

# Quantified Robustness for Structured Rotary Attention

Internal note extending RoFormer and STRING analyses

November 9, 2025

## Abstract

RoFormer [1] and STRING [2] show that commuting skew generators give rise to relative attention logits when paired with structured post-rotations. This note isolates the genuinely novel contribution of our work: explicit robustness guarantees when the commutator and post-rotation assumptions are only approximately satisfied. We keep the notation of the original analyses and prove quantitative bounds for the BCH series, structured rotations, Cayley post-rotations, and, finally, the attention logits.

## 1 Setup and baseline results

We adopt the shared notation used by RoFormer and STRING. Let  $\{L_k\}_{k=1}^{d_c} \subset \mathfrak{so}(d_h)$  be skew-symmetric generators that commute pairwise in the idealized model. For a displacement  $r \in \mathbb{R}^{d_c}$  define

$$A(r) := \sum_{k=1}^{d_c} r_k L_k, \quad R_{\text{STR}}(r) := \exp(A(r)) \in SO(d_h). \quad (1)$$

The joint-block structure [2] provides an orthogonal matrix  $U$  that block-diagonalizes the generators and exposes an active subspace of dimension  $d_{\text{act}} = 2m$ . The corresponding projector is denoted by  $\Pi_{\text{act}} = U \begin{bmatrix} I_{2m} & 0 \\ 0 & 0 \end{bmatrix} U^\top$  and  $\Pi_{\text{null}} = I - \Pi_{\text{act}}$ . Following STRING, we allow for post-rotations  $P_{\text{sp}} \in SO(d_h)$  obtained from a Cayley transform

$$P_{\text{sp}} = (I - S)(I + S)^{-1}, \quad S = S_- + E, \quad (2)$$

where  $S_-^\top = -S_-$  captures the exact skew component induced by the algebra and  $E^\top = -E$  represents a small mixing perturbation.

**Assumption 1** (Baseline relative property). *In the unperturbed setting ( $E = 0$  and  $[L_a, L_b] = 0$ ) the active rotation obeys  $R_{\text{STR}}(r)^\top R_{\text{STR}}(s) = R_{\text{STR}}(s - r)$  and  $\Pi_{\text{act}} P_{\text{sp}} \Pi_{\text{act}} = \Pi_{\text{act}}$ .*

This captures the core theorems proved in RoFormer (Equations (11)–(13)) and the deterministic STRING analysis. The remainder of the note quantifies the deviation from this idealisation.

## Error parameters

We measure the size of the commutator defect by

$$\varepsilon_{\text{comm}} := \max_{a,b} \| [L_a, L_b] \| . \quad (3)$$

The block structure induces weights  $\lambda_{k,u}$  so that  $A(r)$  becomes block-diagonal in the active planes. Define the structural constant

$$\Lambda := \max_{u \in \{1, \dots, m\}} \sum_{k=1}^{d_c} |\lambda_{k,u}|, \quad R := \max_i \|r_i\|_1. \quad (4)$$

For the Cayley map we set  $\rho := \|S\|$  and  $\eta := \|E\|$ , assuming  $\rho < 1$  so that  $(I + S)$  is invertible. When the post-rotation is supplied externally (“model-free” case) we instead work with

$$\eta_{\text{mix}} := \max \{\|\Pi_{\text{null}} P_{\text{sp}} \Pi_{\text{act}}\|, \|\Pi_{\text{act}} P_{\text{sp}} \Pi_{\text{null}}\|\}. \quad (5)$$

## 2 Quantitative BCH control

The first technical ingredient bounds the size of the generator and provides an explicit constant for the Baker–Campbell–Hausdorff (BCH) remainder. This portion recapitulates the deterministic algebra from STRING but makes the constants explicit for later use.

**Lemma 1** (Generator norm). *For every displacement  $r \in \mathbb{R}^{d_c}$  we have  $\|A(r)\| \leq \Lambda \|r\|_1$ . Define  $C_{\text{BCH}}(r, s) := C_\star \exp(2\Lambda(\|r\|_1 + \|s\|_1))$  for an absolute  $C_\star$  that bounds higher-order BCH coefficients.*

**Lemma 2** (Quantitative BCH). *Let  $r, s \in \mathbb{R}^{d_c}$ . Then*

$$\|\log(\exp A(r) \exp A(s)) - (A(r) + A(s))\| \leq \frac{1}{2}\varepsilon_{\text{comm}} \|r\|_1 \|s\|_1 + C_{\text{BCH}}(r, s) \varepsilon_{\text{comm}}^2 \|r\|_1^2 \|s\|_1^2. \quad (6)$$

**Proposition 1** (Structured rotation stability). *With the same hypotheses we obtain*

$$\left\| R_{\text{STR}}(r)^\top R_{\text{STR}}(s) - R_{\text{STR}}(s - r) \right\| \leq C_1 \varepsilon_{\text{comm}} \|r\|_1 \|s\|_1 + C_2 \varepsilon_{\text{comm}}^2 \|r\|_1^2 \|s\|_1^2, \quad (7)$$

where  $C_1, C_2$  are universal (arising from the matrix exponential Lipschitz bound applied to Lemma 2).

## 3 Cayley post-rotations under perturbations

We now deviate from the exact STRING setup and allow  $S$  to acquire an additive perturbation  $E$ . All statements below are new: neither RoFormer nor STRING provide quantitative stability estimates in this regime.

**Lemma 3** (Cayley Lipschitzness). *If  $\|S\| \leq \rho < 1$  and  $\|E\| = \eta$ , then*

$$\|P_{\text{sp}}(S) - P_{\text{sp}}(S_-)\| \leq \frac{2}{(1-\rho)^2} \eta. \quad (8)$$

**Lemma 4** (Quadratic active mixing). *Let  $\Pi_{\text{act}} S \Pi_{\text{act}} = 0$  and  $\Pi_{\text{null}} S \Pi_{\text{null}} = S_{\text{null}}$ , as in the STRING block-structured Cayley factorisation. Then*

$$\|\Pi_{\text{act}} P_{\text{sp}}(S) \Pi_{\text{act}} - \Pi_{\text{act}}\| \leq \frac{2}{(1-\rho)^3} \eta^2. \quad (9)$$

**Lemma 5** (Model-free bound). *For an arbitrary  $P_{\text{sp}} \in SO(d_h)$  written in the STRING basis as  $U \begin{bmatrix} A & B \\ C & D \end{bmatrix} U^\top$  we have*

$$\|\Pi_{\text{act}} P_{\text{sp}} \Pi_{\text{act}} - \Pi_{\text{act}}\| \leq 2 \eta_{\text{mix}}. \quad (10)$$

Combining Lemmas 4 and 5 yields a unified parameter

$$\delta_{\text{mix}} := \begin{cases} \frac{2}{(1-\rho)^3} \eta^2, & \text{structured Cayley / ESPR case,} \\ 2\eta_{\text{mix}}, & \text{model-free case.} \end{cases} \quad (11)$$

**Corollary 1** (Active block approximation). *For every displacement  $r$ ,*

$$\|\Pi_{\text{act}} R_{\text{sp}}(r) \Pi_{\text{act}} - \Pi_{\text{act}} R_{\text{STR}}(r) \Pi_{\text{act}}\| \leq \delta_{\text{mix}}, \quad R_{\text{sp}}(r) := R_{\text{STR}}(r) P_{\text{sp}}. \quad (12)$$

## 4 Robust relative logits

Let  $q_i = W_Q x_i$  and  $k_j = W_K x_j$  denote the query and key projections used in RoFormer/STRING, and set  $q_i^{(\text{act})} = \Pi_{\text{act}} q_i$ ,  $k_j^{(\text{act})} = \Pi_{\text{act}} k_j$ . Define rotated vectors  $\tilde{q}_i = \Pi_{\text{act}} R_{\text{sp}}(r_i) q_i^{(\text{act})}$  and  $\tilde{k}_j = \Pi_{\text{act}} R_{\text{sp}}(r_j) k_j^{(\text{act})}$ . The attention logit at head dimension  $d_{\text{act}} = 2m$  is

$$\alpha_{ij} = \frac{1}{\sqrt{2m}} \tilde{q}_i^\top \tilde{k}_j. \quad (13)$$

Let

$$\alpha_{ij}^* := \frac{1}{\sqrt{2m}} (q_i^{(\text{act})})^\top (\Pi_{\text{act}} R_{\text{STR}}(r_j - r_i) \Pi_{\text{act}}) k_j^{(\text{act})}, \quad (14)$$

the ideal relative logit from the original proofs.

**Lemma 6** (Logit perturbation). *For all  $(i, j)$ ,*

$$|\alpha_{ij} - \alpha_{ij}^*| \leq \|q_i^{(\text{act})}\| \|k_j^{(\text{act})}\| \left( C_1 \varepsilon_{\text{comm}} \|r_i\|_1 \|r_j\|_1 + C_2 \varepsilon_{\text{comm}}^2 \|r_i\|_1^2 \|r_j\|_1^2 + \delta_{\text{mix}} \right). \quad (15)$$

**Theorem 1** (Quantified robust relative logits). *Let  $C_{\text{mix}} = 1$  and  $\delta_{\text{mix}} = \eta_{\text{mix}}$  in the model-free case, and  $C_{\text{mix}} = \frac{2}{(1-\rho)^3}$  together with  $\delta_{\text{mix}} = \eta^2$  in the structured Cayley case. Then*

$$\begin{aligned} |\alpha_{ij} - \alpha_{ij}^*| &\leq \left( \frac{1}{2} + C_{\text{BCH}}(r_i, r_j) \right) \varepsilon_{\text{comm}} \|r_i\|_1 \|r_j\|_1 \|q_i^{(\text{act})}\| \|k_j^{(\text{act})}\| \\ &\quad + C_{\text{mix}} \delta_{\text{mix}} \|q_i^{(\text{act})}\| \|k_j^{(\text{act})}\|. \end{aligned} \quad (16)$$

*Proof.* Combine Lemma 6 with Proposition 1 and Corollary 1.  $\square$

## 5 Discussion: overlap versus novelty

- Definitions of  $A(r)$ ,  $R_{\text{STR}}$ , and the use of joint block diagonalisation mirror the core statements of RoFormer and STRING; they are reintroduced here solely to fix notation.
- Lemma 1 matches the deterministic bounds implicit in STRING, but we expose explicit constants needed for downstream stability.
- Lemmas 3–5, Corollary 1, Lemma 6, and Theorem 1 are new: neither RoFormer nor STRING quantify how attention logits behave when generators only approximately commute and when the Cayley post-rotation leaks outside the active subspace.

## References

- [1] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. RoFormer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021.
- [2] Krzysztof Choromanski and Daniel Hardesty Lewis. STRING: Generalization of RoPE. Internal manuscript, 2023.