

Equivariant Structured Positional Rotations: Complete Algebraic Derivation

Abstract

This document provides a complete, boxed derivation of the relative-position property for equivariant structured positional rotations in attention mechanisms. Each step is recorded as a labelled statement $[T]$ with explicit dependencies so that the argument can be followed mechanically.

Contents

1 Tier 0: Linear-Algebra Primitives	2
2 Tier 1: Structured Generators	4
3 Tier 2: Joint Block Structure	5
4 Tier 3: Cayley Post-Rotations	6
5 Tier 4: Attention Queries, Keys, and Scores	6
6 Tier 5: Stability Estimates	7
7 Tier 6: Conclusions	8
8 Tier 6': Quantified Robustness	8
9 Tier 7: Positive Random Features and Variance Bounds	10
10 Tier 8': Tight Assumptions and OOD Dominance	13

1 Tier 0: Linear-Algebra Primitives

$$\boxed{[T0.1 : \mathbb{N}, \text{Idx}]} \quad \begin{aligned} \mathbb{N} &\stackrel{\mathcal{Q}}{=} \{1, 2, 3, \dots\} & [\text{deps} : \emptyset] \\ d_h, d_c, D &\in \mathbb{N} \end{aligned}$$

$$H \stackrel{\mathcal{Q}}{=} \{1, \dots, d_h\}, \quad K \stackrel{\mathcal{Q}}{=} \{1, \dots, d_c\}, \quad J_D \stackrel{\mathcal{Q}}{=} \{1, \dots, D\}$$

$$\boxed{[T0.2 : I_n, \delta_{ij}]} \quad \delta_{ij} \stackrel{\mathcal{Q}}{=} \begin{cases} 1, & i = j, \\ 0, & i \neq j, \end{cases} \quad [\text{deps} : T0.1]$$

$$I_n \stackrel{\mathcal{Q}}{=} [\delta_{ij}]_{i,j=1}^n$$

$$\boxed{[T0.3 : \mathbb{R}^{m \times n}, (\cdot)^\top, (\cdot)^*, \text{tr}]} \quad \begin{aligned} \mathbb{R}^{m \times n} &\stackrel{\mathcal{Q}}{=} \text{set of } m \times n \text{ real matrices} & [\text{deps} : T0.1] \\ \mathbb{C}^{m \times n} &\stackrel{\mathcal{Q}}{=} \text{set of } m \times n \text{ complex matrices} \\ A^\top &\in \mathbb{R}^{n \times m} \quad (A \in \mathbb{R}^{m \times n}) \\ B^* &\stackrel{\mathcal{Q}}{=} \overline{B}^\top \quad (B \in \mathbb{C}^{m \times n}) \\ \text{tr}(M) &\stackrel{\mathcal{Q}}{=} \sum_{i=1}^n M_{ii} \quad (M \in \mathbb{C}^{n \times n}) \end{aligned}$$

$$\boxed{[T0.4 : \det]} \quad \begin{aligned} \det : \mathbb{C}^{n \times n} &\rightarrow \mathbb{C} & [\text{deps} : T0.3] \\ \det(AB) &= \det(A) \det(B) \\ \det(A^\top) &= \det(A) \\ U \text{ invertible} &\Rightarrow \det(UAU^{-1}) = \det(A) \\ (\text{multiplicativity, transpose invariance, similarity invariance}) \end{aligned}$$

$$\boxed{[T0.5 : \oplus, \text{blk}, \det \oplus]} \quad \begin{aligned} \bigoplus_{u=1}^m M_u &\stackrel{\mathcal{Q}}{=} \text{diag}(M_1, \dots, M_m) \\ \det\left(\bigoplus_{u=1}^m M_u\right) &= \prod_{u=1}^m \det(M_u) \end{aligned}$$

(block triangular matrices have determinant equal to product of block determinants)

$$\boxed{[T0.6 : O, SO, U]} \quad \begin{aligned} O(n) &\stackrel{\mathcal{Q}}{=} \{Q \in \mathbb{R}^{n \times n} \mid Q^\top Q = I_n\} & [\text{deps} : T0.2, T0.3, T0.4] \\ SO(n) &\stackrel{\mathcal{Q}}{=} \{Q \in O(n) \mid \det Q = 1\} \\ U(n) &\stackrel{\mathcal{Q}}{=} \{W \in \mathbb{C}^{n \times n} \mid W^* W = I_n\} \\ Q \in O(n) &\Rightarrow Q^{-1} = Q^\top \end{aligned}$$

$$\boxed{[T0.7 : \det O]} \quad \begin{aligned} Q \in O(n) \Rightarrow (\det Q)^2 &= \det(Q^\top Q) = \det(I_n) = 1 & [\text{deps} : T0.6, T0.4] \\ \Rightarrow \det Q &\in \{\pm 1\} \end{aligned}$$

$[T0.8 : \sigma, \sigma_{\min}, \text{normal}]$

$\sigma(M) \stackrel{\otimes}{=} \{\lambda \in \mathbb{C} \mid \exists x \neq 0 : Mx = \lambda x\}$ [deps : T0.3]

$M^*M = MM^* \Rightarrow M \text{ normal}$

Singular values $s_i(M) \stackrel{\otimes}{=} \sqrt{\sigma_i(M^*M)}$

$$\sigma_{\min}(M) \stackrel{\otimes}{=} \min_i s_i(M)$$

$M \text{ normal} \Rightarrow s_i(M) = |\lambda_i| \text{ for } \lambda_i \in \sigma(M)$

(unitary diagonalisation of normal matrices)

$[T0.9 : \cdot , \cdot _1]$

$|x| \stackrel{\otimes}{=} \left(\sum_u x_u^2 \right)^{1/2} \quad (x \in \mathbb{R}^m)$ [deps : T0.8]

$|x|_1 \stackrel{\otimes}{=} \sum_u |x_u|$

$$|M| \stackrel{\otimes}{=} \sup_{v \neq 0} \frac{|Mv|}{|v|} \quad (M \in \mathbb{R}^{m \times n})$$

$$|AB| \leq |A| |B|, \quad |A + B| \leq |A| + |B|$$

$$|a| \stackrel{\otimes}{=} |a| \quad (a \in \mathbb{R})$$

(overloaded notation chooses scalar, vector, or operator norm from context)

$[T0.10 : [A, B]]$

$[A, B] \stackrel{\otimes}{=} AB - BA$ [deps : T0.3]

$[T0.11 : \exp, \log, \exp \text{ props}]$
--

$\exp(M) \stackrel{\otimes}{=} \sum_{t \geq 0} \frac{1}{t!} M^t$ [deps : T0.3, T0.4, T0.5, T0.6]

$\log(X)$ locally defined as $\exp^{-1}(X)$ near I_n

$$[A, B] = 0 \Rightarrow \exp(A + B) = \exp(A) \exp(B)$$

$$(\exp M)^\top = \exp(M^\top)$$

U invertible $\Rightarrow \exp(UMU^{-1}) = U \exp(M) U^{-1}$

$$\det(\exp M) = \exp(\text{tr } M)$$

$$\exp\left(\bigoplus_{u=1}^m M_u\right) = \bigoplus_{u=1}^m \exp(M_u)$$

(all properties follow from the power-series definition)

$[T0.12 : \mathcal{C}, O(\cdot), \approx_\varepsilon]$
--

$\mathcal{C} \stackrel{\otimes}{=} \{C, C', C'', \dots \mid C \geq 0 \text{ finite constants}\}$ [deps : T0.9]

$f = O(g) \Leftrightarrow \exists C \in \mathcal{C} : f \leq Cg$

$$X \approx_\varepsilon Y \Leftrightarrow |X - Y| \leq \varepsilon$$

$[T0.13 : \text{BCH}]$

$\exp(M) \exp(N) = \exp(M + N + \frac{1}{2}[M, N] + R(M, N))$ [deps : T0.11, T0.10, T0.12, T0.9]

$$|R(M, N)| = O(|[M, N]|^2)$$

(Baker–Campbell–Hausdorff series (local version))

$$\boxed{[T0.14 : \ker, \text{im}, \text{rank}]} \quad \begin{aligned} \ker(M) &\stackrel{\otimes}{=} \{x \mid Mx = 0\} & [\text{deps} : T0.3, T0.1] \\ \text{im}(M) &\stackrel{\otimes}{=} \{Mx \mid x\} \\ \text{rank}(M) &\stackrel{\otimes}{=} \dim(\text{im}(M)) \end{aligned}$$

$$\boxed{[T0.15 : \text{proj}]} \quad \begin{aligned} P^2 = P \text{ and } P^\top = P \Rightarrow P \text{ orthogonal projector} && [\text{deps} : T0.14, T0.3] \\ P, Q \text{ orthogonal projectors and } PQ = 0 \Rightarrow \text{im}(P) \perp \text{im}(Q) \\ && (\text{characterisation via images}) \end{aligned}$$

$$\boxed{[T0.16 : \text{skew}_{\text{spec}}]} \quad \begin{aligned} S^\top = -S \Rightarrow S \text{ normal} && [\text{deps} : T0.3, T0.8] \\ \sigma(S) \subset i\mathbb{R} \\ \lambda = i\mu \in \sigma(S) \Rightarrow -\lambda = -i\mu \in \sigma(S) \\ (\text{skew-symmetric spectra occur in conjugate-sign pairs}) \end{aligned}$$

$$\boxed{[T0.17 : J, R_2]} \quad J \stackrel{\otimes}{=} \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, \quad R_2(\theta) \stackrel{\otimes}{=} \exp(\theta J) \quad [\text{deps} : T0.11]$$

$$\boxed{[T0.18 : R_2 \in SO(2)]} \quad \begin{aligned} R_2(\theta)^\top R_2(\theta) &= I_2 & [\text{deps} : T0.17, T0.16, T0.11] \\ \det R_2(\theta) &= 1 \\ \Rightarrow R_2(\theta) &\in SO(2) \end{aligned}$$

(uses $J^\top = -J$ and properties of \exp and the determinant)

$$\boxed{[T0.19 : \text{Weyl}]} \quad |\sigma_{\min}(A + E) - \sigma_{\min}(A)| \leq |E| \quad [\text{deps} : T0.8, T0.9]$$

$$\boxed{[T0.20 : (X^{-1})^\top]} \quad X \text{ invertible} \Rightarrow (X^{-1})^\top = (X^\top)^{-1} \quad [\text{deps} : T0.3]$$

2 Tier 1: Structured Generators

$$\boxed{[T1.1 : L_k]} \quad \forall k \in K : L_k \in \mathbb{R}^{d_h \times d_h}, \quad L_k^\top = -L_k, \quad [L_a, L_b] = 0 \quad [\text{deps} : T0.1, T0.3, T0.10, T0.16]$$

(model hypothesis: commuting skew-symmetric generators)

$$\boxed{[T1.2 : \text{normal } L_k]} \quad L_k^\top = -L_k \Rightarrow L_k \text{ normal} \quad [\text{deps} : T1.1, T0.16]$$

$$\boxed{[T1.3 : A(r)]} \quad r \in \mathbb{R}^{d_c} \Rightarrow A(r) \stackrel{\otimes}{=} \sum_{k=1}^{d_c} r_k L_k \quad [\text{deps} : T1.1, T0.1]$$

$$\boxed{[T1.4 : A(r)^\top]} \quad A(r)^\top = -A(r) \quad [\text{deps} : T1.3, T1.1, T0.3]$$

$$\boxed{[T1.5 : [A(r), A(s)]]} \quad [A(r), A(s)] = 0 \quad [\text{deps} : T1.3, T1.1, T0.10]$$

$$\boxed{[T1.6 : R_{\text{STR}}]} \quad R_{\text{STR}}(r) \stackrel{\otimes}{=} \exp(A(r)) \quad [\text{deps} : T1.3, T0.11]$$

3 Tier 2: Joint Block Structure

$$\boxed{[T2.1 : \text{sim-unitary}]} \quad \{L_k\}_{k \in K} \text{ commute and are normal} \quad [\text{deps : } T1.2, T1.1, T0.6, T0.16]$$

$$\Rightarrow \exists \mathcal{U}_c \in U(d_h) : \mathcal{U}_c^* L_k \mathcal{U}_c = \text{diag}(i\lambda_{k,1}, \dots, i\lambda_{k,d_h})$$

$$\lambda_{k,j} \in \mathbb{R}$$

$$\boxed{[T2.2 : \text{real 2D planes}]} \quad \begin{aligned} & \text{Fix } j \in H : \mathcal{U}_c^* L_k \mathcal{U}_c e_j = i\lambda_{k,j} e_j & [\text{deps : } T2.1, T0.17, T0.16, T0.3] \\ & \Rightarrow L_k(\Re e_j) = \lambda_{k,j} J(\Re e_j, \Im e_j) \\ & L_k(\Im e_j) = \lambda_{k,j} J(\Re e_j, \Im e_j) \\ & \Rightarrow \text{span}\{\Re e_j, \Im e_j\} \text{ is a real, 2-D, } L_k\text{-invariant plane} \end{aligned}$$

$$\boxed{[T2.3 : \text{JointBlockDiag}]} \quad \begin{aligned} & \exists m, d_{\text{null}} \in \mathbb{N} : 2m + d_{\text{null}} = d_h & [\text{deps : } T2.2, T0.6, T0.16, T0.5] \\ & \exists U \in O(d_h) : U^\top L_k U = \left(\bigoplus_{u=1}^m \lambda_{k,u} J \right) \oplus 0_{d_{\text{null}} \times d_{\text{null}}} \end{aligned}$$

$$\boxed{[T2.4 : \beta, \theta]} \quad \begin{aligned} \beta_u &\stackrel{\otimes}{=} (\lambda_{1,u}, \dots, \lambda_{d_c,u})^\top \in \mathbb{R}^{d_c} & [\text{deps : } T2.3, T1.1, T0.1] \\ \theta_u(r) &\stackrel{\otimes}{=} \beta_u^\top r \in \mathbb{R} \end{aligned}$$

$$\boxed{[T2.5 : U^\top A(r)U]} \quad U^\top A(r)U = \left(\bigoplus_{u=1}^m \theta_u(r) J \right) \oplus 0_{d_{\text{null}} \times d_{\text{null}}} \quad [\text{deps : } T1.3, T2.3, T2.4, T0.17, T0.5]$$

$$\boxed{[T2.6 : \text{exp blk}]} \quad \exp(U^\top A(r)U) = \left(\bigoplus_{u=1}^m R_2(\theta_u(r)) \right) \oplus I_{d_{\text{null}}} \quad [\text{deps : } T2.5, T0.17, T0.11, T0.5]$$

$$\boxed{[T2.7 : R_{\text{STR}} \text{ form}]} \quad R_{\text{STR}}(r) = U \left[\left(\bigoplus_{u=1}^m R_2(\theta_u(r)) \right) \oplus I_{d_{\text{null}}} \right] U^\top \quad [\text{deps : } T1.6, T2.6, T0.11, T2.3]$$

$$\boxed{[T2.8 : R_{\text{STR}} \in SO]} \quad R_{\text{STR}}(r) \in SO(d_h) \quad [\text{deps : } T2.7, T0.11]$$

(uses $R_2(\theta) \in SO(2)$, block determinant multiplicativity, and $U \in O(d_h)$)

$$\boxed{[T2.9 : \text{rel } R_{\text{STR}}]} \quad R_{\text{STR}}(r_i)^\top R_{\text{STR}}(r_j) = R_{\text{STR}}(r_j - r_i) \quad [\text{deps : } T1.4, T1.5, T1.6, T0.11]$$

$$\boxed{[T2.10 : \Pi_{\text{act}}, \Pi_{\text{null}}, d_{\text{act}}]} \quad \begin{aligned} \Pi_{\text{act}} &\stackrel{\otimes}{=} U \text{diag}(I_{2m}, 0) U^\top & [\text{deps : } T2.3, T0.15, T0.2, T0.14, T0.6] \\ \Pi_{\text{null}} &\stackrel{\otimes}{=} U \text{diag}(0, I_{d_{\text{null}}}) U^\top \\ \Pi_{\text{act}}^\top &= \Pi_{\text{act}}, \quad \Pi_{\text{act}}^2 = \Pi_{\text{act}} \\ \Pi_{\text{null}}^\top &= \Pi_{\text{null}}, \quad \Pi_{\text{null}}^2 = \Pi_{\text{null}} \\ \Pi_{\text{act}} \Pi_{\text{null}} &= 0 \Rightarrow \text{im}(\Pi_{\text{act}}) \perp \text{im}(\Pi_{\text{null}}) \\ d_{\text{act}} &\stackrel{\otimes}{=} \text{rank}(\Pi_{\text{act}}) = 2m \end{aligned}$$

$$\boxed{[T2.11 : \Pi_{\text{act}} \text{ comm}]} \quad R_{\text{STR}}(r) \Pi_{\text{act}} = \Pi_{\text{act}} R_{\text{STR}}(r) \quad [\text{deps : } T2.7, T2.10, T0.5, T0.6]$$

4 Tier 3: Cayley Post-Rotations

$$\boxed{[T3.1 : \text{Cayley}]} \quad S \in \mathbb{R}^{d_h \times d_h}, \quad S^\top = -S \Rightarrow (I + S) \text{ invertible} \quad [\text{deps : } T0.16, T0.3, T0.2]$$

$$P_{\text{sp}} \stackrel{\mathcal{D}}{=} (I - S)(I + S)^{-1}$$

$$(I - S)(I + S) = (I + S)(I - S) = I - S^2$$

$$(\text{uses } \sigma(S) \subset i\mathbb{R} \text{ so } 1 + i\mu \neq 0)$$

$$\boxed{[T3.2 : P_{\text{sp}} \in SO, \mathcal{B}]} \quad P_{\text{sp}}^\top = ((I + S)^{-1})^\top (I - S)^\top = I \quad [\text{deps : } T3.1, T0.16, T0.20, T0.7, T0.6, T0.4]$$

$$\det(P_{\text{sp}}) = \prod_u \frac{1 - i\mu_u}{1 + i\mu_u} = 1$$

$$\mathcal{B} \stackrel{\mathcal{D}}{=} \{(I - S)(I + S)^{-1} \mid S^\top = -S\} \subseteq SO(d_h)$$

$$\boxed{[T3.3 : \text{Cayley surj}]} \quad Q \in SO(d_h), \quad -1 \notin \sigma(Q) \Rightarrow S_Q \stackrel{\mathcal{D}}{=} (I - Q)(I + Q)^{-1} \quad [\text{deps : } T0.16, T0.20, T0.4]$$

$$S_Q^\top = -S_Q, \quad (I - S_Q)(I + S_Q)^{-1} = Q$$

(Cayley transform bijection on $SO(d_h)$ with no -1 eigenvalues)

$$\boxed{[T3.4 : \mathcal{A}]} \quad \mathcal{A} \stackrel{\mathcal{D}}{=} \left\{ U \text{diag}(I_{2m}, R_{\text{null}}) U^\top \mid R_{\text{null}} \in SO(d_{\text{null}}) \right\} \subseteq SO(d_h) \quad [\text{deps : } T2.3, T0.6, T0.7, T0.5, T0.4, T0.2]$$

$$R_{\text{null}} \in SO(d_{\text{null}}), \quad -1 \notin \sigma(R_{\text{null}}) \Rightarrow U \text{diag}(I_{2m}, R_{\text{null}}) U^\top \in \mathcal{B}$$

$$\boxed{[T3.5 : R_{\text{sp}}]} \quad R_{\text{sp}}(r) \stackrel{\mathcal{D}}{=} R_{\text{STR}}(r)P_{\text{sp}} \quad [\text{deps : } T2.8, T3.2, T0.6, T0.4]$$

$$R_{\text{STR}}(r) \in SO(d_h), \quad P_{\text{sp}} \in SO(d_h) \Rightarrow R_{\text{sp}}(r) \in SO(d_h)$$

$$\boxed{[T3.6 : \Pi_{\text{act}} R_{\text{sp}} \Pi_{\text{act}}]} \quad P_{\text{sp}} \in \mathcal{A} \Rightarrow \Pi_{\text{act}} P_{\text{sp}} = \Pi_{\text{act}} = P_{\text{sp}} \Pi_{\text{act}} \quad [\text{deps : } T3.5, T2.10, T2.11, T3.4, T0.5, T0.6]$$

$$R_{\text{STR}}(r) \Pi_{\text{act}} = \Pi_{\text{act}} R_{\text{STR}}(r)$$

$$\Rightarrow \Pi_{\text{act}} R_{\text{sp}}(r) \Pi_{\text{act}} = \Pi_{\text{act}} R_{\text{STR}}(r) \Pi_{\text{act}}$$

$$\boxed{[T3.7 : \Pi\text{-rel } R_{\text{sp}}]} \quad P_{\text{sp}} \in \mathcal{A} \Rightarrow \Pi_{\text{act}} R_{\text{sp}}(r_i)^\top R_{\text{sp}}(r_j) \Pi_{\text{act}} = \Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}} \quad [\text{deps : } T3.6, T2.9, T3.5, T3.2]$$

5 Tier 4: Attention Queries, Keys, and Scores

$$\boxed{[T4.1 : W, q, k, v]} \quad W_Q, W_K, W_V \in \mathbb{R}^{d_h \times D} \quad [\text{deps : } T0.1, T0.3]$$

$$x_i \in \mathbb{R}^D, \quad r_i \in \mathbb{R}^{d_c}$$

$$q_i = W_Q x_i, \quad k_j = W_K x_j, \quad v_j = W_V x_j$$

$$\boxed{[T4.2 : q^{(\text{act})}, k^{(\text{act})}, d_{\text{act}}]} \quad q_i^{(\text{act})} \stackrel{\mathcal{D}}{=} \Pi_{\text{act}} q_i, \quad k_j^{(\text{act})} \stackrel{\mathcal{D}}{=} \Pi_{\text{act}} k_j \quad [\text{deps : } T2.10, T4.1]$$

$$d_{\text{act}} \stackrel{\mathcal{D}}{=} 2m$$

$$\boxed{[T4.3 : \tilde{q}, \tilde{k}]} \quad \tilde{q}_i = \Pi_{\text{act}} R_{\text{sp}}(r_i) q_i^{(\text{act})} \quad [\text{deps : } T3.5, T4.2, T2.10] \\ \tilde{k}_j = \Pi_{\text{act}} R_{\text{sp}}(r_j) k_j^{(\text{act})}$$

$$\boxed{[T4.4 : \alpha_{ij} \text{ def}]} \quad \alpha_{ij} \stackrel{\otimes}{=} \frac{\tilde{q}_i^\top \tilde{k}_j}{\sqrt{d_{\text{act}}}} \quad [\text{deps : } T4.3, T4.2, T2.10, T0.3] \\ = \frac{1}{\sqrt{d_{\text{act}}}} (q_i^{(\text{act})})^\top (\Pi_{\text{act}} R_{\text{sp}}(r_i) \Pi_{\text{act}})^\top (\Pi_{\text{act}} R_{\text{sp}}(r_j) \Pi_{\text{act}}) k_j^{(\text{act})} \\ (\text{uses } \Pi_{\text{act}}^\top = \Pi_{\text{act}} = \Pi_{\text{act}}^2)$$

$$\boxed{[T4.5 : \alpha_{ij} \text{ rel}]} \quad P_{\text{sp}} \in \mathcal{A} \Rightarrow \alpha_{ij} = \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top (\Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}}) k_j^{(\text{act})} \quad [\text{deps : } T4.4, T3.6, T3.7, T4.2, T2.10]$$

6 Tier 5: Stability Estimates

$$\boxed{[T5.1 : \varepsilon \text{ comm}]} \quad \varepsilon_{ab} \stackrel{\otimes}{=} |[L_a, L_b]|, \quad \varepsilon \stackrel{\otimes}{=} \max_{a,b} \varepsilon_{ab} \quad [\text{deps : } T1.3, T1.1, T0.10, T0.9] \\ |[A(r), A(s)]| \leq \sum_{a,b} |r_a| |s_b| \varepsilon_{ab} \leq |r|_1 |s|_1 \varepsilon$$

$$\boxed{[T5.2 : \text{BCH err}]} \quad |\log(\exp(A(r)) \exp(A(s))) - (A(r) + A(s))| \leq \frac{1}{2} |[A(r), A(s)]| + C |[A(r), A(s)]|^2 \quad [\text{deps : } T0.1]$$

$$\boxed{[T5.3 : R_{\text{STR}} \text{ rel approx}]} \quad |R_{\text{STR}}(r)^\top R_{\text{STR}}(s) - R_{\text{STR}}(s - r)| \leq C \varepsilon |r|_1 |s|_1 + O(\varepsilon^2) \quad [\text{deps : } T1.6, T1.4, T1.5, T5.3]$$

$$\boxed{[T5.4 : S = S_- + E]} \quad S = S_- + E, \quad S_-^\top = -S_-, \quad |E| \leq \eta, \quad 0 \leq \eta < 1 \quad [\text{deps : } T0.16, T0.9]$$

$$\boxed{[T5.5 : \sigma_{\min}(I + S_-)]} \quad S_-^\top = -S_- \Rightarrow \sigma(S_-) \subset i\mathbb{R} \quad [\text{deps : } T5.4, T0.16, T0.8] \\ \Rightarrow |1 + i\mu| \geq 1 \text{ for } \lambda = i\mu \in \sigma(S_-) \\ \sigma_{\min}(I + S_-) \geq 1$$

$$\boxed{[T5.6 : \sigma_{\min}(I + S)]} \quad \sigma_{\min}(I + S) \geq \sigma_{\min}(I + S_-) - |E| \geq 1 - \eta > 0 \quad [\text{deps : } T5.4, T5.5, T0.19, T0.9] \\ \Rightarrow I + S \text{ invertible}$$

$$\boxed{[T5.7 : \text{Cayley Lip}]} \quad P_{\text{sp}}(S) \stackrel{\otimes}{=} (I - S)(I + S)^{-1}, \quad P_{\text{sp}}(S_-) \stackrel{\otimes}{=} (I - S_-)(I + S_-)^{-1} \quad [\text{deps : } T5.6, T5.4, T0.20, T0.9] \\ |P_{\text{sp}}(S) - P_{\text{sp}}(S_-)| \leq C' \eta$$

$$\boxed{[T5.8 : \eta_{\text{mix}}]} P_{\text{sp}} = U \begin{bmatrix} A & B \\ C & D \end{bmatrix} U^\top, \quad \eta_{\text{mix}} \stackrel{\otimes}{=} |B| + |C| \quad [\text{deps : } T3.2, T2.3, T2.10, T0.9, T0.6, T0.5, T0.7]$$

$$P_{\text{sp}} \in SO(d_h) \Rightarrow \begin{bmatrix} A & B \\ C & D \end{bmatrix}^\top \begin{bmatrix} A & B \\ C & D \end{bmatrix} = I$$

$$\Rightarrow |A - I_{2m}| \leq C'' \eta_{\text{mix}}$$

$$\Rightarrow |\Pi_{\text{act}} P_{\text{sp}} \Pi_{\text{act}} - \Pi_{\text{act}}| \leq C'' \eta_{\text{mix}}$$

$$\boxed{[T5.9 : \Pi R_{\text{sp}} \Pi \text{ approx}]} |\Pi_{\text{act}} R_{\text{sp}}(r) \Pi_{\text{act}} - \Pi_{\text{act}} R_{\text{STR}}(r) \Pi_{\text{act}}| \leq C'' \eta_{\text{mix}} \quad [\text{deps : } T3.5, T2.11, T5.8, T2.10, T0.9, T0.7]$$

$$\boxed{[T5.10 : \alpha_{ij} \text{ stab}]} \quad \Delta_{ij} \stackrel{\otimes}{=} C(\varepsilon |r_i|_1 |r_j|_1 + \eta_{\text{mix}}) |q_i^{(\text{act})}| |k_j^{(\text{act})}| \quad [\text{deps : } T4.4, T5.9, T5.3, T4.2, T5.10]$$

$$\left| \alpha_{ij} - \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top (\Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}}) k_j^{(\text{act})} \right| \leq \Delta_{ij}$$

7 Tier 6: Conclusions

$$\boxed{[T6.1 : \text{GOAL}]} (T1.1 \wedge T2.3 \wedge P_{\text{sp}} \in \mathcal{A}) \Rightarrow \alpha_{ij} = \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top (\Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}}) k_j^{(\text{act})} \quad [\text{deps : } T4.5, T5.10]$$

$$\boxed{[T6.2 : \text{GOAL}_{\text{rob}}]} (|[L_a, L_b]| \leq \varepsilon \ \forall a, b, \quad P_{\text{sp}} \in SO(d_h) \text{ with mixing } \eta_{\text{mix}}) \Rightarrow \alpha_{ij} \approx_{\Delta_{ij}} \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top (\Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}}) k_j^{(\text{act})}$$

$$\Delta_{ij} \leq C(\varepsilon |r_i|_1 |r_j|_1 + \eta_{\text{mix}})$$

8 Tier 6': Quantified Robustness

$$\boxed{[T6'.0 : \varepsilon_{\text{comm}}, \Lambda, R, S, S_-, E, \rho, \eta, \eta_{\text{mix}}]} \quad \varepsilon_{\text{comm}} \stackrel{\otimes}{=} \max_{a, b \in K} |[L_a, L_b]| \quad [\text{deps : } T1.1, T2.3, T2.4, T2.5]$$

$$\Lambda \stackrel{\otimes}{=} \max_{u \in \{1, \dots, m\}} \sum_{k=1}^{d_c} |\lambda_{k,u}|$$

$$R \stackrel{\otimes}{=} \max_i |r_i|_1$$

$$S \stackrel{\otimes}{=} S_- + E, \quad S_-^\top = -S_-, \quad E^\top = -E$$

$$\rho \stackrel{\otimes}{=} |S|, \quad \eta \stackrel{\otimes}{=} |E|$$

$$\eta_{\text{mix}} \stackrel{\otimes}{=} \max \left\{ |\Pi_{\text{null}} P_{\text{sp}} \Pi_{\text{act}}|, |\Pi_{\text{act}} P_{\text{sp}} \Pi_{\text{null}}| \right\}$$

$$\boxed{[T6'.1 : |A(r)|, \text{BCH const}]} |A(r)| \stackrel{T2.5}{\leq} \sum_{u=1}^m |\theta_u(r)| \stackrel{T2.4}{=} \sum_{u=1}^m \left| \sum_{k=1}^{d_c} r_k \lambda_{k,u} \right| \leq \Lambda |r|_1 \quad [\text{deps : } T2.5, T2.4, T0.9, T0.12]$$

$$C_{\text{BCH}}(r, s) \stackrel{\otimes}{=} C_\star \exp(2\Lambda(|r|_1 + |s|_1))$$

$$\boxed{[T6'.2 : \text{BCH quant}]} \left| \log(e^{A(r)}e^{A(s)}) - (A(r) + A(s)) \right| \stackrel{T0.13}{\leq} \frac{1}{2}|[A(r), A(s)]| + C|[A(r), A(s)]|^2 \quad [\text{deps : } T0.13, T5.1]$$

$$\stackrel{T5.1}{\leq} \frac{1}{2}\varepsilon_{\text{comm}}|r|_1|s|_1 + C_{\text{BCH}}(r, s)\varepsilon_{\text{comm}}^2|r|_1^2|s|_1^2$$

$$\boxed{[T6'.3 : \text{Rel approx } R_{\text{STR}}]} \quad R_{\text{STR}}(r)^\top R_{\text{STR}}(s) \stackrel{T1.6}{=} e^{-A(r)}e^{A(s)} \quad [\text{deps : } T0.13, T5.1]$$

$$\log(R_{\text{STR}}(r)^\top R_{\text{STR}}(s)) \stackrel{T0.11}{=} \log(e^{-A(r)}e^{A(s)}) \stackrel{T6'.2}{=} A(s) - A(r) + \Delta_{rs}$$

$$A(s) - A(r) \stackrel{T1.3}{=} A(s - r)$$

$$\Rightarrow |R_{\text{STR}}(r)^\top R_{\text{STR}}(s) - R_{\text{STR}}(s - r)| \leq C_1\varepsilon_{\text{comm}}|r|_1|s|_1 + C_2\varepsilon_{\text{comm}}^2|r|_1^2|s|_1^2$$

$$\boxed{[T6'.4 : \sigma_{\min}(I + S), |(I + S)^{-1}|]} \quad \sigma_{\min}(I + S) \stackrel{T0.19}{\geq} \sigma_{\min}(I) - |S| = 1 - \rho \quad [\text{deps : } T0.19, T0.8, T0.9]$$

$$|(I + S)^{-1}| = \sigma_{\min}(I + S)^{-1} \leq (1 - \rho)^{-1}$$

$$\boxed{[T6'.5 : \text{Cayley Lip}]} \quad P_{\text{sp}}(S) \stackrel{\emptyset}{=} (I - S)(I + S)^{-1} \quad [\text{deps : } T3.2]$$

$$P_{\text{sp}}(S_1) - P_{\text{sp}}(S_2) = (S_2 - S_1)(I + S_1)^{-1} + (I - S_2)[(I + S_2)^{-1} - (I + S_1)^{-1}]$$

$$(I + S_2)^{-1} - (I + S_1)^{-1} = (I + S_2)^{-1}(S_1 - S_2)(I + S_1)^{-1}$$

$$\Rightarrow |P_{\text{sp}}(S_1) - P_{\text{sp}}(S_2)| \leq \frac{2}{\sigma_{\min}(I + S_1)\sigma_{\min}(I + S_2)}|S_1 - S_2|$$

$$\Rightarrow |P_{\text{sp}}(S) - P_{\text{sp}}(S_-)| \leq \frac{2}{(1 - \rho)^2}\eta$$

$$\boxed{[T6'.6 : \text{Neumann}]} \quad |S| = \rho < 1 \Rightarrow (I + S)^{-1} = \sum_{t=0}^{\infty}(-S)^t \quad [\text{deps : } T0.9, T0.12]$$

$$P_{\text{sp}}(S) = (I - S) \sum_{t=0}^{\infty}(-S)^t = I - 2 \sum_{u=1}^{\infty}(-1)^{u-1}S^u$$

$$\boxed{[T6'.7 : \text{QuadMix } \Pi_{\text{act}}]} \quad \Pi_{\text{act}}S\Pi_{\text{act}} = 0, \quad \Pi_{\text{null}}S\Pi_{\text{null}} = S_{\text{null}} \quad [\text{deps : } T0.9, T0.12]$$

$$\Pi_{\text{act}}S^2\Pi_{\text{act}} = (\Pi_{\text{act}}S\Pi_{\text{null}})(\Pi_{\text{null}}S\Pi_{\text{act}}) \Rightarrow |\Pi_{\text{act}}S^2\Pi_{\text{act}}| \leq \eta^2$$

$$\Pi_{\text{act}}P_{\text{sp}}(S)\Pi_{\text{act}} - \Pi_{\text{act}} = -2 \sum_{u=1}^{\infty}(-1)^{u-1}\Pi_{\text{act}}S^u\Pi_{\text{act}} = -2 \sum_{u=2}^{\infty}(-1)^{u-1}\Pi_{\text{act}}S^u\Pi_{\text{act}}$$

$$\Rightarrow |\Pi_{\text{act}}P_{\text{sp}}(S)\Pi_{\text{act}} - \Pi_{\text{act}}| \leq 2 \sum_{u=2}^{\infty}|\Pi_{\text{act}}S^u\Pi_{\text{act}}| \leq 2 \sum_{u=2}^{\infty}\rho^{u-2}|\Pi_{\text{act}}S^2\Pi_{\text{act}}|$$

$$\leq \frac{2}{1 - \rho}\eta^2 \leq \frac{2}{(1 - \rho)^3}\eta^2$$

$$\boxed{[T6'.8 : \text{Mix model-free}]} \quad P_{\text{sp}} = U \begin{bmatrix} A & B \\ C & D \end{bmatrix} U^\top, \quad A^\top A + C^\top C = I, \quad AA^\top + BB^\top = I \quad [\text{deps : } T3.2, T2.3, T2.4]$$

$$|A - I| \leq 2 \max\{|B|, |C|\} = 2\eta_{\text{mix}}$$

$$\Rightarrow |\Pi_{\text{act}}P_{\text{sp}}\Pi_{\text{act}} - \Pi_{\text{act}}| \leq 2\eta_{\text{mix}}$$

$$\begin{aligned}
& [T6'.9 : \text{Rob } R_{\text{sp}} \text{ on act}] \quad |\Pi_{\text{act}} R_{\text{sp}}(r) \Pi_{\text{act}} - \Pi_{\text{act}} R_{\text{STR}}(r) \Pi_{\text{act}}| = |\Pi_{\text{act}} R_{\text{STR}}(r) (P_{\text{sp}} - \Pi_{\text{act}} - \Pi_{\text{null}}) \Pi_{\text{act}}| \quad [\text{deps : } T6'.1, T6'.2] \\
& \stackrel{T0.9}{\leq} |P_{\text{sp}} - \Pi_{\text{act}} - \Pi_{\text{null}}| \leq \begin{cases} \frac{2}{(1-\rho)^3} \eta^2, & (\text{ESPR/Cayley/pure mix}), \\ 2\eta_{\text{mix}}, & (\text{model-free}). \end{cases} =: \delta_{\text{mix}}
\end{aligned}$$

$$\begin{aligned}
& [T6'.10 : \text{Logit err}] \quad \alpha_{ij} \stackrel{T4.4}{=} \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top \left(\Pi_{\text{act}} R_{\text{sp}}(r_i)^\top R_{\text{sp}}(r_j) \Pi_{\text{act}} \right) k_j^{(\text{act})} \quad [\text{deps : } T6'.1, T6'.2] \\
& \alpha_{ij}^* \stackrel{\mathcal{D}}{=} \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top \left(\Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}} \right) k_j^{(\text{act})} \\
& |\alpha_{ij} - \alpha_{ij}^*| \leq \frac{1}{\sqrt{2m}} |q_i^{(\text{act})}| |k_j^{(\text{act})}| \\
& \cdot \left(|\Pi_{\text{act}} (R_{\text{sp}}(r_i)^\top R_{\text{sp}}(r_j) - R_{\text{STR}}(r_j - r_i)) \Pi_{\text{act}}| + |\Pi_{\text{act}} R_{\text{sp}}(\cdot) \Pi_{\text{act}} - \Pi_{\text{act}} R_{\text{STR}}(\cdot) \Pi_{\text{act}}| \right) \\
& \leq |q_i^{(\text{act})}| |k_j^{(\text{act})}| \left(C_1 \varepsilon_{\text{comm}} |r_i|_1 |r_j|_1 + C_2 \varepsilon_{\text{comm}}^2 |r_i|_1^2 |r_j|_1^2 + \delta_{\text{mix}} \right)
\end{aligned}$$

$$\begin{aligned}
& [T6'.11 : \text{Rob theorem}] \quad |\alpha_{ij} - \alpha_{ij}^*| \leq \left(\frac{1}{2} + C_{\text{BCH}}(r_i, r_j) \right) \varepsilon_{\text{comm}} |r_i|_1 |r_j|_1 |q_i^{(\text{act})}| |k_j^{(\text{act})}| \quad [\text{deps : } T6'.1, T6'.3, T6'.9] \\
& \quad + C_{\text{mix}} \delta_{\text{mix}} |q_i^{(\text{act})}| |k_j^{(\text{act})}| \\
& \begin{cases} C_{\text{mix}} = 1, \quad \delta_{\text{mix}} = \eta_{\text{mix}}, & (\text{model-free}), \\ C_{\text{mix}} = \frac{2}{(1-\rho)^3}, \quad \delta_{\text{mix}} = \eta^2, & (\text{ESPR/Cayley/pure mix}). \end{cases}
\end{aligned}$$

$$\begin{aligned}
& [T6'.12 : \text{Replacements}] \quad \{T5.2, T5.3\} \rightsquigarrow \{T6'.2, T6'.3\} \quad [\text{deps : } T5.2, T5.3, T5.7, T5.8, T5.9, T6.2, T6'.2, T6'.3] \\
& \quad \{T5.7\} \rightsquigarrow \{T6'.5\} \\
& \quad \{T5.8, T5.9\} \rightsquigarrow \{T6'.7, T6'.9\} \\
& \quad \{T6.2\} \rightsquigarrow \{T6'.11\}
\end{aligned}$$

9 Tier 7: Positive Random Features and Variance Bounds

$$\begin{aligned}
& [T7.0 : \mathcal{N}(0, I_n), \mathbb{E}, \text{Var}] \quad \mathcal{N}(0, I_n) \stackrel{\mathcal{D}}{=} \text{the mean-zero Gaussian law on } \mathbb{R}^n \text{ with density } (2\pi)^{-n/2} \exp(-\frac{1}{2}|w|^2) \\
& w \sim \mathcal{N}(0, I_n) \Rightarrow \mathbb{E}[w] = 0, \quad \mathbb{E}[ww^\top] = I_n \\
& \mathbb{E}[\cdot] \stackrel{\mathcal{D}}{=} \text{expectation with respect to the (joint) law of all Gaussian draws } w \\
& \text{Var}[Z] \stackrel{\mathcal{D}}{=} \mathbb{E}[Z^2] - (\mathbb{E}[Z])^2 \quad \text{for any scalar random variable } Z
\end{aligned}$$

[T7.1 : RF setup]

$$d_{\text{act}} = 2m \text{ from T2.10}$$

For probabilistic estimates we draw $w \sim \mathcal{N}(0, I_{d_{\text{act}}})$ in the sense of T7.0

$$\phi_w(x) \stackrel{\mathcal{D}}{=} \exp\left(w^\top x - \frac{1}{2}|x|^2\right) \quad (x \in \mathbb{R}^{d_{\text{act}}})$$

(positive softmax random feature map as in FAVOR + estimators)

$$Z(x, y; w) \stackrel{\mathcal{D}}{=} \phi_w(x) \phi_w(y) = \exp\left(w^\top (x + y) - \frac{1}{2}(|x|^2 + |y|^2)\right) \quad (x, y \in \mathbb{R}^{d_{\text{act}}})$$

$$\widehat{\text{SM}}_{n_{\text{rf}}}(x, y) \stackrel{\mathcal{D}}{=} \frac{1}{n_{\text{rf}}} \sum_{t=1}^{n_{\text{rf}}} Z(x, y; w_t) \quad \text{for } n_{\text{rf}} \in \mathbb{N}, w_t \text{ i.i.d. } \mathcal{N}(0, I_{d_{\text{act}}})$$

Note: $\widehat{\text{SM}}_{n_{\text{rf}}}(x, y)$ is our Monte Carlo estimator of $\exp(x^\top y)$ (unnormalized softmax kernel)

[T7.2 : $\widehat{q}_i, \widehat{k}_j, \alpha_{ij}$]

\tilde{q}_i, \tilde{k}_j from T4.3,

$$\alpha_{ij} \stackrel{\mathcal{D}}{=} \frac{1}{\sqrt{d_{\text{act}}}} \tilde{q}_i^\top \tilde{k}_j$$

$$\widehat{q}_i \stackrel{\mathcal{D}}{=} d_{\text{act}}^{-1/4}$$

$$\widehat{q}_i^\top \widehat{k}_j = (d_{\text{act}}^{-1/4} \tilde{q}_i)^\top (d_{\text{act}}^{-1/4} \tilde{k}_j)$$

(The attention logit α_{ij} is *exactly* a Euclidean dot product between $\widehat{q}_i, \widehat{k}_j \in \mathbb{R}^{d_{\text{act}}}$, with no renormalization)

[T7.3 : $\mathbb{E}Z = \exp(x^\top y)$, unbiasedness]

Note: $\widehat{\text{SM}}_{n_{\text{rf}}}(\widehat{q}_i, \widehat{k}_j)$ is an *unbiased* positive estimator of $\exp(\alpha_{ij})$, the *unnormalized* softmax kernel.

[T7.4 : $\mathbb{E}Z^2$, $\text{Var}Z$, closed form]

Continue with fixed $x, y \in \mathbb{R}^n$

$$Z(x, y; w) = \exp\left(w^\top(x + y) - \frac{1}{2}(|x|^2 + |y|^2)\right)$$

$$Z(x, y; w)^2 = \exp\left(2w^\top(x + y) - (|x|^2 + |y|^2)\right)$$

$$\mathbb{E}[Z(x, y; w)^2] = \exp\left(\frac{1}{2}|2(x + y)|^2 - (|x|^2 + |y|^2)\right) = \exp\left(2|x + y|^2 - (|x|^2 + |y|^2)\right)$$

$$|x + y|^2 = |x|^2 + |y|^2 + 2x^\top y \text{ (as used in T7.4)}$$

$$\Rightarrow 2|x + y|^2 - (|x|^2 + |y|^2) = 2(|x|^2 + |y|^2 + 2x^\top y) - (|x|^2 + |y|^2) = |x|^2 + |y|^2 + 4x^\top y$$

$$\Rightarrow \mathbb{E}[Z(x, y; w)^2] = \exp\left(|x|^2 + |y|^2 + 4x^\top y\right)$$

$$\mathbb{E}[Z(x, y; w)] = \exp(x^\top y) \quad \text{from T7.4}$$

$$(\mathbb{E}[Z(x, y; w)])^2 = \exp(2x^\top y)$$

$$\text{Var}[Z(x, y; w)] = \mathbb{E}[Z(x, y; w)^2] - (\mathbb{E}[Z(x, y; w)])^2$$

$$= \exp\left(|x|^2 + |y|^2 + 4x^\top y\right) - \exp\left(2x^\top y\right)$$

$$\text{Var}[Z(x, y; w)] = \exp\left(2x^\top y\right) \left[\exp\left(|x|^2 + |y|^2 + 2x^\top y\right) - \exp\left(2x^\top y\right) \right]$$

Note: The single-sample variance is finite and given in closed form for all finite n .

[T7.5 : Var $\widehat{\text{SM}}_{n_{\text{rf}}}$, geom link]

The draws w_t in T7.1 are i.i.d., so the $Z(\cdot; w_t)$ are

$$\widehat{\text{SM}}_{n_{\text{rf}}}(x, y) = \frac{1}{n_{\text{rf}}} \sum_{t=1}^{n_{\text{rf}}} Z(x, y; w_t)$$

$$\text{Var}\left[\widehat{\text{SM}}_{n_{\text{rf}}}(x, y)\right] = \frac{1}{n_{\text{rf}}^2} \sum_{t=1}^{n_{\text{rf}}} \text{Var}[Z(x, y; w_t)] = \frac{1}{n_{\text{rf}}} \text{Var}[Z(x, y; w)] \quad (\text{i.i.d. variance})$$

Insert the expression from

$$\text{Var}\left[\widehat{\text{SM}}_{n_{\text{rf}}}(x, y)\right] = \frac{1}{n_{\text{rf}}} \exp(2x^\top y) \left[\exp(|x|^2 + |y|^2 + 2x^\top y) - 1 \right]$$

Now specialise $x = \widehat{q}_i$, $y = \widehat{k}_j$ from T7.2. Then $x^\top y = \alpha_{ij}$

$$\Rightarrow \text{Var}\left[\widehat{\text{SM}}_{n_{\text{rf}}}(\widehat{q}_i, \widehat{k}_j)\right] = \frac{1}{n_{\text{rf}}} \exp(2\alpha_{ij}) \left[\exp(|\widehat{q}_i|^2 + |\widehat{k}_j|^2 + 2\alpha_{ij}) - 1 \right]$$

Relate $|\widehat{q}_i|, |\widehat{k}_j|$ back to geometric quantities already controlled in T7.2

$$\widehat{q}_i = d_{\text{act}}^{-1/4} \tilde{q}_i, \quad \widehat{k}_j = d_{\text{act}}^{-1/4} \tilde{k}_j \text{ from T7.2} \Rightarrow |\widehat{q}_i|^2 + |\widehat{k}_j|^2 = d_{\text{act}}^{-1/2} (|\tilde{q}_i|^2 + |\tilde{k}_j|^2)$$

$$\tilde{q}_i = \Pi_{\text{act}} R_{\text{sp}}(r_i) q_i^{(\text{act})} \text{ and } \tilde{k}_j = \Pi_{\text{act}} R_{\text{sp}}(r_j) k_j^{(\text{act})} \text{ from T7.2}$$

$$R_{\text{sp}}(r) \in SO(d_h) \text{ (T3.5) so } |R_{\text{sp}}(r)z| = |z| \text{ for all } z \in \mathbb{R}^{d_h}$$

Π_{act} is an orthogonal projector (T2.10) so $|\Pi_{\text{act}}z| \leq |z|$ for all $z \in \mathbb{R}^{d_h}$

$$\Rightarrow |\tilde{q}_i| \leq |q_i^{(\text{act})}|, \quad |\tilde{k}_j| \leq |k_j^{(\text{act})}|$$

$$\Rightarrow |\tilde{q}_i|^2 + |\tilde{k}_j|^2 \leq |q_i^{(\text{act})}|^2 + |k_j^{(\text{act})}|^2$$

(1) The estimator $\widehat{\text{SM}}_{n_{\text{rf}}}(\widehat{q}_i, \widehat{k}_j)$ is unbiased for $\exp(\alpha_{ij})$

10 Tier 8': Tight Assumptions and OOD Dominance

[T8'.0 : $|\cdot|_2, |\cdot|_F, \mathbb{S}^{d_h-1}$] $|M|_2 \stackrel{\mathcal{D}}{=} |M|$ (operator norm from T0.9), $|M|_F \stackrel{\mathcal{D}}{=} \sqrt{\text{tr}(M^\top M)}$ [deps : T0.9, T0.3]
 $\mathbb{S}^{d_h-1} \stackrel{\mathcal{D}}{=} \{z \in \mathbb{R}^{d_h} \mid |z| = 1\}$

[T8'.1 : $(\mathcal{X}, \mathcal{Y}), \mathcal{D}, \mathcal{T}, \Phi_{\#}\mathcal{D}$]

\mathcal{X}, \mathcal{Y} measurable, $(X, Y) \sim \mathcal{D}$

$\mathcal{T} \subseteq \{\Phi : \mathcal{X} \rightarrow \mathcal{X} \text{ bijective, measurable}\}$

$\Phi_{\#}\mathcal{D}$ (pushforward) satisfies $(X', Y') \sim \Phi_{\#}\mathcal{D} \Leftrightarrow (X', Y') = (\Phi \cdot X, Y)$, $(X, Y) \sim \mathcal{D} \Leftrightarrow (X', Y') \sim \Phi_{\#}\mathcal{D}$

$$(\Phi \cdot x)(u) \stackrel{\mathcal{D}}{=} x(\Phi^{-1}u)$$

[T8'.1bis : Parameterization map] $\exists \Pi_\Delta : \mathcal{T} \rightarrow \text{supp}(\nu)$ measurable s.t. $\Pi_{\text{act}} R_{\text{STR}}(\Pi_\Delta(\Phi)) \Pi_{\text{act}}$ matches the active

$$[T8'.2 : \text{Label invariance}] \quad \forall \Phi \in \mathcal{T} : \Pr_{(X,Y) \sim \mathcal{D}} [Y' = Y \text{ where } (X', Y') \sim \Phi \# \mathcal{D}] = 1 \quad [\text{deps : } T8'.1]$$

$$\boxed{[T8'.3 : \psi, \ell, \text{ bounded logits, } L]} \quad \begin{aligned} \psi : \mathbb{R}^{d_h} &\rightarrow \mathbb{R}^C, \quad \ell : \mathbb{R}^C \times \mathcal{Y} \rightarrow \mathbb{R}_+ \quad [\text{deps : } T8'.1] \\ \exists B_\psi > 0 : \sup_{x \in \mathcal{X}} |\psi(z_f(x))|_2 &\leq B_\psi \text{ (spectral control or fixed LayerNorm)} \\ \exists L > 0 : \forall y, \ell(\cdot, y) &\text{ is } L\text{-Lipschitz on } \{u \in \mathbb{R}^C : |u|_2 \leq B_\psi\} \end{aligned}$$

$$\boxed{[T8'.4 : z_f, \text{ measurable, } B_z]} \quad z_f : \mathcal{X} \rightarrow \mathbb{R}^{d_h} \text{ measurable}, \quad \sup_{x \in \mathcal{X}} |z_f(x)| \leq B_z \quad [\text{deps : } T0.9]$$

$$\boxed{[T8'.5 : \nu, \Delta, R_\Delta]} \quad \nu : \text{prob. law on } \mathbb{R}^{d_c}, \quad \Delta \sim \nu, \quad |\Delta|_1 \leq R_\Delta \text{ a.s.} \quad [\text{deps : } T0.1, T0.9]$$

$$\boxed{[T8'.5bis : r\text{-radius}]} \quad r \in \mathbb{R}^{d_c} \text{ ranges over } \{r : |r|_1 \leq R\}, \quad R > 0 \text{ fixed} \quad [\text{deps : } T0.9]$$

$$\boxed{[T8'.6 : \mathcal{H}_{\text{STRING}}, \mathcal{H}_{\text{ESPR}}]} \quad \begin{aligned} \mathcal{H}_{\text{STRING}} &\stackrel{\otimes}{=} \{f : P_{\text{sp}} \in \mathcal{A}, [L_a, L_b] = 0\} \quad [\text{deps : } T3.4, T0.9] \\ \mathcal{H}_{\text{ESPR}} &\stackrel{\otimes}{=} \{f : P_{\text{sp}} \in SO(d_h), |[L_a, L_b]| \leq \varepsilon_{\text{comm}}, \eta_{\text{mix}} \text{ as in T6'.0}\} \\ \mathcal{H}_{\text{STRING}} &\subset \mathcal{H}_{\text{ESPR}} \end{aligned}$$

$$\boxed{[T8'.7 : \text{IR}_{\text{spec}}(f)]} \quad \begin{aligned} \Pi_{\text{act}} \in \mathbb{R}^{d_h \times d_h} &\text{ is the orthogonal projector onto the first } 2m \text{ coordinates (T2.10), } \Pi_{\text{null}} = I \\ R_{\text{STR}}(\cdot) &\text{ acts as block } 2 \times 2 \text{ rotations on } m \text{ planes with angles from } \beta \text{ (identity on null block, T6'.0)} \\ \text{IR}_{\text{spec}}(f) &\stackrel{\otimes}{=} \mathbb{E}_{\Delta \sim \nu} \sup_{|r|_1 \leq R} \left\| \Pi_{\text{act}} R_{\text{sp}}(r)^\top R_{\text{sp}}(r + \Delta) \Pi_{\text{act}} - \Pi_{\text{act}} R_{\text{STR}}(\Delta) \right\| \end{aligned}$$

$$\boxed{[T8'.8 : C_{\text{comm}}, C_{\text{mix}}, C_{\mathcal{T}}, \rho]} \quad \begin{aligned} C_{\text{comm}} &\stackrel{\otimes}{=} \frac{1}{2} \varepsilon_{\text{comm}} (R + R_\Delta)^2 + C_{\text{BCH}} (2R + R_\Delta) \varepsilon_{\text{comm}}^2 (R + R_\Delta)^4 \quad [\text{deps : } T6'.1, T0.9] \\ \rho &\stackrel{\otimes}{=} |S|_2 < 1 \text{ for } P_{\text{sp}} = \text{cayley}(S) \text{ (Neumann regime, T6'.6)} \\ C_{\text{mix}} &\stackrel{\otimes}{=} \begin{cases} \frac{2}{(1-\rho)^3} \eta^2, & (\text{ESPR/Cayley/pure mix}), \\ \eta_{\text{mix}}, & (\text{model-free}) \end{cases} \\ C_{\mathcal{T}} &\stackrel{\otimes}{=} 1 \end{aligned}$$

$$\boxed{[T8'.8bis : \text{Small-angle/BCH regime}]} \quad \varepsilon_{\text{comm}}, \eta_{\text{mix}}, R, R_\Delta \text{ are small so that higher-order BCH terms are absorbed}$$

$$\boxed{[T8'.9 : \Delta_{\mathcal{T}}(f)]} \quad \Delta_{\mathcal{T}}(f) \stackrel{\otimes}{=} \mathbb{E}_{\Phi \sim \mu} \mathbb{E}_{(X,Y) \sim \mathcal{D}} |z_f(\Phi \cdot X) - z_f(X)| \quad [\text{deps : } T8'.1, T8'.4]$$

$$\boxed{[T8'.10 : \text{Geom} \Rightarrow \text{IR link (tight } C_z)]} z' = U'z, z = Uz, |z| \leq B_z \Rightarrow |z' - z| = |(U' - U)z| \leq |U' - U|_2 B_z \\ C_z \stackrel{\otimes}{=} B_z \Rightarrow \Delta_{\mathcal{T}}(f) \leq C_z \text{IR}_{\text{spec}}(f)$$

$$\boxed{[T8'.11 : R_{\text{train}}, R_{\mathcal{T}}]} R_{\text{train}}(f) \stackrel{\otimes}{=} \mathbb{E}_{(X,Y) \sim \mathcal{D}} \ell(\psi(z_f(X)), Y) \quad [\text{deps : } T8'.3, T8'.4, T8'.1] \\ R_{\mathcal{T}}(f) \stackrel{\otimes}{=} \mathbb{E}_{\Phi \sim \mu} \mathbb{E}_{(X,Y) \sim \mathcal{D}} \ell(\psi(z_f(\Phi \cdot X)), Y)$$

$$\boxed{[T8'.12 : \text{Shift-Lipschitz}]} R_{\mathcal{T}}(f) - R_{\text{train}}(f) \leq L \Delta_{\mathcal{T}}(f) \leq L C_z \text{IR}_{\text{spec}}(f) \quad [\text{deps : } T8'.3, T8'.10, T8'.11]$$

$$\boxed{[T8'.13 : \text{Upper control of } \text{IR}_{\text{spec}}]} \text{IR}_{\text{spec}}(f) \leq C_{\text{comm}} + C_{\text{mix}} \quad [\text{deps : } T6'.11, T8'.7, T8'.8, T8'.8bis]$$

$$\boxed{[T8'.14 : c_{\mathcal{T}} \text{ (noncommuting lower bound on active block)}]} m < d_h/2 \wedge \exists a \neq b : [L_a, L_b] \neq 0 \\ \exists \kappa_{\mathcal{T}} > 0 : \mathbb{E}_{\Delta \sim \nu} \inf_{|r|_1 \leq R} \|R_{\text{STR}}(r)^{\top} R_{\text{STR}}(r + \Delta) - R_{\text{STR}}(r)^{\top} R_{\text{STR}}(r)\| \geq \kappa_{\mathcal{T}} \\ \Rightarrow \inf_{f \in \mathcal{H}_{\text{STRING}}} \text{IR}_{\text{spec}}(f) \geq m$$

$$\boxed{[T8'.15 : \text{Reachability in ESPR}]} \forall \delta > 0 \exists f_{\delta} \in \mathcal{H}_{\text{ESPR}} : \varepsilon_{\text{comm}} \leq \delta, \eta_{\text{mix}} \leq \delta \quad [\text{deps : } T6'.11, T8'.6, T8'.7, T8'.8] \\ \Rightarrow \text{IR}_{\text{spec}}(f_{\delta}) \leq C_{\text{comm}}(\delta) + C_{\text{mix}}(\delta) \xrightarrow{\delta \rightarrow 0} 0$$

$$\boxed{[T8'.16 : \text{Infimum gap}]} \inf_{f \in \mathcal{H}_{\text{ESPR}}} R_{\mathcal{T}}(f) \leq \inf_{f \in \mathcal{H}_{\text{ESPR}}} (R_{\text{train}}(f) + L C_z \text{IR}_{\text{spec}}(f)) \quad [\text{deps : } T8'.12, T8'.14, T8'.15, T8'.16] \\ \leq \inf_{f \in \mathcal{H}_{\text{STRING}}} (R_{\text{train}}(f) + L C_z \text{IR}_{\text{spec}}(f)) - L C_z c_{\mathcal{T}} \\ \leq \inf_{f \in \mathcal{H}_{\text{STRING}}} R_{\mathcal{T}}(f) - L C_z c_{\mathcal{T}}$$

$$\boxed{[T8'.17 : \text{Optimization: Polyak--Lojasiewicz (PL)}]} \exists \mu > 0, L_g > 0 : \forall H \in \{\mathcal{H}_{\text{STRING}}, \mathcal{H}_{\text{ESPR}}\}, \mathcal{L}_H(f) \stackrel{\otimes}{=} \widehat{R}_{\text{train}}(f) \\ \text{satisfies PL: } \tfrac{1}{2} \|\nabla \mathcal{L}_H(f)\|^2 \geq \mu(\mathcal{L}_H(f)) \\ \nabla \mathcal{L}_H \text{ is } L_g\text{-Lipschitz} \Rightarrow \text{GD finds global minimizer } \hat{f}_{\lambda, H} \in \mathcal{H}_{\text{STRING}}$$

$$\boxed{[T8'.18 : \mathfrak{C}_n(H, \delta) \text{ (Rademacher; explicit scaling)}]} \exists S_0, \tilde{c} > 0 : \mathfrak{C}_n(H, \delta) \stackrel{\otimes}{=} c \left(\frac{S(H)}{\sqrt{n}} + \sqrt{n} \right) \\ S(\mathcal{H}_{\text{STRING}}) \leq S_0, \quad S(\mathcal{H}_{\text{ESPR}}) \leq S_0 + \tilde{c} \cdot (\text{rank}(E) + \varepsilon_{\text{comm}})$$

$$\boxed{[T8'.19 : \text{Uniform generalization}]} \Pr \left[\forall f \in H : R_{\text{train}}(f) \leq \widehat{R}_{\text{train}}(f) + \mathfrak{C}_n(H, \delta) \right] \geq 1 - \delta \quad [\text{deps : } T8'.18]$$

$$[T8'.20 : \text{Shift risk bound (uniform)}] \boxed{\Pr \left[\forall f \in H : R_{\mathcal{T}}(f) \leq \hat{R}_{\text{train}}(f) + \mathfrak{C}_n(H, \delta) + LC_z \text{IR}_{\text{spec}}(f) \right] \geq 1 - \delta}$$

$$[T8'.21 : \text{Achieved OOD gap}] \quad \epsilon_0 \stackrel{\mathcal{D}}{=} LC_z c_{\mathcal{T}}, \quad \Delta_{\text{gen}} \stackrel{\mathcal{D}}{=} \mathfrak{C}_n(\mathcal{H}_{\text{ESPR}}, \delta) - \mathfrak{C}_n(\mathcal{H}_{\text{STRING}}, \delta) \quad [\text{deps : } T8'.16, T8'.17]$$

$$\Pr \left[R_{\mathcal{T}}(\hat{f}_{\lambda, \mathcal{H}_{\text{ESPR}}}) \leq R_{\mathcal{T}}(\hat{f}_{\lambda, \mathcal{H}_{\text{STRING}}}) - \epsilon_0 + \Delta_{\text{gen}} \right] \geq 1 - \delta$$

$$[T8'.22 : \text{Worst-case corollary}] \quad R_{\text{sup}}(f) \stackrel{\mathcal{D}}{=} \sup_{\Phi \in \mathcal{T}} \mathbb{E}_{(X, Y) \sim \mathcal{D}} \ell(\psi(z_f))$$

$$R_{\text{sup}}(f) - R_{\text{train}}(f) \leq LC_z \sup_{\Delta \in \text{supp}(\nu)} \sup_{|r|_1 \leq R} \left\| \Pi_{\text{act}} R_{\text{sp}}(r)^{\top} R_{\text{sp}}(r + \Delta) \Pi_{\text{act}} - \Pi_{\text{act}} R_{\text{ST}}(r)^{\top} R_{\text{ST}}(r + \Delta) \Pi_{\text{act}} \right\|$$

and $\sup_{\Phi \in \mathcal{T}} = \sup_{\Delta \in \text{supp}(\nu)}$

$$[T8'.23 : \text{No benefit when commuting (sharpness)}] \quad \text{If } \forall \Delta \in \text{supp}(\nu) : [L_a, L_b] = 0 \wedge P_{\text{sp}} \in \mathcal{A} \Rightarrow \text{IR}_{\text{spec}}(f) = 0$$

$$\Rightarrow c_{\mathcal{T}} = 0, \epsilon_0 = 0$$