\beta = 0$, its contribution vanishes when paired with β and with γ , i.e.,

$$\left(-\sum_i (\partial_\delta \pi_i) (d\mathbb{A}_i^\top) (P_i - \mathbb{P}_i) \right) \beta = 0, \quad \left(\cdot \right)^\top \gamma = 0.$$

2. The term $-\sum_i \pi_i \partial_\delta d\mathbb{A}_i^\top (P_i - \mathbb{P}_i)$: Since $d\mathbb{A}_i = e_i^\top \Phi \text{Var}(\mathbb{A}_i)$, we get $\partial_\delta d\mathbb{A}_i = e_i^\top \Phi (\partial_\delta \text{Var}(\mathbb{A}_i))$. For $i = 1$, we have $\partial_\delta \text{Var}(\mathbb{A}_1) = 0$ since $\mathbb{A}_1 = e_1^\top$. For $i \neq 1$,

$$\partial_\delta \text{Var}(\mathbb{A}_i) = (1 - \pi_1) \left[\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} \begin{pmatrix} 0, \frac{1}{2}, \frac{1}{2} \end{pmatrix} - \begin{pmatrix} 0 \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} (0, 1, -1) \right] = 0$$

Hence this term is zero.

3. The term $-\sum_i \pi_i (d\mathbb{A}_i^\top) \partial_\delta (P_i - \mathbb{P}_i)$: since $\partial_\delta (P_i - \mathbb{P}_i) = \partial_\delta P_i$, we obtain

$$-\sum_i \pi_i (d\mathbb{A}_i^\top) \partial_\delta (P_i - \mathbb{P}_i) = \sum_i \pi_i (d\mathbb{A}_i^\top) (0, 1 - \lambda, -(1 - \lambda)).$$

By the symmetric specialization $\beta_2 = \beta_3$ and $\gamma_2 = \gamma_3$, this term also satisfies

$$\left(\cdot \right) \beta = 0, \quad \left(\cdot \right)^\top \gamma = 0.$$

4. The term $-\sum_i (\partial_\delta \pi_i) \mathbb{A}_i^\top (-d\mathbb{P}_i)$: using $\partial_\delta \pi_2 = -\partial_\delta \pi_3$ and $\mathbb{A}_2 = \mathbb{A}_3$ at $\delta = 0$, we get

$$-\sum_i (\partial_\delta \pi_i) \mathbb{A}_i^\top (-d\mathbb{P}_i) = -\mathbb{A}_2^\top (-d\mathbb{P}_2) + \mathbb{A}_3^\top (-d\mathbb{P}_3) = 0.$$

1870 5. The term $-\sum_i \pi_i \partial_\delta \mathbb{A}_i^\top (-d\mathbb{P}_i)$: We have

$$1871 \quad -\sum_i \pi_i \partial_\delta \mathbb{A}_i^\top (-d\mathbb{P}_i) = \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} \mathbb{A}_{i \neq 1} dM \text{Var}(\mathbb{P}_i) \quad (125)$$

1875 6. The term $-\sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (-d\mathbb{P}_i)$:

$$1876 \quad -\sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (-d\mathbb{P}_i) = \sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (d\mathbb{A}_i M \text{Var}(\mathbb{P}_i) + \mathbb{A}_i dM \text{Var}(\mathbb{P}_i))$$

1877 Similar to the previous computation, we have

$$1878 \quad -\sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (-d\mathbb{P}_i) = \mathbb{A}_{i \neq 1}^\top (0, 1, -1) (d\gamma) \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}). \quad (126)$$

1883 The last two terms, $-\sum_i \pi_i (\partial_\delta \mathbb{A}_i^\top) (-d\mathbb{P}_i)$ and $-\sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (-d\mathbb{P}_i)$, produce the only nonzero contribution to $\partial_\delta d(\partial \mathcal{L} / \partial M) \beta$ along the $(2, -3)$ antisymmetric direction. Collecting them gives

$$1884 \quad \partial_\delta d\left(\frac{\partial \mathcal{L}}{\partial M}\right) \beta = \begin{pmatrix} 0 & 0 & 0 & \left| & 0 \\ 0 & \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \beta & 0 & \left| & \gamma_{i \neq 1} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & 0 & -\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \beta & \left| & -\gamma_{i \neq 1} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \end{pmatrix} d(\gamma, \beta), \quad (127)$$

1889 which exactly corresponds to the $J_{1, \gamma \gamma}$ and $J_{1, \gamma \beta}$ blocks in (122).

1891 **Step 2:** $\partial_\delta d(\partial \mathcal{L} / \partial \Phi) = 0$. We differentiate

$$1892 \quad \frac{\partial \mathcal{L}}{\partial \Phi} = -\sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i),$$

1895 and check term by term (product rule) that every contribution vanishes at the symmetric point: the $\partial_\delta \pi_i$ -terms cancel by symmetry and the ∂_δ -dependence of $\text{Var}(\mathbb{A}_i)$ does not contribute at $\delta = 0$. Hence $\partial_\delta d(\partial \mathcal{L} / \partial \Phi) = 0$.

1898 **Step 3: contribution from $\partial_\delta (\partial \mathcal{L} / \partial M) d\beta$.** Using the explicit formula of $\partial_\delta (\partial \mathcal{L} / \partial M)$ (computed previously), we obtain the linear map acting on $d\beta$:

$$1901 \quad \partial_\delta \left(\frac{\partial \mathcal{L}}{\partial M}\right) d\beta = -\begin{pmatrix} 0 & \left| & \pi_1(1-\lambda) & \right. & -\pi_1(1-\lambda) \\ 0 & \left| & \frac{1}{2}(\lambda + (1-\pi_1)(1-\lambda)) & \right. & -\frac{1}{2}(\lambda + (1-\pi_1)(1-\lambda)) \\ 0 & \left| & \frac{1}{2}(\lambda + (1-\pi_1)(1-\lambda)) & \right. & -\frac{1}{2}(\lambda + (1-\pi_1)(1-\lambda)) \end{pmatrix} d\beta, \quad (128)$$

1904 which is exactly the A block in (123) (placed in the (γ, β) off-diagonal).

1906 **Step 4: assembling J_1 from the two dynamics components.** By the chain rule,

$$1907 \quad \partial_\delta d\left(-\frac{\partial \mathcal{L}}{\partial M} \beta\right) = -\left(\partial_\delta d\frac{\partial \mathcal{L}}{\partial M}\right) \beta - \left(\partial_\delta \frac{\partial \mathcal{L}}{\partial M}\right) d\beta,$$

1909 so combining (127) and (128) yields the γ -equation blocks in (121).

1911 Similarly,

$$1912 \quad \partial_\delta d\left(-\left(\frac{\partial \mathcal{L}}{\partial M}\right)^\top \gamma\right) = -\left(\partial_\delta d\frac{\partial \mathcal{L}}{\partial M}\right)^\top \gamma - \left(\partial_\delta \frac{\partial \mathcal{L}}{\partial M}\right)^\top d\gamma,$$

1913 which gives the (β, γ) block $J_{1, \beta \gamma}$ together with the transpose A^\top in (121). \square

1916 Next, we will calculate the identity needed in Theorem D.3.

1917 **Lemma D.10.** Let $\theta = q_2 y_2$ be the leading-order reduction (since $\zeta(0, \delta) \sim \delta q_2$). Then

$$1918 \quad Q_K^\top J_1(q_2 y_2) = -\frac{\sqrt{2}}{\sqrt{\|\gamma\|^2 + \|\beta\|^2}} \begin{pmatrix} 0 \\ 0 \\ \pi_1 \gamma_1 (1-\lambda) + \gamma_{i \neq 1} (\lambda + (1-\pi_1)(1-\lambda)) \end{pmatrix} y_2.$$

1922 *Proof.* This can be verified using the expression in Lem. D.9 and by direct calculation. \square

D.3.6. COMPUTATION OF THE BILINEAR FORM $B(\cdot, \cdot)$

Before proceeding with the specific calculations, let's review the following lemma.

Lemma D.11 (Second differential). *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be C^2 . Then for any $h \in \mathbb{R}^d$,*

$$f(\theta + h) = f(\theta) + df(\theta)[h] + \frac{1}{2} d^2 f(\theta)[h, h] + \mathcal{O}(\|h\|^3), \quad (129)$$

where $d^2 f(\theta)[u, v] = u^\top \nabla^2 f(\theta) v$ is the (symmetric) bilinear form induced by the Hessian. The same expansion applies componentwise to vector-valued maps; in particular, for the gradient map $g(\theta) = \nabla f(\theta)$,

$$g(\theta + h) = g(\theta) + Dg(\theta) h + \frac{1}{2} D^2 g(\theta)[h, h] + \mathcal{O}(\|h\|^3). \quad (130)$$

Also, we have the following lemma which simplifies the computation.

Lemma D.12 (A useful identity: $d \text{Var}(\mathbb{A}_{i \neq 1}) = 0$ for the antisymmetric direction). *At $\mathbb{A}_{i \neq 1} = \hat{e}_1 = (0, \frac{1}{2}, \frac{1}{2})$, if $d\mathbb{A}_{i \neq 1} = (0, a, -a)$ for some a , then $d \text{Var}(\mathbb{A}_{i \neq 1}) = 0$.*

Proof. By definition, $d \text{Var}(\mathbb{A}) = \text{diag}(d\mathbb{A}) - (d\mathbb{A})\mathbb{A}^\top - \mathbb{A}(d\mathbb{A})^\top$. Substituting $\mathbb{A} = \hat{e}_1$ and $d\mathbb{A} = (0, a, -a)$ gives exact cancellation of all entries. \square

Consequently, in our regime the second differential of \mathbb{A}_i simplifies to

$$d^2 \mathbb{A}_i = d^2(\eta \gamma_i \gamma_i^\top) \text{Var}(\mathbb{A}_i), \quad (131)$$

because the potentially present term $d(\eta \gamma_i \gamma_i^\top) d \text{Var}(\mathbb{A}_i)$ vanishes (identically for $i = 1$ since $d\mathbb{A}_1 = 0$, and by Lemma D.12 for $i \neq 1$).

Now we will begin the calculation of the bilinear term B .

Second differential of $\frac{\partial \mathcal{L}}{\partial M}$. Recall

$$\frac{\partial \mathcal{L}}{\partial M} = - \sum_i \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i). \quad (132)$$

Differentiating once (with π_i fixed) yields

$$d \frac{\partial \mathcal{L}}{\partial M} = - \sum_i \pi_i d\mathbb{A}_i^\top (P_i - \mathbb{P}_i) + \sum_i \pi_i \mathbb{A}_i^\top d\mathbb{P}_i, \quad (133)$$

since $d(P_i - \mathbb{P}_i) = -d\mathbb{P}_i$. Differentiating again gives the decomposition

$$d^2 \frac{\partial \mathcal{L}}{\partial M} = - \sum_i \pi_i d^2 \mathbb{A}_i^\top (P_i - \mathbb{P}_i) + 2 \sum_i \pi_i d\mathbb{A}_i^\top d\mathbb{P}_i + \sum_i \pi_i \mathbb{A}_i^\top d^2 \mathbb{P}_i. \quad (134)$$

The coefficient 2 in the middle term is the standard product-rule contribution: it comes once from differentiating $-\sum \pi_i d\mathbb{A}_i^\top (P_i - \mathbb{P}_i)$ and once from differentiating $+\sum \pi_i \mathbb{A}_i^\top d\mathbb{P}_i$.

On $d^2 \mathbb{P}_i$. Using $d\mathbb{P}_i = d(\mathbb{A}_i M) \text{Var}(\mathbb{P}_i)$, we have

$$d^2 \mathbb{P}_i = d^2(\mathbb{A}_i M) \text{Var}(\mathbb{P}_i) + d(\mathbb{A}_i M) d \text{Var}(\mathbb{P}_i). \quad (135)$$

Moreover,

$$d^2(\mathbb{A}_i M) = d^2 \mathbb{A}_i M + 2 d\mathbb{A}_i dM + \mathbb{A}_i d^2 M. \quad (136)$$

From $d^2 \frac{\partial \mathcal{L}}{\partial M}$ to the quadratic term in the vector field In the rank-one dynamics, the γ -component contains the factor $(\frac{\partial \mathcal{L}}{\partial M})\beta$. At the critical point, $\frac{\partial \mathcal{L}}{\partial M} = 0$, hence

$$d^2 \left(\frac{\partial \mathcal{L}}{\partial M} \beta \right) = \left(d^2 \frac{\partial \mathcal{L}}{\partial M} \right) \beta + 2 \left(d \frac{\partial \mathcal{L}}{\partial M} \right) d\beta. \quad (137)$$

An analogous identity holds for $d^2 \left(\left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top \gamma \right)$.

Decomposition into explicit matrix blocks. We decompose the resulting bilinear form $B(\cdot, \cdot)$ into contributions coming from the different terms in (134)–(135) and from $d^2 \frac{\partial \mathcal{L}}{\partial \Phi}$. Concretely, for each output coordinate k ,

$$B_k(\cdot, \cdot) = \sum_{\ell} B_k^{(\ell)}(\cdot, \cdot), \quad (138)$$

where $B^{(1)}\text{--}B^{(5)}$ come from the M -part and $B^{(6)}$ comes from the Φ -part.

We calculate bilinear term for $1 \leq k \leq 3$ and $4 \leq k \leq 6$ respectively.

1. The computation of B_k for $1 \leq k \leq 3$. We take the second differential of $-\frac{\partial \mathcal{L}}{\partial M} \beta - \eta \left(\frac{\partial \mathcal{L}}{\partial \Phi} + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \right) \gamma$,

$$-d^2 \left(\frac{\partial \mathcal{L}}{\partial M} \beta \right) - d^2 \left[\eta \left(\frac{\partial \mathcal{L}}{\partial \Phi} + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \right) \gamma \right]$$

The computation of M -term and Φ -term is computed as follows.

M -term.

- (a) Contribution from $2(d(\partial \mathcal{L}/\partial M)) d\beta$. This produces the blocks denoted by $B_k^{(1)}$:

$$B_1^{(1)}(\cdot, \cdot) = -\pi_1 \begin{pmatrix} 0 & 0 & 0 & \beta^\top \text{Var}(\mathbb{P}_1) \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \text{Var}(\mathbb{P}_1)\beta & 0 & 0 & 2\gamma_1 \text{Var}(\mathbb{P}_1) \end{pmatrix}, \quad (139)$$

$$B_2^{(1)}(\cdot, \cdot) = B_3^{(1)}(\cdot, \cdot) = -(1 - \pi_1) \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{4}\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & 0 & 0 & \frac{1}{4}\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})\beta & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})\beta & (1 - \pi_1)\gamma_{i \neq 1} \text{Var}(\mathbb{P}_{i \neq 1}) \end{pmatrix}.$$

- (b) Contribution from $\sum_i d^2 \mathbb{A}_i^\top (P_i - \mathbb{P}_i)\beta$. This term vanishes due to $(P_i - \mathbb{P}_i)\beta = 0$ for each i .

- (c) Contribution from $-2 \sum_i \pi_i d\mathbb{A}_i^\top d\mathbb{P}_i\beta$. Using the expression of $d\mathbb{A}_i$ and $d\mathbb{P}_i$,

$$\begin{aligned} -2 \sum_i \pi_i d\mathbb{A}_i^\top d\mathbb{P}_i\beta &= -2(1 - \pi_1)\eta\gamma_{i \neq 1} \begin{pmatrix} 0 \\ \frac{1}{4}(d\gamma_2 - d\gamma_3) \\ -\frac{1}{4}(d\gamma_2 - d\gamma_3) \end{pmatrix} d(\mathbb{A}_i M) \text{Var}(\mathbb{P}_i)\beta \\ &= -\frac{1}{2}(1 - \pi_1)\eta\gamma_{i \neq 1} \begin{pmatrix} 0 \\ d\gamma_2 - d\gamma_3 \\ -(d\gamma_2 - d\gamma_3) \end{pmatrix} \mathbb{A}_i dM \text{Var}(\mathbb{P}_i)\beta \end{aligned}$$

This produces the blocks denoted by $B_k^{(2)}$:

$$B_1^{(2)}(\cdot, \cdot) = 0,$$

$$B_2^{(2)}(\cdot, \cdot) = -\frac{1}{2}(1 - \pi_1)\eta\gamma_{i \neq 1} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2}\|\beta\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2 & 0 & \frac{1}{2}\gamma_{i \neq 1}\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & 0 & -\frac{1}{2}\|\beta\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2 & -\frac{1}{2}\gamma_{i \neq 1}\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & \frac{1}{2}\gamma_{i \neq 1} \text{Var}(\mathbb{P}_{i \neq 1})\beta & -\frac{1}{2}\gamma_{i \neq 1} \text{Var}(\mathbb{P}_{i \neq 1})\beta & 0 \end{pmatrix},$$

$$B_3^{(2)}(\cdot, \cdot) = -B_2^{(2)}(\cdot, \cdot).$$

(140)

- (d) Contribution from $-\sum_i \pi_i \mathbb{A}_i^\top d^2 \mathbb{P}_i\beta$. We further split into:

- Terms contributed by $-2 \sum_i \pi_i \mathbb{A}_i^\top d\mathbb{A}_i dM \text{Var}(\mathbb{P}_i)\beta$:

$$\begin{aligned}
 B_1^{(3)}(\cdot, \cdot) &= 0 \\
 B_2^{(3)}(\cdot, \cdot) &= -\frac{1}{4}(1 - \pi_1)\eta\gamma_{i \neq 1} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & \beta^\top \text{Var}(\mathbb{P}_{i \neq 1})\beta & -\beta^\top \text{Var}(\mathbb{P}_{i \neq 1})\beta & 0 \\ 0 & -\beta^\top \text{Var}(\mathbb{P}_{i \neq 1})\beta & \beta^\top \text{Var}(\mathbb{P}_{i \neq 1})\beta & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\
 B_3^{(3)}(\cdot, \cdot) &= B_2^{(3)}(\cdot, \cdot)
 \end{aligned} \tag{141}$$

- Terms contributed by $-\sum_i \pi_i \mathbb{A}_i^\top \mathbb{A}_i d^2 M \text{Var}(\mathbb{P}_i)\beta$. The matrix form is

$$\begin{aligned}
 B_1^{(4)} &= -\pi_1 \begin{pmatrix} 0 & 0 & 0 & \beta^\top \text{Var}(\mathbb{P}_1) \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \text{Var}(\mathbb{P}_1)\beta & 0 & 0 & 0 \end{pmatrix} \\
 B_2^{(4)} = B_3^{(4)} &= -(1 - \pi_1) \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{4}\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & 0 & 0 & \frac{1}{4}\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & \frac{1}{4}\text{Var}(\mathbb{P}_{i \neq 1})\beta & \frac{1}{4}\text{Var}(\mathbb{P}_{i \neq 1})\beta & 0 \end{pmatrix}
 \end{aligned} \tag{142}$$

- Terms contributed by $-\sum_i \pi_i \mathbb{A}_i^\top d(\mathbb{A}_i M) d \text{Var}(\mathbb{P}_i)\beta$. The matrix form is of the shape The matrix form is of the shape

$$B_1^{(5)} = -\pi_1 \begin{pmatrix} c_1 & 0 & 0 & c_2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ c_2^\top & 0 & 0 & c_3 \end{pmatrix} \tag{143}$$

where

$$\begin{aligned}
 c_1 &= \beta^\top \text{diag}(\text{Var}(\mathbb{P}_1))\beta - 2(\mathbb{P}_1\beta)^\top \|\beta\|_{\text{Var}(\mathbb{P}_1)}^2 \\
 c_2 &= \frac{1}{2}\gamma_1 (\beta^{\odot 2, \top} \text{Var}(\mathbb{P}_1) + \beta^\top \odot \beta^\top \text{Var}(\mathbb{P}_1) - 3(\mathbb{P}_1\beta)\beta^\top \text{Var}(\mathbb{P}_1) - \beta^\top \text{Var}(\mathbb{P}_1)\beta\mathbb{P}_1)
 \end{aligned} \tag{144}$$

And

$$c_3(d\beta, d\beta) = \gamma_1^2 (d\beta^\top \text{diag}(d\beta^\top \text{Var}(\mathbb{P}_1))\beta - d\beta^\top \mathbb{P}_1^\top d\beta^\top \text{Var} \mathbb{P}_1 \beta - d\beta^\top \text{Var}(\mathbb{P}_1)d\beta\mathbb{P}_1\beta).$$

Writing into the matrix, we get

$$\begin{aligned}
 c_3 &= \gamma_1^2 \left(\text{diag}(\mathbb{P}_1 \odot \beta) - \frac{1}{2} ((\mathbb{P}_1^\top \mathbb{P}_1 \odot \beta) + (\mathbb{P}_1^\top \mathbb{P}_1 \odot \beta)^\top) \right. \\
 &\quad \left. - \frac{1}{2} ((\mathbb{P}_1^\top \beta^\top \text{Var}(\mathbb{P}_1)) + (\mathbb{P}_1^\top \beta^\top \text{Var}(\mathbb{P}_1))^\top) - (\mathbb{P}_1\beta) \text{Var}(\mathbb{P}_1) \right)
 \end{aligned} \tag{145}$$

And

$$B_2^{(5)} = B_3^{(5)} = -\frac{1}{2}(1 - \pi_1) \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & c_4 & c_4 & c_5 \\ 0 & c_4 & c_4 & c_5 \\ 0 & c_5^\top & c_5^\top & c_6 \end{pmatrix} \tag{146}$$

where $c_6(d\beta, d\beta) = \gamma_{i \neq 1}^2 (d\beta^\top \text{diag}(d\beta^\top \text{Var}(\mathbb{P}_{i \neq 1}))\beta - d\beta^\top \mathbb{P}_{i \neq 1}^\top d\beta^\top \text{Var} \mathbb{P}_{i \neq 1} \beta - d\beta^\top \text{Var}(\mathbb{P}_1)d\beta\mathbb{P}_{i \neq 1}\beta)$ and the matrix form is:

$$\begin{aligned}
 c_6 &= \text{diag}(\mathbb{P}_{i \neq 1} \odot \beta) - \frac{1}{2} \left((\mathbb{P}_{i \neq 1}^\top \mathbb{P}_{i \neq 1} \odot \beta) + (\mathbb{P}_{i \neq 1}^\top \mathbb{P}_{i \neq 1} \odot \beta)^\top \right) \\
 &\quad - \frac{1}{2} \left((\mathbb{P}_{i \neq 1}^\top \beta^\top \text{Var}(\mathbb{P}_{i \neq 1})) + (\mathbb{P}_{i \neq 1}^\top \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}))^\top \right) - (\mathbb{P}_{i \neq 1}\beta) \text{Var}(\mathbb{P}_{i \neq 1})
 \end{aligned} \tag{147}$$

Φ -term. Using the fact that $\frac{\partial \mathcal{L}}{\partial \Phi} = 0$ and $d\frac{\partial \mathcal{L}}{\partial \Phi} = 0$, we get

$$-d^2 \left[\eta \left(\frac{\partial \mathcal{L}}{\partial \Phi} + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \right) \gamma \right] = -\eta d^2 \left(\frac{\partial \mathcal{L}}{\partial \Phi} + \left(\frac{\partial \mathcal{L}}{\partial \Phi} \right)^\top \right) \gamma$$

Differentiating twice $\frac{\partial \mathcal{L}}{\partial \Phi} = -\sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i)$ gives

$$d^2 \frac{\partial \mathcal{L}}{\partial \Phi} = 2 \sum_i \pi_i e_i d\mathbb{P}_i dM^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) d^2 M^\top \text{Var}(\mathbb{A}_i), \quad (148)$$

where the second term vanishes after contraction with γ at the critical point by the same symmetry used in the first-order analysis. The first term yields $B^{(6)}$ blocks:

$$B_1^{(6)}(\cdot, \cdot) = 0,$$

$$B_2^{(6)}(\cdot, \cdot) = -B_3^{(6)}(\cdot, \cdot) = -\frac{1}{2} \eta \gamma_{i \neq 1} (1 - \pi_1) \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2} \|\beta\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2 & 0 & \frac{1}{2} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & 0 & -\frac{1}{2} \|\beta\|_{\text{Var}(\mathbb{P}_{i \neq 1})}^2 & -\frac{1}{2} \beta^\top \text{Var}(\mathbb{P}_{i \neq 1}) \\ 0 & \frac{1}{2} \text{Var}(\mathbb{P}_{i \neq 1}) \beta & -\frac{1}{2} \text{Var}(\mathbb{P}_{i \neq 1}) \beta & 0 \end{pmatrix}. \quad (149)$$

2. The computation of B_k for $4 \leq k \leq 6$. Similarly, we compute the $d^2 \left(\left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top \gamma \right)$

$$\begin{aligned} d^2 \left(\left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top \gamma \right) &= d^2 \left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top \gamma + 2d \left(\frac{\partial \mathcal{L}}{\partial M} \right)^\top d\gamma \\ &= \left(2 \sum_i \pi_i d\mathbb{A}_i^\top d\mathbb{P}_i + \sum_i \pi_i \mathbb{A}_i^\top d^2 \mathbb{P}_i \right)^\top \gamma + 2 \left(\sum_i \pi_i \mathbb{A}_i^\top \mathbb{A}_i dM \text{Var}(\mathbb{P}_i) \right)^\top d\gamma \\ &= \sum_i \pi_i d^2 \mathbb{P}_i^\top \mathbb{A}_i \gamma + 2 \sum_i \pi_i \text{Var}(\mathbb{P}_i) dM^\top \mathbb{A}_i^\top \mathbb{A}_i d\gamma \end{aligned} \quad (150)$$

For the term $2 \sum_i \pi_i \text{Var}(\mathbb{P}_i) dM^\top \mathbb{A}_i^\top \mathbb{A}_i d\gamma$, we get

$$2 \sum_i \pi_i \text{Var}(\mathbb{P}_i) dM^\top \mathbb{A}_i^\top \mathbb{A}_i d\gamma = 2\pi_1 \text{Var}(\mathbb{P}_1) d(\beta \gamma_1) d\gamma_1 + 2(1 - \pi_1) \text{Var}(\mathbb{P}_{i \neq 1}) d\left(\frac{1}{2}(\gamma_2 + \gamma_3)\beta\right) d\left(\frac{1}{2}(\gamma_2 + \gamma_3)\right) \quad (151)$$

Writing into the matrix form, we get

$$\begin{aligned} B_k^{(1)} &= -2 \begin{pmatrix} \pi_1 \text{Var}(\mathbb{P}_1)_{k-3} \beta & 0 & 0 & \frac{1}{2} \pi_1 \gamma_1 \text{Var}(\mathbb{P}_1)_{k-3} \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} \pi_1 \gamma_1 \text{Var}(\mathbb{P}_1)_{k-3}^\top & 0 & 0 & 0 \end{pmatrix} \\ &\quad - 2 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{4} (1 - \pi_1) \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & \frac{1}{4} (1 - \pi_1) \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \\ 0 & \frac{1}{4} (1 - \pi_1) \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & \frac{1}{4} (1 - \pi_1) \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \\ 0 & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3}^\top & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3}^\top & 0 \end{pmatrix} \end{aligned} \quad (152)$$

Recall that $d^2 \mathbb{P}_i = 2d\mathbb{A}_i dM \text{Var}(\mathbb{P}_i) + \mathbb{A}_i d^2 M \text{Var}(\mathbb{P}_i) + d(\mathbb{A}_i M) d \text{Var}(\mathbb{P}_i)$, we have

$$\sum_i \pi_i d^2 \mathbb{P}_i^\top \mathbb{A}_i \gamma = \sum_i \pi_i (2d\mathbb{A}_i dM \text{Var}(\mathbb{P}_i) + \mathbb{A}_i d^2 M \text{Var}(\mathbb{P}_i) + d(\mathbb{A}_i M) d \text{Var}(\mathbb{P}_i))^\top \mathbb{A}_i \gamma \quad (153)$$

The first term contributes to

$$B_k^{(2)} = -\frac{1}{2} (1 - \pi_1) \eta \gamma_{i \neq 1}^2 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & -\text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & 0 \\ 0 & -\text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \beta & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (154)$$

The second term contributes to

$$\begin{aligned}
 B_k^{(3)} = & -2\pi_1\gamma_1 \begin{pmatrix} 0 & 0 & 0 & \frac{1}{2} \text{Var}(\mathbb{P}_1)_{k-3} \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} \text{Var}(\mathbb{P}_1)_{k-3}^\top & 0 & 0 & 0 \end{pmatrix} \\
 & - 2(1 - \pi_1)\gamma_{i \neq 1} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \\ 0 & 0 & 0 & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3} \\ 0 & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3}^\top & \frac{1}{4} \text{Var}(\mathbb{P}_{i \neq 1})_{k-3}^\top & 0 \end{pmatrix}
 \end{aligned} \tag{155}$$

We consider the action of the last term on $d\beta^2$:

$$\begin{aligned}
 \begin{pmatrix} B_4^{(4)} \\ B_5^{(4)} \\ B_6^{(4)} \end{pmatrix} (d\beta, d\beta) = & -\pi_1\gamma_1^3 (\text{diag}(d\beta^\top \text{Var}(\mathbb{P}_1)) - \mathbb{P}_1^\top (d\beta^\top \text{Var}(\mathbb{P}_1)) - \text{Var}(\mathbb{P}_1) d\beta \mathbb{P}_1) d\beta \\
 & - (1 - \pi_1)\gamma_{i \neq 1}^3 (\text{diag}(d\beta^\top \text{Var}(\mathbb{P}_{i \neq 1})) - \mathbb{P}_{i \neq 1}^\top (d\beta^\top \text{Var}(\mathbb{P}_{i \neq 1})) - \text{Var}(\mathbb{P}_{i \neq 1}) d\beta \mathbb{P}_{i \neq 1}) d\beta
 \end{aligned} \tag{156}$$

Lemma D.13 (Kernel-equation identities used in Theorem D.3). *Let $\theta = q_2 y_2$ be the leading-order reduction (since $\zeta(0, \delta) \sim \delta q_2$). Then*

$$\frac{1}{2} Q_K^\top B(q_2 y_2, q_2 y_2) = \frac{1}{\sqrt{\|\gamma\|^2 + \|\beta\|^2}} \begin{pmatrix} 0 \\ 0 \\ \pi_1 \gamma_1^2 \mathbb{P}_{1,2} + (1 - \pi_1) \gamma_2^2 \mathbb{P}_{i \neq 1,2} \end{pmatrix} y_2^2.$$

Proof. We first calculate $B(q_2 y_2, q_2 y_2)$, and then calculate its projection onto the kernel basis. We calculate the cases where $1 \leq k \leq 3$ and $4 \leq k \leq 6$ respectively, and then combine them into the form we want.

1. The computation of the cases where $1 \leq k \leq 3$. From Eq. (139), we get

$$\begin{pmatrix} B_1^{(1)}(q_2 y_2, q_2 y_2) \\ B_2^{(1)}(q_2 y_2, q_2 y_2) \\ B_3^{(1)}(q_2 y_2, q_2 y_2) \end{pmatrix} = -\frac{1}{2} \begin{pmatrix} 4\pi_1 \gamma_1 \mathbb{P}_{1,2} \\ 2(1 - \pi_1) \gamma_{i \neq 1} \mathbb{P}_{i \neq 1,2} \\ 2(1 - \pi_1) \gamma_{i \neq 1} \mathbb{P}_{i \neq 1,2} \end{pmatrix} y_2^2 \tag{157}$$

For l from 2 to 4 and $l = 6$, their contribution vanishes. For $l = 5$, the contribution is

$$\begin{pmatrix} B_1^{(5)}(q_2 y_2, q_2 y_2) \\ B_2^{(5)}(q_2 y_2, q_2 y_2) \\ B_3^{(5)}(q_2 y_2, q_2 y_2) \end{pmatrix} = -\frac{1}{2} \begin{pmatrix} 2\pi_1 \gamma_1^2 (\mathbb{P}_{1,2} \beta_2 - \mathbb{P}_{1,2} \mathbb{P}_1 \beta) \\ (1 - \pi_1) \gamma_{i \neq 1}^2 (\mathbb{P}_{i \neq 1,2} \beta_2 - \mathbb{P}_{i \neq 1,2} \mathbb{P}_{i \neq 1} \beta) \\ (1 - \pi_1) \gamma_{i \neq 1}^2 (\mathbb{P}_{i \neq 1,2} \beta_2 - \mathbb{P}_{i \neq 1,2} \mathbb{P}_{i \neq 1} \beta) \end{pmatrix} y_2^2 \tag{158}$$

2. The computation of the cases where $4 \leq k \leq 6$. For $l = 1, 2, 3$, their contribution vanishes. Then contribution of $l = 4$ case is

$$-\frac{1}{2} \pi_1 \gamma_1^3 \begin{pmatrix} 0 \\ 0 \\ 0 \\ \begin{pmatrix} 0 \\ \mathbb{P}_{1,2} \\ \mathbb{P}_{1,2} \end{pmatrix} - 2\mathbb{P}_{1,2} \mathbb{P}_1^\top \end{pmatrix} - \frac{1}{2} (1 - \pi_1) \gamma_{i \neq 1}^3 \begin{pmatrix} 0 \\ 0 \\ 0 \\ \begin{pmatrix} 0 \\ \mathbb{P}_{i \neq 1,2} \\ \mathbb{P}_{i \neq 1,2} \end{pmatrix} - 2\mathbb{P}_{i \neq 1,2} \mathbb{P}_{i \neq 1}^\top \end{pmatrix} \tag{159}$$

Sum them up, we get $B(q_2 y_2, q_2 y_2)$. Then by direct computation, we get the projection onto the kernel directions.

$$\frac{1}{2} Q_K^\top B(q_2 y_2, q_2 y_2) = \frac{1}{4\sqrt{\|\gamma\|^2 + \|\beta\|^2}} \begin{pmatrix} 0 \\ 0 \\ 4\pi_1 \gamma_1^2 \mathbb{P}_{1,2} + 4(1 - \pi_1) \gamma_2^2 \mathbb{P}_{i \neq 1,2} \end{pmatrix} y_2^2. \tag{160}$$

□

D.3.7. COMPUTATION OF SECOND ORDER DERIVATIVE f_2

The calculation about $f_2 = \partial_\delta^2(-\nabla\mathcal{L})$ is summarized by following lemma.

Lemma D.14. *The second derivative of $-\nabla\mathcal{L}$ with respect to perturbation parameter δ vanishes, i.e. $f_2 = 0$.*

Proof. We compute $\partial_\delta^2 \frac{\partial\mathcal{L}}{\partial M}$ and $\partial_\delta^2 \frac{\partial\mathcal{L}}{\partial\Phi}$ as follows.

1. The computation of $\partial_\delta^2 \frac{\partial\mathcal{L}}{\partial M}$. By definition,

$$\begin{aligned} \partial_\delta^2 \frac{\partial\mathcal{L}}{\partial M} &= \partial_\delta^2 \left(- \sum_i \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) \right) \\ &= \partial_\delta \left(- \sum_i \partial_\delta \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i \partial_\delta \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) \right) \\ &= - \sum_i \partial_\delta^2 \pi_i \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - 2 \sum_i \partial_\delta \pi_i \partial_\delta \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - 2 \sum_i \partial_\delta \pi_i \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) \\ &\quad - \sum_i \pi_i \partial_\delta^2 \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - 2 \sum_i \pi_i \partial_\delta \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) - \sum_i \pi_i \mathbb{A}_i^\top \partial_\delta^2 (P_i - \mathbb{P}_i) \\ &= -2 \sum_i \partial_\delta \pi_i \partial_\delta \mathbb{A}_i^\top (P_i - \mathbb{P}_i) - 2 \sum_i \partial_\delta \pi_i \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) - 2 \sum_i \pi_i \partial_\delta \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) \end{aligned}$$

The last equality uses the second order derivative of π_i , \mathbb{A}_i , and $P_i - \mathbb{P}_i$ with respect to δ vanishes. Using the fact that $\partial_\delta \mathbb{A}_i$ is of the shape like $(0, a, -a)$ and $(P_i - \mathbb{P}_i)\beta = 0$. The contribution of the first term vanishes. Similarly, the third term vanishes. For the second term,

$$-2 \sum_i \partial_\delta \pi_i \mathbb{A}_i^\top \partial_\delta (P_i - \mathbb{P}_i) = -2 \mathbb{A}_{i \neq 1}^\top (\partial_\delta P_2 - \partial_\delta P_3) = 0$$

Thus, this term makes no contribution.

2. The computation of $\partial_\delta^2 \frac{\partial\mathcal{L}}{\partial\Phi}$. By definition, $\partial_\delta^2 \frac{\partial\mathcal{L}}{\partial\Phi} = \partial_\delta^2 \left(- \sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) \right)$. In particular,

$$\begin{aligned} &\partial_\delta \left(- \sum_i \partial_\delta \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i \partial_\delta (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \partial_\delta \text{Var}(\mathbb{A}_i) \right) \\ &= - \sum_i \partial_\delta^2 \pi_i e_i (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - 2 \sum_i \partial_\delta \pi_i e_i \partial_\delta (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - 2 \sum_i \partial_\delta \pi_i e_i (P_i - \mathbb{P}_i) M^\top \partial_\delta \text{Var}(\mathbb{A}_i) \\ &\quad - \sum_i \pi_i e_i \partial_\delta^2 (P_i - \mathbb{P}_i) M^\top \text{Var}(\mathbb{A}_i) - 2 \sum_i \pi_i e_i \partial_\delta (P_i - \mathbb{P}_i) M^\top \partial_\delta \text{Var}(\mathbb{A}_i) - \sum_i \pi_i e_i (P_i - \mathbb{P}_i) M^\top \partial_\delta^2 \text{Var}(\mathbb{A}_i) \end{aligned}$$

By direct computation, this term is zero.

□

D.3.8. STABILITY ON THE KERNEL DIRECTIONS

To account for stability on the manifold, we need to calculate the perturbed Hessian matrix. By the expansion of $-\nabla\mathcal{L}$, we get

$$-\nabla_\theta^2 \mathcal{L}(\theta, \delta) = J_0 + \frac{1}{2} \nabla_\theta B(\theta, \theta) + \delta J_1 + \mathcal{O}(\delta^2)$$

Substitute $\theta = q_2 y_2$ into the expression, we get the perturbed hessian matrix

Lemma D.15 (Perturbed hessian matrix). *The expression of the perturbed hessian matrix is*

$$J_{pert} = J_0 + H_1 + \mathcal{O}(\delta^2), \tag{161}$$

where

$$H_1 = -c\delta \left(\begin{pmatrix} 0 & B \\ B^\top & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & C \end{pmatrix} \right) + \delta J_1, \quad (162)$$

$$\text{in which } c = \frac{\lambda\gamma_{i \neq 1} + (1-\lambda)(\pi_1\gamma_1 + (1-\pi_1)\gamma_{i \neq 1})}{\pi_1\gamma_1^2\mathbb{P}_{1,2} + (1-\pi_1)\gamma_{i \neq 1}^2\mathbb{P}_{i \neq 1,2}},$$

$$B = \begin{pmatrix} 0 & 2\pi_1\gamma_1\mathbb{P}_{1,2} & -2\pi_1\gamma_1\mathbb{P}_{1,2} \\ 0 & (1-\pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2} & -(1-\pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2} \\ 0 & (1-\pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2} & -(1-\pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2} \end{pmatrix} \\ + \begin{pmatrix} 0 & \pi_1\gamma_1^2(\mathbb{P}_{1,2}\beta_2 - \mathbb{P}_{1,2}(\mathbb{P}_1\beta)) & -\pi_1\gamma_1^2(\mathbb{P}_{1,2}\beta_2 - \mathbb{P}_{1,2}(\mathbb{P}_1\beta)) \\ 0 & \frac{1}{2}(1-\pi_1)\gamma_{i \neq 1}^2(\mathbb{P}_{i \neq 1,2}\beta_2 - \mathbb{P}_{i \neq 1,2}(\mathbb{P}_{i \neq 1}\beta)) & -\frac{1}{2}(1-\pi_1)\gamma_{i \neq 1}^2(\mathbb{P}_{i \neq 1,2}\beta_2 - \mathbb{P}_{i \neq 1,2}(\mathbb{P}_{i \neq 1}\beta)) \\ 0 & \frac{1}{2}(1-\pi_1)\gamma_{i \neq 1}^2(\mathbb{P}_{i \neq 1,2}\beta_2 - \mathbb{P}_{i \neq 1,2}(\mathbb{P}_{i \neq 1}\beta)) & -\frac{1}{2}(1-\pi_1)\gamma_{i \neq 1}^2(\mathbb{P}_{i \neq 1,2}\beta_2 - \mathbb{P}_{i \neq 1,2}(\mathbb{P}_{i \neq 1}\beta)) \end{pmatrix},$$

and

$$C = \pi_1\gamma_1^3 \begin{pmatrix} 0 & 0 \\ 0 & \begin{pmatrix} 0 & -\mathbb{P}_{1,1}\mathbb{P}_{1,2} & \mathbb{P}_{1,1}\mathbb{P}_{1,2} \\ -\mathbb{P}_{1,1}\mathbb{P}_{1,2} & \mathbb{P}_{1,2} - 2\mathbb{P}_{1,2}^2 & 0 \\ \mathbb{P}_{1,1}\mathbb{P}_{1,2} & 0 & -(\mathbb{P}_{1,2} - 2\mathbb{P}_{1,2}^2) \end{pmatrix} \end{pmatrix} \\ + (1-\pi_1)\gamma_{i \neq 1}^3 \begin{pmatrix} 0 & 0 \\ 0 & \begin{pmatrix} 0 & -\mathbb{P}_{i \neq 1,1}\mathbb{P}_{i \neq 1,2} & \mathbb{P}_{i \neq 1,1}\mathbb{P}_{i \neq 1,2} \\ -\mathbb{P}_{i \neq 1,1}\mathbb{P}_{i \neq 1,2} & \mathbb{P}_{i \neq 1,2} - 2\mathbb{P}_{i \neq 1,2}^2 & 0 \\ \mathbb{P}_{i \neq 1,1}\mathbb{P}_{i \neq 1,2} & 0 & -(\mathbb{P}_{i \neq 1,2} - 2\mathbb{P}_{i \neq 1,2}^2) \end{pmatrix} \end{pmatrix}$$

Proof. We just need to compute $B(q_2y_2)$ and substitute y_2 as the solution of the range equation. Similar to previous computation, we divide into the cases of $1 \leq k \leq 3$ and $4 \leq k \leq 6$.

The cases of $1 \leq k \leq 3$.

1. The contribution from $l = 1$. Using the matrix form defined in Eq. (139), we get

$$B_1^{(1)}(q_2y_2) = -\frac{y_2}{\sqrt{2}}(0, 0, 0, 0, 2\pi_1\gamma_1\mathbb{P}_{1,2}, -2\pi_1\gamma_1\mathbb{P}_{1,2}) \\ B_2^{(1)}(q_2y_2) = -\frac{y_2}{\sqrt{2}}(0, 0, 0, 0, (1-\pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2}, -(1-\pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2}) \\ B_3^{(1)}(q_2y_2) = B_2^{(1)}(q_2y_2)$$

2. The contributions from $l = 2, 3, 4, 6$ vanish.

3. The contribution from $l = 5$. Using the matrix defined in Eq. (145), we get

$$c_3 \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} = -\pi_1\gamma_1^2(0, \mathbb{P}_{1,2}\beta_2 - \mathbb{P}_{1,2}(\mathbb{P}_1\beta), -\mathbb{P}_{1,2}\beta_2 - \mathbb{P}_{1,2}(\mathbb{P}_1\beta)),$$

which implies

$$B_1^{(5)}(q_2y_2) = -\frac{y_2}{\sqrt{2}}\pi_1\gamma_1^2(0, 0, 0, 0, \mathbb{P}_{1,2}\beta_2 - \mathbb{P}_{1,2}(\mathbb{P}_1\beta), -(\mathbb{P}_{1,2}\beta_2 - \mathbb{P}_{1,2}(\mathbb{P}_1\beta)))$$

Similarly,

$$B_2^{(5)}(q_2y_2) = B_3^{(5)}(q_2y_2) = -\frac{y_2}{\sqrt{2}}\frac{1}{2}(1-\pi_1)\gamma_{i \neq 1}^2(0, 0, 0, 0, \mathbb{P}_{i \neq 1,2}\beta_2 - \mathbb{P}_{i \neq 1,2}(\mathbb{P}_{i \neq 1}\beta), -\mathbb{P}_{i \neq 1,2}\beta_2 - \mathbb{P}_{i \neq 1,2}(\mathbb{P}_{i \neq 1}\beta))$$

The cases of $4 \leq k \leq 6$.

1. The contribution from $l = 1$. By direct computation,

$$B_4^{(1)}(q_2 y_2) = (0, 0, 0, 0, 0, 0)$$

$$B_5^{(1)}(q_2 y_2) = -\frac{y_2}{\sqrt{2}}(\pi_1 \gamma_1 \mathbb{P}_{1,2}, \frac{1}{2}(1 - \pi_1) \gamma_{i \neq 1} \mathbb{P}_{i \neq 1,2}, \frac{1}{2}(1 - \pi_1) \gamma_{i \neq 1} \mathbb{P}_{i \neq 1,2}, 0, 0, 0)$$

$$B_6^{(1)}(q_2 y_2) = -B_5^{(1)}$$

2. The contribution from $l = 2$. This term vanishes.

3. The contribution from $l = 3$. This term makes the same contribution as the first term.

4. The contribution from $l = 4$. By definition,

$$-\sum_i \pi_i (\mathbb{d}(\mathbb{A}_i M) \mathbb{d} \text{Var}(\mathbb{P}_i))^\top \mathbb{A}_i \gamma = -\pi_1 \gamma_1 (\mathbb{A}_1 \mathbb{d} M \mathbb{d} \text{Var}(\mathbb{P}_1))^\top - (1 - \pi_1) \gamma_{i \neq 1} (\mathbb{A}_{i \neq 1} \mathbb{d} M \mathbb{d} \text{Var}(\mathbb{P}_{i \neq 1}))^\top$$

Take the first term as an example, the cross term in $(\mathbb{A}_1 \mathbb{d} M \mathbb{d} \text{Var}(\mathbb{P}_1))^\top$ is

$$\begin{aligned} & \gamma_1 \mathbb{d} \gamma_1 (\text{diag}(\beta^\top \text{Var}(\mathbb{P}_1)) - \mathbb{P}_1^\top \beta^\top \text{Var}(\mathbb{P}_1) - \text{Var}(\mathbb{P}_1) \beta \mathbb{P}_1) \mathbb{d} \beta \\ & + \gamma_1 \mathbb{d} \gamma_1 (\text{diag}(\mathbb{d} \beta^\top \text{Var}(\mathbb{P}_1)) - \mathbb{P}_1^\top \mathbb{d} \beta^\top \text{Var}(\mathbb{P}_1) - \text{Var}(\mathbb{P}_1) \mathbb{d} \beta \mathbb{P}_1) \beta \end{aligned}$$

Write in entry form, for $4 \leq k \leq 6$, we get

$$\begin{aligned} & \gamma_1 \mathbb{d} \gamma_1 ((\beta^\top \text{Var}(\mathbb{P}_1))_k \mathbb{d} \beta_k - \mathbb{P}_{1,k} \beta^\top \text{Var}(\mathbb{P}_1) \mathbb{d} \beta - (\text{Var}(\mathbb{P}_1) \beta)_k \mathbb{P}_1 \mathbb{d} \beta) \\ & + \gamma_1 \mathbb{d} \gamma_1 (\mathbb{d} \beta^\top \text{Var}(\mathbb{P}_1)_k \beta_k - \mathbb{P}_{1,k} \mathbb{d} \beta^\top \text{Var}(\mathbb{P}_1) \beta - \text{Var}(\mathbb{P}_1)_k \mathbb{d} \beta \mathbb{P}_1 \beta) \end{aligned}$$

Writing into the matrix form, we get

$$\begin{aligned} & \begin{pmatrix} 0 & \frac{1}{2}(\beta^\top \text{Var}(\mathbb{P}_1))_k E_{1,k} \\ \frac{1}{2}(\beta^\top \text{Var}(\mathbb{P}_1))_k E_{1,k}^\top & 0 \end{pmatrix} + \begin{pmatrix} 0 & \frac{1}{2} \beta_k e_k \text{Var}(\mathbb{P}_1)_k \\ \frac{1}{2} \beta_k \text{Var}(\mathbb{P}_1)_k e_k^\top & 0 \end{pmatrix} \\ & - \mathbb{P}_{1,k} \begin{pmatrix} 0 & e_1 \beta^\top \text{Var}(\mathbb{P}_1) \\ \text{Var}(\mathbb{P}_1) \beta e_1^\top & 0 \end{pmatrix} - \begin{pmatrix} 0 & \frac{1}{2}(\text{Var}(\mathbb{P}_1) \beta)_k e_1 \mathbb{P}_1 \\ \frac{1}{2}(\text{Var}(\mathbb{P}_1) \beta)_k \mathbb{P}_1 e_1^\top & 0 \end{pmatrix} \\ & - \begin{pmatrix} 0 & \frac{1}{2}(\mathbb{P}_1 \beta) e_1 \text{Var}(\mathbb{P}_1)_k \\ \frac{1}{2}(\mathbb{P}_1 \beta) \text{Var}(\mathbb{P}_1)_k e_1^\top & 0 \end{pmatrix} \end{aligned}$$

Multiplying the matrix form by $(0, 0, 0, 0, 1, -1)$ on the right, we get

$$B_4(q_2 y_2) = 0$$

$$\begin{aligned} B_5(q_2 y_2) &= -\frac{y_2}{\sqrt{2}} \pi_1 \gamma_1^2 \left(\frac{1}{2}(\beta^\top \text{Var}(\mathbb{P}_1))_2 + \frac{1}{2} \beta_2 \mathbb{P}_{1,2} - \frac{1}{2}(\mathbb{P}_1 \beta) \mathbb{P}_{1,2}, 0, \dots, 0 \right) \\ &= -\frac{y_2}{\sqrt{2}} \pi_1 \gamma_1^2 (\beta_2 \mathbb{P}_{1,2} - (\mathbb{P}_1 \beta) \mathbb{P}_{1,2}, 0, \dots, 0) \end{aligned}$$

$$B_6(q_2 y_2) = -B_5(q_2 y_2)$$

Similarly, the cross term in $-(1 - \pi_1) \gamma_{i \neq 1} (\mathbb{A}_{i \neq 1} \mathbb{d} M \mathbb{d} \text{Var}(\mathbb{P}_{i \neq 1}))^\top$ contributes to

$$B_4(q_2 y_2) = 0$$

$$B_5(q_2 y_2) = -\frac{y_2}{\sqrt{2}} \frac{1}{2} (1 - \pi_1) \gamma_{i \neq 1}^2 (0, \beta_2 \mathbb{P}_{i \neq 1,2} - (\mathbb{P}_{i \neq 1} \beta) \mathbb{P}_{i \neq 1,2}, \beta_2 \mathbb{P}_{i \neq 1,2} - (\mathbb{P}_{i \neq 1} \beta) \mathbb{P}_{i \neq 1,2}, 0, \dots, 0)$$

$$B_6(q_2 y_2) = -B_5(q_2 y_2)$$

Finally, we compute the contribution from quadratic form. Take the first term as an example,

$$\gamma_1^2 (\text{diag}(\mathbb{d} \beta^\top \text{Var}(\mathbb{P}_1)) - \mathbb{P}_1^\top \text{Var} \mathbb{P}_1 - \text{Var}(\mathbb{P}_1) \mathbb{d} \beta \mathbb{P}_1) \mathbb{d} \beta$$

Writing into the matrix form, we get

$$B_k = \begin{pmatrix} 0 & 0 \\ 0 & \frac{1}{2} (e_k \text{Var}(\mathbb{P}_1)_k + \text{Var}(\mathbb{P}_1)_k^\top e_k^\top) \end{pmatrix} - \begin{pmatrix} 0 & 0 \\ 0 & \mathbb{P}_{1,k} \text{Var}(\mathbb{P}_1) \end{pmatrix} - \begin{pmatrix} 0 & 0 \\ 0 & \frac{1}{2} (\mathbb{P}_1^\top \text{Var}(\mathbb{P}_1)_k + \text{Var}(\mathbb{P}_1)_k^\top \mathbb{P}_1) \end{pmatrix}$$

By direct computation, we get the contribution to the perturbed hessian is

$$\begin{aligned} & -\frac{y_2}{\sqrt{2}} \pi_1 \gamma_1^3 \begin{pmatrix} 0 & 0 \\ 0 & \begin{pmatrix} 0 & -\mathbb{P}_{1,1} \mathbb{P}_{1,2} & \mathbb{P}_{1,1} \mathbb{P}_{1,2} \\ -\mathbb{P}_{1,1} \mathbb{P}_{1,2} & \mathbb{P}_{1,2} - 2\mathbb{P}_{1,2}^2 & 0 \\ \mathbb{P}_{1,1} \mathbb{P}_{1,2} & 0 & -(\mathbb{P}_{1,2} - 2\mathbb{P}_{1,2}^2) \end{pmatrix} \end{pmatrix} \\ & -\frac{y_2}{\sqrt{2}} (1 - \pi_1) \gamma_{i \neq 1}^3 \begin{pmatrix} 0 & 0 \\ 0 & \begin{pmatrix} 0 & -\mathbb{P}_{i \neq 1,1} \mathbb{P}_{i \neq 1,2} & \mathbb{P}_{i \neq 1,1} \mathbb{P}_{i \neq 1,2} \\ -\mathbb{P}_{i \neq 1,1} \mathbb{P}_{i \neq 1,2} & \mathbb{P}_{i \neq 1,2} - 2\mathbb{P}_{i \neq 1,2}^2 & 0 \\ \mathbb{P}_{i \neq 1,1} \mathbb{P}_{i \neq 1,2} & 0 & -(\mathbb{P}_{i \neq 1,2} - 2\mathbb{P}_{i \neq 1,2}^2) \end{pmatrix} \end{pmatrix} \end{aligned}$$

We obtain the result by summing all non-zero terms. \square

Lemma D.16 (Vanishing of the first-order perturbation on the kernel). *On the symmetric rank-one manifold ($\gamma_2 = \gamma_3$ and $\beta_2 = \beta_3$), we have $Q_K^\top H_1 Q_K = 0$.*

Proof. We write any $z \in \mathbb{R}^6$ as $z = (z_\gamma, z_\beta)$ with $z_\gamma, z_\beta \in \mathbb{R}^3$. Let

$$s := (0, 1, -1)^\top, \quad u := (1, 1, 1)^\top.$$

Then $k_{1,\gamma} = \frac{1}{\sqrt{2}}s$, $k_{1,\beta} = 0$; $k_{2,\gamma} = 0$, $k_{2,\beta} = \frac{1}{\sqrt{3}}u$; and $k_{3,\gamma} \propto -\gamma$, $k_{3,\beta} \propto \beta$.

Step 1: the off-diagonal block with respect to B For any $x = (x_\gamma, x_\beta)$ and $y = (y_\gamma, y_\beta)$,

$$x^\top \begin{pmatrix} 0 & B \\ B^\top & 0 \end{pmatrix} y = x_\gamma^\top B y_\beta + x_\beta^\top B^\top y_\gamma.$$

From Eq. (171), B has the structural identities

$$B_{\cdot,1} = 0, \quad B_{\cdot,3} = -B_{\cdot,2}, \quad B_{2,\cdot} = B_{3,\cdot}.$$

Hence

$$Bu = B(e_1 + e_2 + e_3) = 0, \quad B\beta = \beta_1 B_{\cdot,1} + \beta_2 (B_{\cdot,2} + B_{\cdot,3}) = 0 \quad (\text{since } \beta_2 = \beta_3),$$

and

$$s^\top B = (0, 1, -1)B = B_{2,\cdot} - B_{3,\cdot} = 0.$$

Combining these, every matrix element $k_i^\top \begin{pmatrix} 0 & B \\ B^\top & 0 \end{pmatrix} k_j$ vanishes: either a factor $k_{1,\beta} = 0$ appears, or a factor $s^\top B = 0$ appears, or a factor $Bu = 0 / B\beta = 0$ appears.

Step 2: the lower-right block C . Here

$$x^\top \begin{pmatrix} 0 & 0 \\ 0 & C \end{pmatrix} y = x_\beta^\top C y_\beta.$$

Each summand of C in Eq. (172) has the pattern

$$C^{(\ell)} = \begin{pmatrix} 0 & -p_\ell & p_\ell \\ -p_\ell & a_\ell & 0 \\ p_\ell & 0 & -a_\ell \end{pmatrix} \quad \text{for some } (p_\ell, a_\ell), \quad C = \sum_{\ell} \omega_\ell C^{(\ell)}.$$

A direct expansion gives, for any $x, y \in \mathbb{R}^3$,

$$x^\top C^{(\ell)} y = p_\ell (x_3 - x_2) y_1 + p_\ell x_1 (y_3 - y_2) + a_\ell (x_2 y_2 - x_3 y_3).$$

Therefore, if $x_2 = x_3$ and $y_2 = y_3$, then $x^\top C^{(\ell)} y = 0$ and hence $x^\top C y = 0$. On the symmetric manifold we have

$(k_{2,\beta})_2 = (k_{2,\beta})_3$ and $(k_{3,\beta})_2 = (k_{3,\beta})_3$ (since $\beta_2 = \beta_3$), while $k_{1,\beta} = 0$. Thus, $k_i^\top \begin{pmatrix} 0 & 0 \\ 0 & C \end{pmatrix} k_j = 0$ for all i, j .

Step 3: the J_1 term. By Proposition 7.4,

$$J_1 = \tilde{J}_1 + \begin{pmatrix} 0 & A \\ A^\top & 0 \end{pmatrix},$$

where A in Eq. (109) satisfies the same cancellation identities as B :

$$A_{\cdot,1} = 0, \quad A_{\cdot,3} = -A_{\cdot,2}, \quad A_{2,\cdot} = A_{3,\cdot}.$$

Hence $Au = 0$, $A\beta = 0$ (since $\beta_2 = \beta_3$), and $s^\top A = 0$. Repeating Step 1 with B replaced by A , we get

$$Q_K^\top \begin{pmatrix} 0 & A \\ A^\top & 0 \end{pmatrix} Q_K = 0.$$

For the remaining part \tilde{J}_1 , the explicit computation shows that its action on the kernel directions has no kernel component, i.e. $Q_K^\top \tilde{J}_1 Q_K = 0$. Therefore $Q_K^\top J_1 Q_K = 0$.

Conclusion. Combining Step 1–3 yields $Q_K^\top H_1 Q_K = 0$. □

D.3.9. FAST TRANSVERSE INSTABILITY: $\Theta(\delta)$ EIGENVALUE

Lemma D.17 (A transverse eigenvalue of order $\Theta(\delta)$). *Let $c = \frac{\lambda\gamma_{i \neq 1} + (1-\lambda)(\pi_1\gamma_1 + (1-\pi_1)\gamma_{i \neq 1})}{\pi_1\gamma_1^2\mathbb{P}_{1,2} + (1-\pi_1)\gamma_{i \neq 1}^2\mathbb{P}_{i \neq 1,2}}$ and assume that*

$$c(1 - \pi_1)\gamma_{i \neq 1}\mathbb{P}_{i \neq 1,2} - (\lambda + (1 - \pi_1)(1 - \lambda)) \neq 0. \quad (163)$$

At the perturbed point, $\partial\mathcal{L}/\partial\Phi = \mathcal{O}(\delta^2)$ while $\partial\mathcal{L}/\partial M = \Theta(\delta)$. Consequently, linearizing the full dynamics (Eq. (107)) yields a transverse positive eigenvalue of size $\Theta(\delta)$.

Proof. Since $d\frac{\partial\mathcal{L}}{\partial\Phi} = 0$ and $\partial_\delta\frac{\partial\mathcal{L}}{\partial\Phi} = 0$, we get $\frac{\partial\mathcal{L}}{\partial\Phi} = \mathcal{O}(\delta^2)$ after perturbation.

Next we compute perturbed $\frac{\partial\mathcal{L}}{\partial M}$. By Lemma D.5, we get the new term is

$$c\pi_1\gamma_1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (0, \mathbb{P}_{1,2}, -\mathbb{P}_{1,2}) + c(1 - \pi_1)\gamma_{i \neq 1} \begin{pmatrix} 0 \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} (0, \mathbb{P}_{i \neq 1,2}, -\mathbb{P}_{i \neq 1,2})$$

where $c = \frac{\lambda\gamma_{i \neq 1} + (1-\lambda)(\pi_1\gamma_1 + (1-\pi_1)\gamma_{i \neq 1})}{\pi_1\gamma_1^2\mathbb{P}_{1,2} + (1-\pi_1)\gamma_{i \neq 1}^2\mathbb{P}_{i \neq 1,2}} \delta$.

However, from Eq. (119), we find that

$$\frac{\partial}{\partial\delta} \frac{\partial\mathcal{L}}{\partial M} = - \begin{pmatrix} 0 \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} (0, \lambda, -\lambda) - \begin{pmatrix} \pi_1 \\ \frac{1}{2}(1 - \pi_1) \\ \frac{1}{2}(1 - \pi_1) \end{pmatrix} (0, 1 - \lambda, -(1 - \lambda)).$$

The two terms cannot cancel each other out by our assumption (A parameter that does not meet the condition is a zero test set), thus resulting in an $\Theta(\delta)$ term. □

E. Detailed Experiment Setup

E.1. Detailed Synthetic Experiment Setup

Dataset To better induce an exponential decay in the stationary distribution, and to more clearly illustrate the phase transition, we adopt an exponentially decaying form

$$\pi_0 = (1/2, 1/4, \dots, 1/2^d), \quad \pi = \pi_0 / \|\pi_0\|.$$

Following (Makkuva et al., 2025), diagonal dominant transition matrices are unfavorable local minima during optimization, we set $\lambda = 0.8$ to ensure diagonal dominance. For a fixed sequence length of 20, we sample 100,000 sequences $\{X_i\}_{i=1}^{100,000}$ from the resulting Markov chain, and use the last token $X_i[-1]$ as the training label. By the Markov property, this label is completely determined by the second-to-last token $X_i[-2]$. Accordingly, we group both the training and test sets by the value of $X_i[-2]$, denoting the group indexed by state k as S_k .

Model We follow exactly the model specification in Def. 2.2, with embedding dimension $m = 256$. Since our theoretical analysis is derived under small initialization, we adopt the initialization scheme of (Zhang et al., 2024b; 2025b), initializing each weight independently as $\mathcal{N}(0, 1/m^2)$.

Training We train the model using the Adam optimizer with a fixed learning rate of 1.5×10^{-4} and do not use any learning-rate scheduler.

E.2. Analysis Tools

Condensation Heatmap To quantify parameter condensation, we compute the pairwise cosine similarity between the input-weight vectors of neurons in the weight matrix W . Specifically, for the i -th and j -th neurons, we define

$$C(i, j) = \frac{W[i, :] \cdot W[j, :]}{\|W[i, :]\|_2 \|W[j, :]\|_2}.$$

For clearer visualization, we permute the rows and columns of the similarity matrix C and display the reordered matrix in Fig. 2(A).

Embedding Visualization Let $W_{0,t}$ denote the embedding parameters at training epoch t for $t = 0, \dots, T$. We form the collection of embedding snapshots $\{W_{0,0}, \dots, W_{0,T}\}$ and apply principal component analysis (PCA) to obtain the leading eigen-directions \hat{e}_1 and \hat{e}_2 . We then project the embedding vectors onto \hat{e}_1 and \hat{e}_2 to produce the two-dimensional visualization shown in Fig. 2(B).