JOINT STRUCTURE SEARCH FOR TENSOR NETWORK OPERATORS INSPIRED BY SYMMETRY BREAKING

Anonymous authorsPaper under double-blind review

ABSTRACT

Tensor networks (TNs) offer a compact representation for high-dimensional operators in physics and machine learning. While TN structure search (TN-SS) has advanced model selection, prior work is limited to a *single* operator. Yet real systems, such as transformers and quantum circuits, would contain multiple coupled operators, where treating them independently or enforcing a single shared structure is fundamentally limiting. We introduce *joint TN-SS*, the first framework for multioperator structure search. Our physics-inspired algorithm runs in two phases: a symmetry phase, where standard TN-SS finds a shared structure capturing common inductive bias; and a symmetry-breaking phase, where operator-specific diversity emerges through greedy core masking, guided by task-explainable loss tolerances. Across tensor decomposition, parameter-efficient fine-tuning of LLMs, and quantum circuit optimization, joint TN-SS delivers more compact representations with equal or better accuracy than state-of-the-arts, with affordable search cost. These results demonstrate that symmetry-driven diversification offers a simple, general, and scalable solution to TN structure selection in multi-operator systems.

1 Introduction

Linear operators are the foundation of modern computation. In machine learning (ML) for example, multiple linear operators appear as fully connected, convolutional, and attention layers, but their high dimensionality often results in prohibitive computational and memory costs (Desislavov et al., 2023). Tensor networks (TNs) offer a structured and parameter-efficient way to represent linear operators (Orús, 2019; Memmel et al., 2024), driving growing interest in *tensor network operators* (TNOs) across ML (Novikov et al., 2015; Stoudenmire & Schwab, 2016; Richter et al., 2021; Yang et al., 2024; George et al., 2024), quantum physics, and beyond.

This growing adoption raises the central challenge of *tensor network structure search (TN-SS)*: how to select TNs' structure-related hyperparameters such as ranks (Kodryan et al., 2023; Zheng et al., 2024), topologies (Li & Sun, 2020), or permutations (Li et al., 2022; Zeng et al., 2024a). Because TN-SS is NP-hard (Hillar & Lim, 2013) and highly combinatorial, existing methods rely on heuristics and therefore focus almost exclusively on optimizing a *single* TN, or at most imposing one *shared* structure across multiple operators.

However, numerous studies have shown that the structural complexity of linear operators in deep neural networks varies substantially across layers. For instance, transformer weights often exhibit low-rank structure, but the low-rankness differs dramatically between layers (Jaiswal et al., 2024; Wang et al., 2024). Similar findings in adaptive pruning and low-rank adaptation further highlight that different operators respond unequally to factorization or sparsity (Frantar & Alistarh, 2023; Yang et al., 2024). These results expose a key weakness of existing TN-SS methods: *they ignore inter-layer heterogeneity by forcing either isolated optimization or rigid sharing across operators*. The open question is how to search for diverse yet efficient TN structures across multiple operators without exploding the search cost.

To address this challenge, we introduce the *first* formal formulation of joint TN-SS, together with a simple yet efficient algorithm inspired by the principle of symmetry breaking in physics (Lee, 1974). See Figure 1 for illustration. Just as symmetry breaking lowers energy in physical systems, we observe, in the context of joint TN-SS, a shared structure (i.e., a symmetric solution space) always exists across operators but is typically *suboptimal*. More efficient solutions emerge when

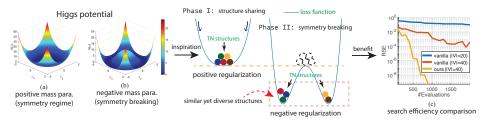


Figure 1: **Symmetry breaking (SB) in joint TN-SS.** Subfigures (a-b) illustrates the SB in physics (cms, 2022) and the inspiration for the proposed two-phase algorithm. Subfigure (c) demonstrates the search efficiency of the proposed methods compared to *vanilla* TN-SS methods; experimental details are provided in Section 4.1.

introducing controlled diversity close to that shared structure. Guided by this principle, we develop a two-phase algorithm: Phase I enforces structural sharing for search efficiency, and Phase II introduces operator-specific diversity through greedy-like core masking. A unified loss (5) governs the transition, with task-specific tolerances replacing opaque hyperparameters for practical explainability. Extensive numerical results across domains, including joint tensor decomposition, parameter-efficient fine-tuning of LLMs, and quantum circuit decomposition, confirm that our approach consistently produces more compact and effective TN representations than the existing TN-SS methods with minimal overhead. Our contributions are summarized as follows:

- We introduce the *first* formal framework for joint TN-SS, extending TN-SS to multi-operator systems and addressing heterogeneity ignored by prior work;
- We propose a symmetry-breaking algorithm that balances efficiency with expressiveness through structured diversification;
- We demonstrate broad applicability and consistent gains in parameter efficiency, scalability, and accuracy across domains.

1.1 RELATED WORKS

Tensor networks (TNs) in representing linear operators. TNs generalize classical tensor decompositions (Hitchcock, 1927; Tucker, 1966) (see reviews (Kolda & Bader, 2009; Cichocki et al., 2017)) and have become a standard tool for representing high-dimensional operators in both physics and ML (Novikov et al., 2015; Hou et al., 2019; Kossaifi et al., 2020; Chen et al., 2024a; Wang et al., 2024), where multiple operators are often coupled through nonlinearities or a unified loss. To capture diverse correlation patterns, a wide range of variants have been developed, from tensor train (also known as MPO) (Oseledets, 2011), tensor ring (TR) (Zhao et al., 2016), and tubal-SVD (Kilmer & Martin, 2011) to more flexible *circuit-like* TNs such as tree (Hackbusch & Kühn, 2009), stairs (Rudolph et al., 2023), brick-wall (Bensa & Žnidarič, 2021), and random circuits. These circuit-like designs are particularly relevant for quantum computing and show strong potential in recent ML studies (Chen et al., 2024a; Li et al., 2025). With this explosion of architectures, tensor network structure search has emerged as a central challenge in the field.

Tensor network structure search (TN-SS). TN-SS extends classical TN rank selection (Babacan et al., 2012; Rai et al., 2014; Zhao et al., 2015; Yokota et al., 2016) to richer structure-related hyperparameters such as topology and tensor permutations (Cheng et al., 2020; Mickelin & Karaman, 2020; Li et al., 2021a; Kodryan et al., 2023; Hayashi et al., 2019; Hashemizadeh et al., 2020; Li & Sun, 2020; Haberstich et al., 2023; Chen et al., 2024b; Zheng et al., 2024; Li et al., 2023; Zeng et al., 2024a; Guo et al., 2025). Despite steady progress, two major gaps remain. First, the structural modeling of those circuit-like TNs has never been systematically formulated. Second, existing TN-SS methods almost exclusively target a single TN or enforce one shared structure across operators, ignoring inter-operator heterogeneity. This work closes both gaps. We provide the first formal formulation for the structural representation of circuit-like TNs and introduce joint TN-SS, the first framework for simultaneous structure search across multiple operators within a unified system.

2 Basics of Joint TN-SS

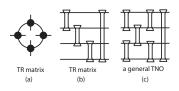
We introduce the key concepts of tensor networks, present the new formulation of joint TN-SS, and close with a brief review of symmetry breaking, the central inspiration for our search algorithm.

Notations. We use \mathbb{R} , \mathbb{C} , and $\mathbb{Z}_{>0}$ to denote the sets of real numbers, complex numbers, and positive integers, respectively. When both fields are admissible, we use $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$. Bold lowercase and uppercase letters (e.g., $\mathbf{x} \in \mathbb{F}^N$, $\mathbf{A} \in \mathbb{F}^{I \times J}$) represent vectors and matrices, respectively, and calligraphic letters (e.g., $\mathcal{A}, \mathcal{B} \in \mathbb{F}^{I_1 \times \cdots \times I_N}$) denote tensors. For a vector $\mathbf{x} \in \mathbb{F}^n$, $\operatorname{diag}(\mathbf{x}) \in \mathbb{F}^{n \times n}$ denotes the diagonal matrix whose diagonal entries are given by \mathbf{x} . The symbol \otimes is used for the tensor product, and \bigotimes denotes a sequential application of the products. We use \times for the Cartesian product *between sets* and \circ for function actions on graphs. The space $\mathbb{H} := \mathbb{F}^{I_1} \otimes \cdots \otimes \mathbb{F}^{I_N}$ is written as $\mathbb{F}^{I_1 \times \cdots \times I_N}$ or $\bigotimes_{n=1}^N \mathbb{F}^{I_n}$ for brevity. The notation $|\cdot|$ denotes norm-like quantities depending on context: e.g., $|\phi|$ is the absolute value for $\phi \in \mathbb{C}$, and |G| denotes the order (number of vertices) of a graph G. Last, we define $[K] := \{1, 2, \ldots, K\}$ for convenience.

2.1 Tensor Network Operators

Tensor network (TN). We adopt the formal definition of TNs by Ye & Lim (2019), where a TN, denoted by $tns_{\mathbb{H}}(G,r,p)$, is defined as a set of tensors in a space \mathbb{H} . Here, G=(V,E) is a graph that defines the TN topology (Li & Sun, 2020), where the vertex set V corresponds to a collection of core tensors, while each edge in E indicates a tensor contraction between connected tensors. The function $r:E\to\mathbb{Z}_{>0}$ assigns a positive integer to each edge, specifying the TN ranks. Additionally, the mapping $p:\{\mathbb{F}^{I_n}\}_{n=1}^N\to V$ associates each input subspace \mathbb{F}^{I_n} with a vertex of G, specifying the TN permutation (Li et al., 2022). For a more detailed introduction to TNs, we refer readers to the overview (Cichocki et al., 2016).

TNOs thus compactly represent high-dimensional operators of size $I^Q \times I^Q$ through the contraction of a sequence of lower-dimensional operators of size $I^K \times I^K$. Because all core tensors are constrained to be isomorphic to square matrices, *circuit graphs* (see right figures for example) offer a more intuitive visualization of complex TNOs than traditional TN diagrams (Li & Sun, 2020; Li et al., 2022; 2023). As illustrated on the right, subfigures (a) and (b) show a TR matrix



in both traditional and circuit representations, while subfigure (c) depicts a more general TNO model.

TNO system. Multiple TNOs build up a TNO system. Given Q, K, I, we formulate a M-dimensional TNO system as follows:

$$ts(G_1, G_2, \dots, G_M; Q, K, I) := tno(G_1) \times tno(G_2) \times \dots \times tno(G_M), \tag{1}$$

where $tno(G_m) := tno(G_m; Q, K, I)$ for all m. This construction creates a unified combinatorial space of operators that share the same parameters (Q, K, I) but may differ in topology through graphs $\{G_m\}_{m=1}^M$. We remark that each element in $ts(G_1, G_2, \ldots, G_M; Q, K, I)$ can be viewed as a vector of operators, which makes the framework able to model multiple, distinct linear operators that jointly define a system. For example, in computer vision, it can capture multi-view video tensors where each view is represented by a tensor; in deep learning, it can model a collection of linear layers in a neural network, where each layer may benefit from a different structure; and in quantum computing, it represents different unitary blocks in a circuit, which must be optimized jointly but are not identical. This perspective highlights that a TNO system encodes heterogeneity across operators in a principled way, going beyond the treatment of one TN in isolation.

2.2 Joint TN-SS: Extending TN-SS to TNO systems

Building on prior studies in TN-SS (Li & Sun, 2020; Li et al., 2022; 2023), we move beyond the conventional task of optimizing a *single* TN to the more general and practically relevant setting of *joint TN-SS*. Formally, given parameters Q, K, I, the goal is to solve the bi-level optimization problem:

$$\min_{\{G_m\}_{m=1}^M} \min_{\mathcal{X}^M} \pi_D(\mathcal{X}^M)
s.t. \quad G_m \in \mathbb{G}_m, \ m \in [M], \quad \mathcal{X}^M \in ts(G_1, G_2, \dots, G_M; Q, K, I)$$
(2)

where \mathbb{G}_m denotes the set of valid graphs satisfying the TNO constraints in Section 2.1, and $\pi_D(\cdot)$: $\mathbb{F}^{I^Q \times I^Q} \to \mathbb{R}$ is a task-specific loss function defined on some data D. Note that in the existing TN-SS literature (Li & Sun, 2020; Li et al., 2022; 2023; Zeng et al., 2024b;a) the graph G is typically represented by adjacency matrices. However, such a representation makes it difficult to enforce the required TNO structural constraints. To overcome this, we introduce a *vertex-indexed incidence (VI) matrix* representation, which naturally encodes circuit-like structures. Specifically, for each TNO with graph G_m , we define:

$$\mathbf{I}_{m} = \begin{pmatrix} i_{1,1}^{m} & i_{1,2}^{m} & \cdots & i_{1,L_{m}}^{m} \\ i_{2,1}^{m} & i_{2,2}^{m} & \cdots & i_{2,L_{m}}^{m} \\ \vdots & \vdots & \ddots & \vdots \\ i_{K,1}^{m} & i_{K,2}^{m} & \cdots & i_{K,L_{m}}^{m} \end{pmatrix},$$
(3)

where L_m is the order of G_m , and the entries satisfy $i_{k,l}^m \in [Q]$ and $i_{k_1,l}^m \neq i_{k_2,l}^m$ for all $l \in [L_m]$, $k_1, k_2 \in [K]$ with $k_1 \neq k_2$.

2.3 A Brief Review of Symmetry Breaking and Connection to Joint TN-SS

Our approach to joint TN-SS is motivated by the physical principle of symmetry breaking (SB), where a system transitions from a symmetric to an asymmetric state (Lee, 1974). A classical example is the Higgs potential (Melo, 2017):

$$V(\phi) = \lambda \phi^4 + \mu^2 \phi^2,\tag{4}$$

where $\phi = \sqrt{\phi_x^2 + \phi_y^2}$ is the field magnitude, $\lambda > 0$ ensures stability, and μ^2 controls the curvature. The system's behavior depends critically on the sign of μ^2 . When $\mu^2 > 0$, the potential has a unique symmetric minimum at $\phi = 0$, where all trajectories collapse to the same state. When $\mu^2 < 0$ ($\mu \in \mathbb{C}$), the potential becomes a "Mexican hat" (see Figure 1 (b)) with infinitely many degenerate minima along a ring, forcing the system to spontaneously choose one, thereby breaking symmetry.

Connection to joint TN-SS. This physics metaphor directly informs our algorithmic design. In joint TN-SS, the challenge is the combinatorial explosion of possible structures as the number of TNOs M grows. SB suggests a natural way to manage this complexity:

- Symmetry regime ($\mu^2 > 0$) all TNOs share a common structure, collapsing the search to standard TN-SS and giving a tractable starting point;
- SB regime ($\mu^2 < 0$) diversity is gradually introduced as the system "selects" different operator-specific structures, analogous to particles settling in different minima of the Mexican hat.

This phase transition, from enforced symmetry to controlled asymmetry, provides a principled mechanism to reduce the search space initially, then progressively enrich it. Instead of facing the full exponential complexity of unconstrained joint search, our algorithm leverages SB to unlock diversity only when needed, enabling efficient exploration of heterogeneous TNO structures.

3 THE PROPOSED APPROACH

3.1 REFORMULATION INSPIRED BY SYMMETRY BREAKING

For brevity, we write $qs(G_1, \ldots, G_M) := qs(G_1, \ldots, G_M; Q, K, I)$. The joint TN-SS problem (2) is reformulated by imposing a common TN structure and introducing controlled structure diversity, in direct analogy to the Higgs potential (4):

$$\min_{G \in \mathbb{G}_q, \{\epsilon_m\}_{m=1}^M} \min_{\mathcal{X}^M \in ts(G_1, G_2, \dots, G_M)} \pi_D(\mathcal{X}^M) + \lambda \sum_{m=1}^M |\epsilon_m|,$$

$$s.t. G_m = G \circ \epsilon_m, \forall m \in [M].$$

$$(5)$$

Here, \mathbb{G}_q denotes the set of valid graphs satisfying the TNO constraints. Each $\epsilon_m:\mathbb{G}_q\to\mathbb{G}_q$ represents a graph operation applied to \mathbb{G} , while $\lambda\in\mathbb{R}$ controls the strength of a regularization term penalizing the "degree" of the structure diversity. Compared to the original formulation (2), (5) assumes that all structures $\{G_m\}$ are generated from a *common graph* G with small perturbations $\{\epsilon_m\}$. The regularization balances the trade-off between structural sharing (when ϵ_m is close to identity) and diversity (when "big" ϵ_m 's are encouraged), thereby mirroring the transition from symmetry to symmetry breaking.

It is worth emphasizing that imposing the common graph G does *not* limit expressiveness. By the universal approximation property of circuit-like TNs (Barenco et al., 1995; Möttönen et al., 2004), we know there always exists a sufficiently high order G that can represent all operator-specific graphs $\{G_m\}$. While highly diverse G_m may in principle require a large G, our experiments show that in practice one can almost always find an affordable G close to optimal. From this shared structure, the perturbations ϵ_m efficiently specialize into diverse G_m with more compact representation for a TNO system. Hence, the common graph is not a restriction but a trick for efficiency: it reduces the initial combinatorial burden while still enabling the emergence of heterogeneous, task-adapted structures.

Practical choice of ϵ_m . Empirically, we use $core\ masking$ as a simple yet powerful mechanism for modeling ϵ_m in (5), offering both efficiency and effectiveness at negligible cost. As illustrated on the right, core masking corresponds to removing selected core tensors from a TNO, and the regularization term $|\epsilon_m|$ in (5) is defined as the number of core tensors masked by ϵ_m . Formally, let $\mathbf{w}_m \in \{0,1\}^{L_m}$ be a binary mask over the L_m core tensors of the m-th TNO of \mathcal{X}^M with graph $G_m = (V_m, E_m)$ and VI matrix $\mathbf{I}_m \in [Q]^{K \times L_m}$. We have $\epsilon_m(\mathbf{w}_m) : (\mathbf{I}_m, \{\mathcal{G}_{m,\ell}\}) \mapsto (\mathbf{I}_m(\mathrm{diag}(\mathbf{1} - \mathbf{w}_m)), \{\tilde{\mathcal{G}}_{m,\ell}\})$, where $\tilde{\mathcal{G}}_{m,\ell} = \mathcal{G}_{m,\ell}$ if $w_{m,\ell} = 0$ and $\tilde{\mathcal{G}}_{m,\ell}$ equals identity otherwise, and thus $|\epsilon_m| = \|\mathbf{w}_m\|_0 = \sum_{\ell=1}^{L_m} w_{m,\ell}$.

Hyperparameter λ and phase transition. Note that the tuning parameter λ in (5) plays the same role as the mass parameter in the Higgs potential (4), shaping the loss landscape and governing the transition between the symmetry and SB regimes. When the common graph G is sufficiently expressive to have a smaller value of π_D and $\lambda \geq 0$, the system remains in a symmetry regime: the optimal graphs $\{G_m\}_{m=1}^M$ must collapse to the single common graph G, since any non-trivial core masking would be penalized. In this case, each ϵ_m^* equals the identity, and the joint TN-SS problem reduces to vanilla TN-SS. This phase provides an efficient initialization by enforcing structural sharing.

When $\lambda < 0$, the system enters the SB regime. While negative regularization is unusual in ML, in physics it is precisely what drives the Mexican-hat potential to favor multiple nonzero minima. In our setting, it encourages independent, operator-specific core masking across the $\{G_m\}$, thereby breaking symmetry and introducing structural diversity. This transition expands the search space in a controlled manner, guiding the system toward more compact, heterogeneous configurations that often achieve lower loss.

3.2 Algorithm

Motivated by symmetry breaking, we propose a two-phase algorithm that first enforces a shared structure for efficiency and then gradually introduces diversity via core masking.

Algorithm 1 The proposed algorithm for joint TN-SS (sketch)

```
271
             1: Initialize: parameters Q, K, I, tolerance \eta_1, \eta_2 with \eta_1 \leq \eta_2.
272
            2: Define f(\cdot) = \min_{\mathcal{X}^M \in ts(\cdot)} \pi_D(\mathcal{X}^M).
273
            3: Obtain common structure G \leftarrow \text{TN-SS}(f) with f(G) \leq \eta_1
                                                                                                            ▶ Phase I: symmetry
274
            4: Set G_m \leftarrow G with G_m = (V_m, E_m) for all m \in [M].
275
                                                                                             ▶ Phase II: symmetry breaking!
276
                     Pick one core tensor v from \bigcup_{m=1}^{M} V_m at random.
Construct masked candidates G'_m = G_m \circ \epsilon(v) for all m. \triangleright v is masked from G_m.
            6:
277
            7:
278
                     Evaluate f' = f(\{G'_m\}).
            8:
279
            9:
                     if f' \leq \eta_2 then
                          Accept: G_m \leftarrow G'_m for all m.
           10:
           11:
                     end if
281
           12: until no further improvement is found
282
           13: Output: \{G_m\}_{m=1}^{M}
283
```

As analyzed in Sec. 2.3, the optimal solution of (5) in the symmetry regime collapses to a common graph G. Thus, in Phase I we set ϵ_m^* to be identity for all $m \in [M]$ and run any TN-SS method to find G with $f(G) \leq \eta_1$. A small η_1 ensures G is expressive, though possibly redundant. Once G is fixed, the algorithm enters Phase II (the SB regime), where structural diversity is introduced by iteratively masking cores: at each step, a vertex v (corresponding to a core tensor) is sampled from $\bigcup_{m=1}^M V_m$, across all TNOs, and the updated graphs are evaluated. If the task loss remains within tolerance η_2 , the change is accepted; otherwise, a new masked candidate is tried. The process stops when no further beneficial masks are found. Note that the gap $\eta_2 - \eta_1$ quantifies the performance budget for pruning redundant core tensors and encouraging operator-specific diversity.

Crucially, λ in (5) is never tuned directly in the algorithm. Its effect is implicitly encoded by (η_1, η_2) , which provide an explainable, task-specific rule for when symmetry should be enforced and when diversity is allowed. This yields a practical phase transition mechanism that avoids the exponential cost of unconstrained joint search. Moreover, the greedy masking ensures the number of evaluations grows only linearly with M in Phase II. Experiments in Sec. 4 confirm that our method achieves more compact TN representations with little overhead.

Finally, the following perturbation bound illustrates that the loss increase from masking a core is controlled by its operator–Schmidt rank:

Proposition 3.1 (Perturbation analysis). Let $\mathcal{Z}^M = \{Z_1, \dots, Z_M\}$ be a system of TNOs with $\mathcal{Z}^M \in \arg\min_{\mathcal{Y}^M \in ts(G_1, \dots, G_M; Q, K, I)} \pi_D(\mathcal{Y}^M)$. Suppose one operator Z_j contains a maskable core \mathcal{G} with $\|\mathcal{G}\|_F = 1$ and operator–Schmidt rank S. Let $\mathcal{Z}^{M,\star}$ denote the TNO system obtained by replacing Z_j with its best masked version $Y^\star \in \arg\min_{Y \in tno(G_{\text{mask}}; Q, K, I)} \pi_D(\{Z_1, \dots, Y, \dots, Z_M\})$, where $G_{\text{mask}} = G_j \circ \epsilon(\mathcal{G})$. If π_D is L-Lipschitz in $\|\cdot\|_F$ with respect to each operator, then

$$|\pi_D(\mathcal{Z}^{M,\star}) - \pi_D(\mathcal{Z}^M)| \le LC\sqrt{S-1},$$

where C depends only on the contracted environment of G. In particular, if S=1, masking G does not increase the loss.

The proof is in the Appendix D. Prop. 3.1 shows that rank-one cores can be pruned without loss, while higher ranks incur bounded perturbations scaling with $\sqrt{S-1}$. This provides a sufficient criterion for identifying redundant core tensors, though we argue in the Appendix that low Schmidt rank is not necessary, owing to the universality of TN approximations.

4 Numerical Results

We evaluate the proposed algorithm on diverse domains: joint tensor decomposition, parameter-efficient fine-tuning of LLMs, and quantum circuit optimization, to demonstrate its effectiveness. Due to space limits, we report only key settings in the main text; full experimental details are deferred to the Appendix.

¹Appendix A.3 provides guidelines for adapting TNGA (Li & Sun, 2020), TNLS (Li et al., 2022), TnALE (Li et al., 2023), and Greedy (Hashemizadeh et al., 2020) to TNOs.

4.1 Joint Tensor Decomposition on Synthetic Data

We evaluate first the performance of our algorithm on tensor decomposition using synthetic tensors.

Data preparation. We choose Q=5, K=2, I=2 and M=4 for convenience, following (1), This choice corresponds to optimizing structures for tensors of order-10, which is *markedly larger* than in most prior TN-SS benchmarks, providing a challenging yet tractable testbed for evaluating our joint TN-SS algorithm. Additional results on varying Q, I, and M are reported in Appendix A.6. Next, we choose four circuit-like TN structures widely used in ML or physics, including tensor train (TT, Oseledets 2011), "stairs-like" TN (Stairs, Rudolph et al. 2023), "brick-wall" TN (Brick, Bensa & Žnidarič 2021), and randomly-connected TN (Rand.), to construct the common graph G in $(5)^2$, and then randomly mask core tensors from G individually to obtain the diverse $\{G_m\}_{m=1}^M$ for each TNO in the system. Once the structures are determined, we follow the real-valued and i.i.d. Gaussian distribution $\mathcal{N}(0,1)$ to generate the values of core tensors. Synthetic tensors are thus obtained by contracting the core tensors.

Setup. For our algorithm, we set $\pi_D(\mathcal{X}^M)$ to be the relative squared error (RSE, Li & Sun 2020) calculated globally among all TNOs, and use Adam (Kingma & Ba, 2014) to optimize the value of core tensors. The setting of L is determined differently for each data, following a basic principle that L should be relatively large to achieve a satisfactory approximation error in Phase I. The tolerances η_1, η_2 are set to 5×10^{-10} and 10^{-6} , respectively (the sensitivity analysis of L and η_2 are given in Appendix A.6). We regard the approximation as successful if $RSE \leq 10^{-6}$.

Baselines. For comparison, we adopt four state-of-the-art TN-SS algorithms: TNGA (Li & Sun, 2020), TNLS (Li et al., 2022), TnALE (Li et al., 2023), and Greedy (Hashemizadeh et al., 2020). They also serve as the backbone optimizers for the common graph G in Phase I of our approach. To handle TNO structures, we introduce minimal adaptations on the baselines while keeping each algorithm's original search dynamics intact. Details of these modifications are provided in Appendix A.3.

Results. Figure 2 (a–c) reports results across different TN-SS baselines. We first see from Figure 2 (a) that Phase II significantly reduces the total number of core tensors, resulting in more compact TN representations than the existing TN-SS methods. Furthermore, we see from Figure 2 (b) that such compression preserves accuracy. The approximation error remains at $RSE \leq 10^{-6}$ in nearly all cases. Third, Figure 2 (c) illustrates that the evaluation cost in Phase II (dark bars) is consistently lower than in Phase II, confirming the efficiency of symmetry breaking.

For fairness, we also tested vanilla TN-SS with the assumption of sharing structure among TNOs with the same |V| as our methods, as done in many existing works. Figure 2(b) shows, this severely limits performance: none of the methods in this setting reach $RSE \leq 10^{-6}$, except TNLS on TT. This contrast highlights that structural diversity is important for achieving the full expressive power of TNOs.

Phase transition during the search. Figure 2 (d–f) shows the search dynamics on **Rand.** with our approach (that uses TNGA in Phase I). A clear phase transition emerges: once entering Phase II (the SB regime), the number of core tensors steadily decreases (Figure 2(d)) while RSE remains stable (Figure 2(e)). At the same time, the per-iteration evaluation cost drops sharply (Figure 2(f)), highlighting both efficiency and accuracy of the transition.

Experimental details for cover image. We further test a plausible baseline where vanilla TN-SS tackles joint TN-SS by composing all TN structures together into a single large graph. As shown in Figure 1(c), vanilla TN-SS (TNGA on **Rand.**, blue/red curves) is significantly less efficient than our approach (yellow curve). This highlights the power of the shared-structure prior in (5), which guides the search toward diverse yet compact solutions at far lower cost.

4.2 PARAMETER EFFICIENT FINE-TUNING (PEFT) FOR LLMS

We next demonstrate that our approach improves parameter efficiency in the PEFT task for LLMs. For disclaimer, this experiment is *not* intended to propose a new practical PEFT method, but rather to highlight how joint TN-SS itself contributes to parameter-efficient representations.

²Note that our focus is on TN structures tailored for high-order tensors; thus, classical tensor decomposition models such as CP, Tucker and tSVD are excluded from our experiments.

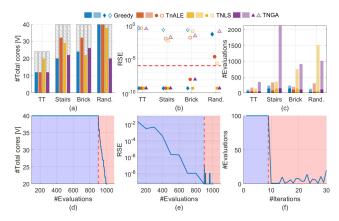


Figure 2: Experimental results of tensor decomposition. (a) Shaded and colored bars show the total number of core tensors in Phase I and Phase II, respectively; (b) Solid and hollow markers indicate the RSE achieved by the proposed algorithm after Phase II and by the vanilla TN-SS with |V| matched to that obtained by the proposed algorithm. The red dotted line marks the $RSE=10^{-6}$ threshold; (c) Light and dark bars denote the total number of evaluations in Phase I and Phase II, respectively; (d)-(f) Search dynamics of **Rand.** using TNGA, where blue and red regions correspond to Phase I and Phase II.

Table 1: Test performance of PEFT. The results marked with "*" are from Li et al. (2024). The details for QuanTA-6/4/2 are described in Appendix B.2.

		1	1				
Method	#Params	PIQA	SIQA	OBQA	ARC-e	ARC-c	Avg.
LoRA*	3.200%	82.1	69.9	80.4	73.8	50.9	71.4
DoRA*	3.200%	82.7	74.1	80.6	76.5	59.8	74.7
QuanTA-6	0.041%	79.9	75.9	80.0	84.8	63.3	76.8
Ours	0.031%	80.4	76.2	80.2	84.7	63.7	77.0
QuanTA-4	0.024%	79.2	73.5	77.2	84.4	62.6	75.4
Ours	0.024%	80.3	75.1	78.4	84.2	63.4	76.3
QuanTA-2	0.017%	78.1	72.5	74.8	82.1	57.1	72.9
Ours	0.017%	79.7	72.7	78.0	83.1	59.9	74.7

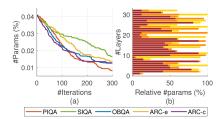


Figure 3: Param. reduction dynamics. (a) Number of param. over iterations; (b) Layer-wise parameter allocation on ARC-c, where the three colors indicate different stages in search.

Setup. We build upon QuanTA (Chen et al., 2024a), a state-of-the-art PEFT method for LLMs utilizing TNO. In this experiment, we enhance it by applying our algorithms to optimize TNO structures, for the goal of reducing the number of fine-tuned parameters. We conduct the experiment on five commonsense reasoning datasets with Llama2-7B model (Touvron et al., 2023) (involving the TNO system of dimension M=32), and follow most of the settings as the work (Chen et al., 2024a), where each fine-tuned weight with size 4096×4096 is recognized as a tensor of order-8 with size $16\times8\times8\times4\times16\times8\times8\times4$. In our algorithm, we directly use the structure in QuanTA as the common graph G, from which we deploy the proposed algorithm from Phase II.

Results. Table 1 summarizes the test results, including baseline comparisons with LoRA (Hu et al., 2021) and DoRA (Liu et al., 2024a). Our approach achieves superior performance with the same or fewer fine-tuning parameters. Figure 3(a) shows the search dynamics across five datasets, where joint TN-SS consistently reduces the parameter count throughout the search. Figure 3(b) further illustrates the layer-wise parameter allocation on ARC-c, demonstrating that our approach allows *non-uniform* parameter distribution across layers, maximizing overall parameter efficiency. This observation is consistent with findings discussed in Section 1, which suggest that different weights within a system often exhibit varying structural complexities. More results are shown in Appendix B.4.

 $^{^{3}}$ Note that QuanTA slightly generalizes the definition of TNOs by varying the parameter I. We adopt the same experimental settings, as the variation in I does not impact the search process of our algorithm.

Table 2: Result for quantum circuit optimization using joint TN-SS.

Case	#Params of Brick-wall (baseline)		Classic QFT circuit		Ours		
Cusc	original opreator	#Params	Fidelity [↑]	#Params	Fidelity [†]	#Params	Fidelity [↑]
Synthetic $(Q = 8)$ Synthetic $(Q = 12)$	$65536 \\ 1.68 \times 10^{7}$	560 528	$> 1 - 10^{-4}$ $> 1 - 10^{-4}$	/	/ /	208 176	$> 1 - 10^{-4}$ $> 1 - 10^{-4}$
$\begin{array}{c} \hline \text{QFT} \ (Q=4) \\ \text{QFT} \ (Q=6) \end{array}$	256 4096	240 800	$> 1 - 10^{-4}$ 0.9859	96 240	$> 1 - 10^{-4}$ $> 1 - 10^{-4}$	96 288	$> 1 - 10^{-4}$ $> 1 - 10^{-4}$

4.3 MEMORY-EFFICIENT QUANTUM CIRCUIT OPTIMIZATION

Quantum circuit optimization aims to represent high-dimensional unitary operators with a compact circuit of lower-dimensional unitary components. Standard approaches often use multi-layer brick-wall structures, which can be highly redundant and inefficient. Here we conduct an illustrative study showing that joint TN-SS can find more compact circuits with fewer components, enhancing resource efficiency or potentially shallower architecture.

Setup. We evaluate both synthetic unitary operators (generated with Haar-random cores, K=2, I=2, and Q=8,12) and the quantum Fourier transform (QFT) (Camps et al., 2021; Chen et al., 2023),where the same settings of K,I,Q are applied for QFT. For baselines, we use brick-wall TNOs and, for QFT, the classic hand-crafted QFT circuit. In our method, repeated brick-wall layers are treated as the common graph G, and structural diversity emerges through symmetry breaking.

Results. Table 2 shows that our approach consistently produces more compact TNO representations than the brick-wall baseline while maintaining fidelity above $1-10^{-4}$. For synthetic operators, we achieve up to a $3\times$ reduction in parameters. For QFT, our method recovers circuits with parameter counts on par with the classic design, yet without SWAP gates, a critical advantage for noise mitigation on quantum hardware. Figure 4 further illustrates this: for Q=4, the result matches the conventional circuit (Camps et al., 2021; Chen et al., 2023), while for Q=6 it discovers a distinctly different and more compact architecture using only local operators. These findings demonstrate that joint TN-SS not only advances efficiency in ML tasks but also provides a new path to discovering efficient quantum circuits. Full experimental details are given in Appendix C.

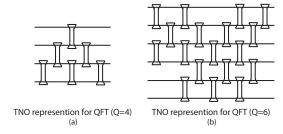


Figure 4: Optimized TNO representation for the QFT operators.

5 CONCLUDING REMARKS

In this work, we introduced *joint TN-SS*, a previously unexplored extension of tensor network structure search. Inspired by the principle of symmetry breaking, we proposed a simple yet effective algorithm to tackle the combinatorial complexity inherent in joint TN-SS. Extensive experiments across diverse tasks, including tensor decomposition, LLM fine-tuning, and quantum circuit optimization, demonstrate the effectiveness and efficiency of our approach.

Limitations. Due to resource constraints, our experiments on quantum circuit optimization were limited to small-scale settings. Future work will extend to larger scales. Additionally, we observed instability in widely used optimization methods, such as gradient-based and SVD-based approaches (Schollwöck, 2005), when applied to tensor decomposition and quantum circuit optimization. Enhancing the robustness of joint TN-SS will therefore be an important direction for future study.

REPRODUCIBILITY STATEMENT

We have made substantial efforts to ensure the reproducibility of our work. Detailed descriptions of our proposed algorithm are provided in Section 3.2. Detailed settings for each experiment—including datasets, hyperparameter configurations, and implementation details—are presented in Appendix A, B and C. We have included the source code for tensor decomposition and LLM fine-tuning in the supplementary materials (the code for quantum circuit optimization will be made available after publication). For theoretical results (i.e., Proposition 3.1), the complete proof is presented in Appendix D. We encourage readers to refer to these sections for comprehensive information necessary to reproduce our work.

REFERENCES

- A portrait of the Higgs boson by the CMS experiment ten years after the discovery. *Nature*, 607 (7917):60–68, 2022.
- S. Derin Babacan, Martin Luessi, Rafael Molina, and Aggelos K. Katsaggelos. Sparse Bayesian Methods for Low-Rank Matrix Estimation. *IEEE Transactions on Signal Processing*, 60(8): 3964–3977, 2012. doi: 10.1109/TSP.2012.2197748.
- Adriano Barenco, Charles H Bennett, Richard Cleve, David P DiVincenzo, Norman Margolus, Peter Shor, Tycho Sleator, John A Smolin, and Harald Weinfurter. Elementary gates for quantum computation. *Physical review A*, 52(5):3457, 1995.
- Jaš Bensa and Marko Žnidarič. Fastest local entanglement scrambler, multistage thermalization, and a non-hermitian phantom. *Physical Review X*, 11(3):031019, 2021.
- Daan Camps, Roel Van Beeumen, and Chao Yang. Quantum fourier transform revisited. *Numerical Linear Algebra with Applications*, 28(1):e2331, 2021.
- Jielun Chen, EM Stoudenmire, and Steven R White. Quantum fourier transform has small entanglement. *PRX Quantum*, 4(4):040318, 2023.
- Zhuo Chen, Rumen Dangovski, Charlotte Loh, Owen Dugan, Di Luo, and Marin Soljačić. QuanTA: Efficient High-Rank Fine-Tuning of LLMs with Quantum-Informed Tensor Adaptation. *arXiv* preprint arXiv:2406.00132, 2024a.
- Ziang Chen, Jianfeng Lu, and Anru Zhang. One-dimensional Tensor Network Recovery. *SIAM Journal on Matrix Analysis and Applications*, 45(3):1217–1244, 2024b.
- Zhiyu Cheng, Baopu Li, Yanwen Fan, and Yingze Bao. A novel rank selection scheme in tensor ring decomposition based on reinforcement learning for deep neural networks. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3292–3296. IEEE, 2020.
- Andrzej Cichocki, Namgil Lee, Ivan Oseledets, Anh-Huy Phan, Qibin Zhao, Danilo P Mandic, et al. Tensor networks for dimensionality reduction and large-scale optimization: Part 1 low-rank tensor decompositions. *Foundations and Trends® in Machine Learning*, 9(4-5):249–429, 2016.
- Andrzej Cichocki, Anh-Huy Phan, Qibin Zhao, Namgil Lee, Ivan Oseledets, Masashi Sugiyama, Danilo P Mandic, et al. Tensor networks for dimensionality reduction and large-scale optimization: Part 2 applications and future perspectives. *Foundations and Trends® in Machine Learning*, 9(6): 431–673, 2017.
- Radosvet Desislavov, Fernando Martínez-Plumed, and José Hernández-Orallo. Trends in ai inference energy consumption: Beyond the performance-vs-parameter laws of deep learning. *Sustainable Computing: Informatics and Systems*, 38:100857, 2023.
- Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International conference on machine learning*, pp. 10323–10337. PMLR, 2023.

- Robert Joseph George, David Pitt, Jiawei Zhao, Jean Kossaifi, Cheng Luo, Yuandong Tian, and Anima Anandkumar. Tensor-GaLore: Memory-Efficient Training via Gradient Tensor Decomposition. In OPT 2024: Optimization for Machine Learning, 2024.
 - Zheng Guo, Aditya Deshpande, Brian Kiedrowski, Xinyu Wang, and Alex Gorodetsky. Tensor network structure search using program synthesis. *arXiv* preprint arXiv:2502.02711, 2025.
 - Cécile Haberstich, Anthony Nouy, and Guillaume Perrin. Active learning of tree tensor networks using optimal least squares. *SIAM/ASA Journal on Uncertainty Quantification*, 11(3):848–876, 2023.
 - Wolfgang Hackbusch and Stefan Kühn. A new scheme for the tensor representation. *Journal of Fourier analysis and applications*, 15(5):706–722, 2009.
 - Meraj Hashemizadeh, Michelle Liu, Jacob Miller, and Guillaume Rabusseau. Adaptive learning of tensor network structures. *arXiv preprint arXiv:2008.05437*, 2020.
 - Kohei Hayashi, Taiki Yamaguchi, Yohei Sugawara, and Shin-ichi Maeda. Exploring Unexplored Tensor Network Decompositions for Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, pp. 5553–5563, 2019.
 - Christopher J Hillar and Lek-Heng Lim. Most tensor problems are NP-hard. *Journal of the ACM (JACM)*, 60(6):45, 2013.
 - Frank L Hitchcock. The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics*, 6(1-4):164–189, 1927.
 - Ming Hou, Jiajia Tang, Jianhai Zhang, Wanzeng Kong, and Qibin Zhao. Deep multimodal multilinear fusion with high-order polynomial pooling. In *Advances in Neural Information Processing Systems*, pp. 12113–12122, 2019.
 - Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
 - Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Lee. Llm-adapters: An Adapter Family for Parameter-Efficient Fine-Tuning of Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 5254–5276, 2023.
 - Ajay Jaiswal, Lu Yin, Zhenyu Zhang, Shiwei Liu, Jiawei Zhao, Yuandong Tian, and Zhangyang Wang. From GaLore to WeLore: How Low-Rank Weights Non-uniformly Emerge from Low-Rank Gradients. *arXiv preprint arXiv:2407.11239*, 2024.
 - Misha E Kilmer and Carla D Martin. Factorization strategies for third-order tensors. *Linear Algebra and its Applications*, 435(3):641–658, 2011.
 - Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - Maxim Kodryan, Dmitry Kropotov, and Dmitry Vetrov. Mars: Masked automatic ranks selection in tensor decompositions. In *International Conference on Artificial Intelligence and Statistics*, pp. 3718–3732. PMLR, 2023.
 - Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3): 455–500, 2009.
- Jean Kossaifi, Zachary C Lipton, Arinbjörn Kolbeinsson, Aran Khanna, Tommaso Furlanello, and
 Anima Anandkumar. Tensor regression networks. *Journal of Machine Learning Research*, 21:
 1–21, 2020.
 - Joseph M Landsberg, Yang Qi, and Ke Ye. On the geometry of tensor network states. *Quantum Information & Computation*, 12(3-4):346–354, 2012.

Tsung-Dao Lee. CP nonconservation and spontaneous symmetry breaking. *Physics Reports*, 9(2): 143–177, 1974.

- Binghua Li, Ziqing Chang, Tong Liang, Chao Li, Toshihisa Tanaka, Shigeki Aoki, Qibin Zhao, and Zhe Sun. Parameter-efficient fine-tuning of 3d ddpm for mri image generation using tensor networks. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 382–392. Springer, 2025.
- Chao Li and Zhun Sun. Evolutionary topology search for tensor network decomposition. In International Conference on Machine Learning, pp. 5947–5957. PMLR, 2020.
 - Chao Li, Junhua Zeng, Zerui Tao, and Qibin Zhao. Permutation search of tensor network structures via local sampling. In *International Conference on Machine Learning*, pp. 13106–13124. PMLR, 2022.
 - Chao Li, Junhua Zeng, Chunmei Li, Cesar F Caiafa, and Qibin Zhao. Alternating local enumeration (tnale): Solving tensor network structure search with fewer evaluations. In *International Conference on Machine Learning*, pp. 20384–20411. PMLR, 2023.
 - Dengchun Li, Yingzi Ma, Naizheng Wang, Zhiyuan Cheng, Lei Duan, Jie Zuo, Cal Yang, and Mingjie Tang. Mixlora: Enhancing large language models fine-tuning with lora based mixture of experts. *arXiv* preprint arXiv:2404.15159, 2024.
 - Nannan Li, Yu Pan, Yaran Chen, Zixiang Ding, Dongbin Zhao, and Zenglin Xu. Heuristic rank selection with progressively searching tensor ring network. *Complex & Intelligent Systems*, pp. 1–15, 2021a.
 - Sujie Li, Feng Pan, Pengfei Zhou, and Pan Zhang. Boltzmann machines as two-dimensional tensor networks. *Physical Review B*, 104(7):075154, 2021b.
 - Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv* preprint *arXiv*:2402.09353, 2024a.
 - Yipeng Liu, Yingcong Lu, Weiting Ou, Zhen Long, and Ce Zhu. Adaptively Topological Tensor Network for Multi-view Subspace Clustering. *IEEE Transactions on Knowledge and Data Engineering*, 2024b.
 - Ivan Melo. Higgs potential and fundamental physics. *European Journal of Physics*, 38(6):065404, 2017.
 - Eva Memmel, Clara Menzen, Jetze Schuurmans, Frederiek Wesel, and Kim Batselier. Position: Tensor networks are a valuable asset for green ai. In *International Conference on Machine Learning*, pp. 35340–35353. PMLR, 2024.
 - Oscar Mickelin and Sertac Karaman. On algorithms for and computing with the tensor ring decomposition. *Numerical Linear Algebra with Applications*, 27(3):e2289, 2020.
 - Mikko Möttönen, Juha J Vartiainen, Ville Bergholm, and Martti M Salomaa. Quantum circuits for general multiqubit gates. *Physical review letters*, 93(13):130502, 2004.
 - Alexander Novikov, Dmitrii Podoprikhin, Anton Osokin, and Dmitry P Vetrov. Tensorizing neural networks. In *Advances in Neural Information Processing Systems*, pp. 442–450, 2015.
 - Román Orús. Tensor networks for complex quantum systems. *Nature Reviews Physics*, 1(9):538–550, 2019.
 - Ivan V Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5): 2295–2317, 2011.
 - Piyush Rai, Yingjian Wang, Shengbo Guo, Gary Chen, David Dunson, and Lawrence Carin. Scalable bayesian low-rank decomposition of incomplete multiway tensors. In *International Conference on Machine Learning*, pp. 1800–1808. PMLR, 2014.

- Lorenz Richter, Leon Sallandt, and Nikolas Nüsken. Solving high-dimensional parabolic pdes using the tensor train format. In *International Conference on Machine Learning*, pp. 8998–9009. PMLR, 2021.
 - Manuel S Rudolph, Jing Chen, Jacob Miller, Atithi Acharya, and Alejandro Perdomo-Ortiz. Decomposition of matrix product states into shallow quantum circuits. *Quantum Science and Technology*, 9(1):015012, 2023.
 - Ulrich Schollwöck. The density-matrix renormalization group. *Reviews of modern physics*, 77(1): 259–315, 2005.
 - Edwin Stoudenmire and David J Schwab. Supervised learning with tensor networks. In *Advances in Neural Information Processing Systems*, pp. 4799–4807, 2016.
 - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
 - Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3): 279–311, 1966.
 - Andong Wang, Chao Li, Mingyuan Bai, Zhong Jin, Guoxu Zhou, and Qibin Zhao. Transformed low-rank parameterization can help robust generalization for tensor neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Zi Yang, Samridhi Choudhary, Xinfeng Xie, Cao Gao, Siegfried Kunzmann, and Zheng Zhang. CoMERA: Computing-and Memory-Efficient Training via Rank-Adaptive Tensor Optimization. *arXiv* preprint arXiv:2405.14377, 2024.
 - Ke Ye and Lek-Heng Lim. Tensor network ranks. arXiv preprint arXiv:1801.02662, 2019.
 - Tatsuya Yokota, Qibin Zhao, and Andrzej Cichocki. Smooth PARAFAC decomposition for tensor completion. *IEEE Transactions on Signal Processing*, 64(20):5423–5436, 2016.
 - Junhua Zeng, Chao Li, Zhun Sun, Qibin Zhao, and Guoxu Zhou. tnGPS: Discovering Unknown Tensor Network Structure Search Algorithms via Large Language Models (LLMs). In *Forty-first International Conference on Machine Learning*, 2024a.
 - Junhua Zeng, Guoxu Zhou, Yuning Qiu, Chao Li, and Qibin Zhao. Bayesian tensor network structure search and its application to tensor completion. *Neural Networks*, pp. 106290, 2024b.
 - Qibin Zhao, Liqing Zhang, and Andrzej Cichocki. Bayesian CP factorization of incomplete tensors with automatic rank determination. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1751–1763, 2015.
 - Qibin Zhao, Guoxu Zhou, Shengli Xie, Liqing Zhang, and Andrzej Cichocki. Tensor ring decomposition. *arXiv preprint arXiv:1606.05535*, 2016.
 - Yu-Bang Zheng, Xi-Le Zhao, Junhua Zeng, Chao Li, Qibin Zhao, Heng-Chao Li, and Ting-Zhu Huang. SVDinsTN: A Tensor Network Paradigm for Efficient Structure Search from Regularized Modeling Perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26254–26263, 2024.

A TECHNICAL APPENDICES WITH ADDITIONAL RESULTS FOR SYNTHETIC DATA

A.1 COMPUTATIONAL ANALYSIS OF THE PROPOSED ALGORITHM

As illustrated in Algo. 1, each iteration requires B evaluations in the best case and |V| evaluations in the worst case, where |V| is the total number of core tensors in the TNO system. This value typically scales linearly with the number of TNO in the system, *i.e.*, M. Thus, even in the worst case, the number of evaluations grows linearly with the size of the system. Numerical results presented in the experiments also demonstrate that Algo. 1 requires significantly fewer evaluations compared to solving joint TN-SS directly following equation 2 directly.

A.2 Details of synthetic data generation

Fig. 5 illustrates the four common TNO graphs G with Q=5. Table 3 and 4 summarize the structures (in VI matrix form) for generating the TNO systems used in our experiment. Specifically, each TNO system contains 4 TNOs, with each TNO generated by masking core tensors with locations 1 to 4 described in Table 3 and 4 from the given common graph. For example, the 2nd TNO for the stairs structure is generated by masking the 1st, 2nd, and 8th core tensors, which corresponds to removing the 1st, 2nd, and 8th columns from the VI matrix of the original stairs structure.

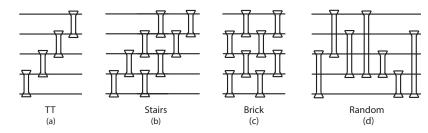


Figure 5: Illustration of common graphs used in our experiment.

Once the TNO system is constructed, we sample each element in core tensors from i.i.d Gaussian distribution with zero mean and unit variance, and the (observed) synthetic tensor data is obtained by contracting the core tensors together.

Table 3: TNO generation for **TT**, **Stairs**, and **Brick** structures with Q = 5, K = 2

Structure	Common graph (VI matrix)	Mask locations 1	Mask locations 2	Mask locations 3	Mask locations 4
TT	$\begin{pmatrix} 4 & 3 & 2 & 1 \\ 5 & 4 & 3 & 2 \end{pmatrix}$	1	2	3	4
Stairs	$ \begin{pmatrix} 4 & 3 & 2 & 1 & 4 & 3 & 2 & 1 \\ 5 & 4 & 3 & 2 & 5 & 4 & 3 & 2 \end{pmatrix} $	1, 2, 3	1, 2, 8	3, 4, 6	1, 4, 7
Brick	$ \begin{pmatrix} 4 & 2 & 3 & 1 & 4 & 2 & 3 & 1 \\ 5 & 3 & 4 & 2 & 5 & 3 & 4 & 2 \end{pmatrix} $	1, 2, 8	1, 6, 8	2, 4, 7	6, 7, 8

Table 4: TNO generation for **Rand.** structures with different Q(K=2)

$\overline{\mathbf{Q}}$	Common graph (VI matrix)		Mask locations 2		Mask locations 4
4	$\begin{pmatrix} 3 & 1 & 1 & 2 & 2 \\ 4 & 2 & 3 & 4 & 3 \end{pmatrix}$	1	2	3	4
5	$ \begin{pmatrix} 3 & 1 & 2 & 1 & 3 & 4 & 2 \\ 5 & 3 & 4 & 4 & 4 & 5 & 5 \end{pmatrix} $	1, 3	2, 3	1, 4	2, 4
6	$ \begin{pmatrix} 3 & 1 & 1 & 3 & 2 & 4 & 3 & 1 \\ 5 & 6 & 4 & 6 & 4 & 5 & 4 & 3 \end{pmatrix} $	1, 2, 3	2, 3, 6	1, 4, 5	4, 6, 8
8	$ \begin{pmatrix} 6 & 5 & 2 & 1 & 5 & 1 & 4 & 3 & 1 & 3 \\ 7 & 6 & 5 & 2 & 8 & 4 & 7 & 4 & 3 & 6 \end{pmatrix} $	1, 2, 3	2, 3, 6	1, 4, 5	4, 6, 8

A.3 TN-SS ALGORITHMS FOR PHASE I

The five joint TN-SS algorithms with default parameter settings used in Phase I are summarized in Algo. 2-5, which can be seen as extensions of the TNGA (Li & Sun, 2020), TNLS (Li et al., 2022), TnALE (Li et al., 2023) and Greedy (Hashemizadeh et al., 2020) methods, respectively. To align with the TNO models targeted in this work, we slightly modified these algorithms but tried our best to keep the searching dynamics unchanged.

In TNO, the adjacent core tensors can be merged if they have the same connectivity (i.e., the adjacent columns in Eq. (3) are the same). For the proposed TN-SS algorithms, once a new TNO graph is generated, we always check if the adjacent core tensors can be merged. If they are mergeable, one

808 809

Algorithm 2 Genetic Algorithm for joint TN-SS (TNGA) in Phase I

```
769
             1: Input:
770
             2: \mathcal{X}^{\overline{M}}
                                                                                                               771
             3: Q, K, I
                                                                                                                  > parameters for TNO
772
             4: L
                                                                                                       ⊳ order for the common graph
773
             5: S = 100
                                                                            > population size of individuals in each generation
774
             6: C_{max} = 20
                                                                                                > maximum number of generations
775
             7: f(\cdot) = \min_{\mathcal{X}^M \in tno(\cdot)} \pi_D(\mathcal{X}^M)
                                                                                                                            776
             8: \eta_1 = 5 \times 10^{-10}

    b tolerance

777
             9: \varepsilon = 0.2
                                                                                                                10: p(r) = \max\{0.01, \ln(200/(10^{-2} + 5r))\}
778
779
           11: Algorithm:
           12: randomly generate \{G_s\}_{s=1}^S \in \{G(V, E) : G \in \mathbb{G}, |V| = L\}
                                                                                                                  ▷ parent Initialization
780
           13: for t = 1 to C_{max} do
781
                      f^s = f(G_s) for all s \in [S]

    b fitness evaluation

782
                      \hat{G} = \arg\min_{G_s: s \in [S]} l^s
           15:
783
                      \hat{f} = f(\hat{G})
784
           16:
785
                      if f(\hat{G}) < \eta_1 then
           17:
                           break
786
           18:
           19:
                      end if
787
                      \begin{array}{ll} \{r_s\}_{s=1}^S = \operatorname{rank}(\{f^s\}_{s=1}^S) & \rhd \text{ get rank of the individuals by fitness evaluation } \\ \mathbb{G}_p = \{(G_s, p(r_s)) : s \in [S], r_s \leq \lceil (1-\varepsilon)S \rceil \} & \rhd \text{ eliminate the } \varepsilon \times 100\% \text{ individuals with } \end{cases}
           20:
789
                 worst fitness and compute the sampling probability p on remaining individuals
790
                      for s = 1 to S do
           22:
791
                                                                \triangleright select parents with probabilities p and with replacement
                           G_{p1}, G_{p2} \leftarrow \text{sample}(\mathbb{G}_p)
           23:
792
                           G_s \leftarrow \text{crossover\&mutate}(G_{p1}, G_{p2})
           24:

    be generate child from the parents

793
           25:
794
                      \{G_s\}_{s=1}^S \leftarrow \text{deduplicate\&fill}(\{G_s\}_{s=1}^S) \ \triangleright \text{remove individuals with duplicated graphs and}
                 generate new individuals to replace them
796
           27: end for
           28: Output: G, f
797
```

863

```
811
           Algorithm 3 Local Sampling algorithm for joint TN-SS (TNLS) in Phase I
812
            1: Input:
813
            2: \mathcal{X}^{\bar{M}}, Q, K, I, f(\cdot)
                                                                           814
            3: L, C_{max} = 40, \theta = 0.5, \eta_1 = 5 \times 10^{-10}
                                                                                  \triangleright algorithm parameters. \theta: sampling ratio
815
            4: Algorithm:
816
            5: initial \hat{G} \in \{G(V, E) : G \in \mathbb{G}, |V| = L\}
                                                                                                             ▷ Initialize the TNO
817
            6: \hat{f} = f(\hat{G})
                                                                                                    818
            7: for t = 1 to C_{max} do
819
                    \mathbb{G}_s = \emptyset
                                                                                                    ▶ Initialize the sampling set
820
                    for s=1 to L do
821
                         \mathbb{G}_s \leftarrow \mathbb{G}_s \cup \{G: G(V,E) \in \mathbb{G}, |V| = L, G \in \mathcal{N}(\hat{G},s)\} \triangleright \text{Add neighborhood TNOs to}
           10:
822
                the sampling set
823
                    end for
           11:
824
                    \mathbb{G}_s \leftarrow \operatorname{random\_sampling}(\mathbb{G}_s, \theta) \triangleright \operatorname{Randomly sample} \theta \times 100\% \operatorname{TNOs} from the sampling
           12:
825
                set
                    if \min_{G \in \mathbb{G}_s} f(G) < \tilde{f} then \triangleright Evaluation on sampled TNOs and update the best TNO if the
           13:
826
                condition satisfies
827
                         \hat{G} = \arg\min_{G \in \mathbb{G}_s} f(G)
           14:
828
                         \hat{f} = f(\hat{G})
                                                             ▶ Update the estimated TNO graph and corresponding loss
           15:
829
                    end if
           16:
830
                    if f < \eta_1 then
           17:
831
                         break
           18:
832
           19:
                    end if
833
           20: end for
834
           21: Output: \hat{G}, \hat{f}
835
```

Algorithm 4 Alternating Local Enumeration algorithm for joint TN-SS (TnALE) in Phase I

```
840
          1: Input:
          2: \mathcal{X}^{\bar{M}}, Q, K, I, f(\cdot)
841
                                                                  3: L, D = 2, \eta_1 = 5 \times 10^{-10}
                                                                  ▶ algorithm parameters. D: round-trips of ALE
842
          4: Algorithm:
843
          5: initial \hat{G} \in \{G(V, E) : G \in \mathbb{G}, |V| = L\}
                                                                                                ▷ Initialize the TNO
844
          6: \hat{f} = f(\hat{G})
                                                                                       845
          7: for t = 1 to D do
846
                  for s=1 to L do
                                                                                                      ⊳ Forward trip
847
                      \mathbb{G}_s = \{G : G(V, E) \in \mathbb{G}, |V| = L, G \in \mathcal{N}(\hat{G}, s)\} \triangleright Add neighborhood TNOs to the
848
              sampling set
849
                      if \min_{G \in \mathbb{G}_s} f(G) < \hat{f} then \triangleright Evaluation on sampled TNOs and update the best TNO if
         10:
850
              the condition satisfies
851
                          G = \arg\min_{G \in \mathbb{G}_s} f(G)
         11:
852
                          \hat{f} = f(\hat{G})
         12:
853
                      end if
         13:
854
         14:
                      if f < \eta_1 then
855
         15:
                          break
856
                      end if
         16:
857
         17:
                  end for
858
         18:
                  for s = L - 1 to 2 do
                                                                                                    859
                      Repeat steps 9-15
         19:
860
         20:
                  end for
         21: end for
861
         22: Output: \hat{G}, \hat{f}
862
```

Algorithm 5 Greedy algorithm for joint TN-SS in Phase I

```
1: Input:
 2: \mathcal{X}^{M}, Q, K, I, f(\cdot)
3: L, \eta_1 = 5 \times 10^{-10}

    ▶ algorithm parameters

 4: Algorithm:
 5: initial \hat{G} = \emptyset
                                                                                    6: for s = 1 to L do
         \mathbb{G}_s = \{G : G(V, E) \in \mathbb{G}, |V| = s, G \in \mathcal{N}^+(\hat{G})\}

    Add neighborhood TNOs to the

    sampling set
         \hat{G} = \arg\min_{G \in \mathbb{G}_s} f(G)
                                                                    ▶ Find the best TNO with minimum loss
         \hat{f} = f(\hat{G})
 9:

    □ Update the loss

         if \hat{f} < \eta_1 then
             break
11:
12:
         end if
13: end for
14: Output: \hat{G}, \hat{f}
```

of the duplicated core tensors will be replaced with a new one with different connectivity. This procedure is repeated until all core tensors in the new TNO graph can not be merged.

For TNLS and TnALE, the neighborhood of a TNO G is denoted as $\mathcal{N}(G,s)$, and is defined as the set of graphs that has the same connectivity of all core tensors except for the s-th one. In the view of the VI matrix, $\mathcal{N}(G,s)$ is generated by replacing the s-th column of the VI matrix of G with columns that represent all other possible connectivity.

Further, in the Greedy algorithm, the TNO is gradually generated with the graph order from 0 to L. At each iteration, starting from the previous TNO graph G, we add a core tensor to the right of the previous TNO, which is equivalent to adding a new column to the right of the VI matrix of G. The TNO set with all possible connectivity is denoted as $\mathcal{N}^+(G)$ in Algo. 5, and the best TNO for the next iteration is chosen as the one with the lowest RSE.

A.4 TENSOR DECOMPOSITION FOR TNO

Given the TNO graph, we apply Adam optimizer (Kingma & Ba, 2014) to perform tensor decomposition, which is commonly used in existing TN-SS methods. Specifically, given a TNO graph and the observed tensor \mathcal{X} , we initialize each core tensor from Gaussian distribution with zero mean and variance 0.3. The contraction expression of core tensors is constructed using einsum function similar to that in (Chen et al., 2024a). The relative squared error (RSE) is used as the loss function in Eq. (2). Specifically, given the generated synthetic tensors $D := \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_M\}$ and the approximated tensors $\mathcal{X}^M := \{\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_M\}$, the RSE is defined as $\pi_D(\mathcal{X}^M) = \sum_{m=1}^M \|\mathcal{X}_m - \mathcal{D}_m\|_F^2 / \sum_{m=1}^M \|\mathcal{D}_m\|_F^2$, where $\|\cdot\|_F$ is the Frobenius norm. The learning rate of Adam is set to 10^{-2} and the maximum iteration number is set to 1500.

For each evaluation, the RSE is measured over 10 runs with different initializations of core tensors. Specifically, we select the minimum squared approximation error for each tensor over 10 runs and compute the overall RSE of four tensors. In Phase I, the tolerance η_1 for Algo. 2-5 is set to 5×10^{-10} to determine if the common graph G is searched, and the initialization strategies are the same as existing TN-SS methods. The setting of common graph order L is determined differently for each data, following a basic principle that L should be relatively large to achieve a satisfactory approximation error in Phase I. While in Phase II, the RSE is compared with tolerance $\eta_2 = 10^{-6}$ to determine if the core tensor should be masked.

A.5 IMPLEMENTATION

In our experiments on synthetic data, we run all algorithms on a cluster of NVIDIA V100 GPUs alongside an Intel Xeon E5-2690 CPU node. Specifically, the CPU node handles data input, applies the TNO generation procedures, and distributes the sampled TNOs across the GPUs. Each GPU then

Table 5: Results of the TD experiment. The RSE, total number of core tensors (shown in round brackets), and the corresponding number of evaluations (shown in square brackets) are presented for $Phase\ I\ (top\ block)$ and $Phase\ II\ (middle\ block)$. The bottom block shows the RSE of TN-SS methods for the same graph orders as $Phase\ II$.

Method	TT	Stairs	Brick	Rand.
GREEDY	$< 5 \times 10^{-10} (24) [55]$	$2.5 \times 10^{-7} (40) [91]$	$7.8 \times 10^{-6} (40) [91]$	0.0426 (40) [91]
TnALE	$< 5 \times 10^{-10} (24) [109]$	$< 5 \times 10^{-10} (40) [250]$	$< 5 \times 10^{-10} (40) [94]$	2.12×10^{-5} (40) [265]
TNLS	$< 5 \times 10^{-10} (24) [89]$	$< 5 \times 10^{-10} (40) [217]$	$< 5 \times 10^{-10} (40) [613]$	2.25×10^{-6} (40) [1441]
TNGA	$< 5 \times 10^{-10} (24) [300]$	$4 \times 10^{-8} (40) [2000]$	$< 5 \times 10^{-10} (40) [800]$	$< 5 \times 10^{-10} (40) [900]$
GREEDY-SB	$< 5 \times 10^{-10} (12) [49]$	$< 5 \times 10^{-10} (20) [140]$	$< 5 \times 10^{-10} (24) [142]$	0.0426 (40) [40]
TnALE-SB	$< 5 \times 10^{-10} (12) [61]$	$< 5 \times 10^{-10} (32) [82]$	1×10^{-8} (32) [103]	2.12×10^{-5} (40) [40]
TNLS-SB	$< 5 \times 10^{-10} (20) [50]$	$< 5 \times 10^{-10} (29) [136]$	$< 5 \times 10^{-10}$ (22) [136]	$< 5 \times 10^{-10} (38) [66]$
TNGA-SB	$< 5 \times 10^{-10} (12) [48]$	$< 5 \times 10^{-10} (22) [147]$	1×10^{-8} (26) [110]	$< 5 \times 10^{-10} (20) [112]$
	RSE of only optim	izing G with the $\left V\right $ aligne	d to the results in Phase	II
GREEDY	0.3568 (12) [28]	0.2032 (20) [46]	0.1838 (24) [55]	0.0426 (40) [91]
TnALE	0.3637 (12) [66]	0.0094 (32) [208]	0.0133 (32) [208]	2.12×10^{-5} (40) [265]
TNLS	$< 5 \times 10^{-10} (20) [111]$	0.0045 (32) [1161]	0.1118 (24) [881]	2.25×10^{-6} (40) [1441]
TNGA	0.2907 (12) [2000]	0.0145 (24) [2000]	0.0243 (28) [2000]	0.1050 (20) [2000]

Table 6: Results on different Q (random structure with I=2)

Method	Q = 4	Q = 5	Q = 6	Q = 7	Q = 8
TNGA (S=100)	$<5 \times 10^{-10} (40) [300]$	$<5 \times 10^{-10} (50) [600]$	$< 7.4 \times 10^{-5} (50) [600]$	0.0048 (60) [1500]	0.0312 (60) [600]
TNGA-SB	$<5 \times 10^{-10} (16) [62]$	$<5 \times 10^{-10} (39) [40]$	$< 7.4 \times 10^{-5} (50) [50]$	0.0048 (60) [60]	0.0312 (60) [60]
TNGA (S=200)	-	-	$<5 \times 10^{-10} (50) [1200]$	$<5 \times 10^{-10} (60) [4000]$	0.0237 (60) [3400]
TNGA-SB	-	-	$<5 \times 10^{-10} (25) [178]$	$<5 \times 10^{-10} (32) [152]$	0237 (60) [60]
TNGA (S=300)	_	_	_	-	$<5 \times 10^{-10} (60) [2100]$
TNGA-SB	_	_	_		$<5 \times 10^{-10} (34) [143]$

performs the tensor decomposition for a given TNO and returns its loss. After each iteration, the CPU node collects these loss values and generates new TNOs according to the specific algorithm for the subsequent procedure.

A.6 ADDITIONAL EXPERIMENTAL RESULTS ON SYNTHETIC DATA

As a detailed version of Figure 2 (a)-(c), Table 5 reports all the RSE, the total number of core tensors, and the total number of evaluations for synthetic data. The algorithms with suffix '-SB' denote that Algo.1 is applied for Phase II. We should remark that, although Greedy and TNLS fail to meet the $RSE \leq 10^{-6}$ for **Stairs** and **Rand.** in Phase II, respectively, their RSE values are close to 10^{-6} such that the Phase II can still work to reach the RSE lower than 10^{-6} .

We conduct additional experiments to verify the performance on different TNO parameters Q and I. The **Rand.** structures with Q = [4,5,6,7,8] (see Table.4) are applied as common graphs for data generation, and the TNGA method is applied for Phase I. Table. 6 shows the performance under different Q. We gradually increase the population size parameter S from 100 until the RSE reaches $\leq 10^{-6}$. The results verify the effectiveness of the proposed method for variant selections of Q. Further, the S needs to be set to a larger value as Q increases due to the increasing searching space of the common graph. Consequently, the total number of evaluations in both Phase I and Phase II grows with Q increases. Further, the decomposition performance under TNO with different I is depicted in Fig. 7. It can be seen that the proposed algorithm works well on TNO with varying selections of I.

We also investigate the performance of different common graph orders L used in Phase I. Fig. 6 depicts the performance under different L for the proposed method. The results imply that L works well when L is slightly larger than groundtruth $L^*=7$. It can be also seen that if the L is underestimated (i.e. L<7), the approximation always fails as they can not find the common graph with $L^*=7$. When L=7, the approximation still fails as the TNGA method may fail to find the common graph with a limited number of evaluations. On the other hand, if the L is too large (i.e., $L\geq 14$ in this experiment), the approximation also fails due to the optimization problem. In practice,

Table 7: RSE, total number of core tensors (round brackets), and the corresponding number of evaluations (square brackets) on different I (**Rand.** structure with Q = 5).

Method	I=2	I=4	I=6
TNGA	$< 5 \times 10^{-10} (32) [100]$	$< 5 \times 10^{-10} (32) [1000]$	$< 5 \times 10^{-10} (32) [800]$
TNGA-SB	$< 5 \times 10^{-10} (16) [70]$	$< 5 \times 10^{-10} (31) [32]$	$6.4 \times 10^{-7} (22) [69]$

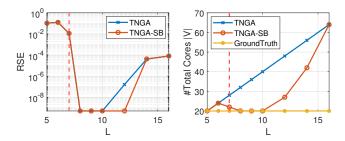


Figure 6: RSE (left) and the total number of core tensors (right) under **Rand.** structure (Q = 5, K = 2, I = 2) using different parameter L. Dotted red line: order of groundtruth **Rand.** structure.

as the L^* is always not known, we may conduct TNGA on various selections of L and choose a proper one that achieves a satisfactory approximation error.

Further, the efficiency of the proposed joint TN-SS method under different M is verified. TNGA is employed as the TN-SS algorithm for Phase I. The maximum number of evaluations C_{max} is $1000 \times M$. Table 8 presents the results of tensor decomposition for various values of M. We add the vanilla TN-SS as the comparison, in which each TNO is optimized independently without searching a common graph. As shown, for M=4, the vanilla TN-SS requires significantly more evaluations than the proposed method to search TNOs that can successfully approximate the tensors. As M increases to 6 and 8, the vanilla TN-SS reaches the maximum evaluation C_{max} but still fails to achieve the desired approximation accuracy of $RSE \leq 10^{-6}$. In contrast, although increasing slightly raises computational complexity for the proposed method, it substantially enhances search efficiency and final performance due to the presence of shared structural patterns across TNOs. This confirms that symmetry breaking enables more effective exploration of the TNO structure space—especially when the TNOs are related (e.g., originating from similar data or architectural settings).

Finally, the sensitivity of different selections of the tolerance parameter η_2 in Phase II is analyzed. The experiment is conducted on **Rand.** structure with Q=5, and TNGA is applied for Phase I. Table 9 reports the performance in terms of different η_2 . These results show that the method is robust in a wide range of η_2 as long as the task has an internal tensor network structure. As expected, larger values give a more compact structure, and smaller values improve accuracy with more structure preserved.

Table 8: RSE and the corresponding number of evaluations (square brackets) on different M (Rand. structure with Q=5).

Method	M = 2	M = 4	M = 6	M = 8
Vanilla TN-SS Ours		$< 5 \times 10^{-10} [3800]$ $< 5 \times 10^{-10} [823]$	$2.9 \times 10^{-6} [6000]$ < $5 \times 10^{-10} [2027]$	$ \begin{array}{r} 1.5 \times 10^{-4} \\ < 5 \times 10^{-10} \text{ [2226]} \end{array} $

Table 9: RSE, total number of core tensors (round brackets), and the corresponding number of evaluations (square brackets) on different η_2 (**Rand.** structure with Q=5).

	$\eta_2 = 1$	$\eta_2 = 0.01$	$\eta_2 = 10^{-4}$	$\eta_2 = 10^{-6}$	$\eta_2 = 10^{-8}$
TNGA-SB	0.9926 (4) [41]	$< 5 \times 10^{-10} (20) [83]$	$< 5 \times 10^{-10} (20) [102]$	$< 5 \times 10^{-10} (20) [112]$	$< 5 \times 10^{-10} (20) [87]$

B TECHNICAL APPENDICES WITH ADDITIONAL RESULTS FOR LLMS FINE-TUNING

B.1 COMMONSENSE REASONING DATASETS

Table 10 presents detailed information about the datasets used in our experiments. Specifically, (Hu et al., 2023) originally collected all datasets including the training set **Train** and the test set **Test**. In our experiments, we randomly select 3000 and 400 data from the original training set as the training data **Train split** and validation data **Valid split**, respectively. We should note that, for ARC-Easy and ARC-Challenge datasets, the training data for our experiments are 1825 and 720 due to insufficient original training data. The fine-tuning performance is evaluated on the test set **Test**.

Table 10: Description of common sense reasoning datasets used in experiments.

Dataset name	Domain	# Train	# Train split	# Valid split	# Test
PIQA	Physical Interaction	16113	3000	400	1838
SIQA	Social Interaction	33410	3000	400	1954
OBQA	Science Facts	4957	3000	400	500
ARC-Easy (ARC-e)	Natural Science	2251	1851	400	2376
ARC-Challenge (ARC-c)	Natural Science	1119	719	400	1172

B.2 PEFT FOR TNO

For each dataset, the fine-tuning is performed on **Train split** described in Table 10. The QuanTA model (Chen et al., 2024a) is adapted to our experiments with modifications that can handle different TNOs for each transformer layer. The training parameters used in experiments are shown in Table 11. The baseline algorithms QuanTA-6/4/2 that use the same TNO for all layers are shown in Table 12. Specifically, QuanTA-6 corresponds to the original QuanTA algorithm proposed in (Chen et al., 2024a).

In our experiment, we directly use the TNO of QuanTA-6 as the common graph G and deploy the proposed algorithm for Phase II. For each evaluation, given a TNO system consisting of the TNOs for all layers, a training procedure with 3 epochs is conducted on **Train split**, following the inference on both **Valid split** and **Test**. Instead of using training loss as the loss function for $\pi_D(\cdot)$ in Algo. 1, we directly set the loss $f(\cdot)$ as $1-acc_v$ where acc_v is the accuracy on **Valid split**. The tolerance η_2 for each dataset is set according to the validation accuracy on TNO using QuanTA-6. Specifically, the η_2 for PIQA, SIQA, OBQA, ARC-e and ARC-c are set to 0.20, 0.23, 0.21, 0.15, 0.38, respectively.

Table 11: Hyperparameter Settings for QuanTA fine-tuning.

Hyperparameters	Values		
Num of Epochs	3		
Batch Size	4		
Optimizer	AdamW		
Scheduler	Linear Scheduler		
Learning Rate	1e-4		
Weight Decay	0		
Dropout	0		
Modules	(q_proj v_proj)		

1	80	0
1	80	1

Table 12: Different TNO structures of QuanTA used in experiments

Method	TNO for all layers (in VI matrix form)
QuanTA-6	$ \begin{pmatrix} 3 & 2 & 1 & 2 & 1 & 1 \\ 4 & 4 & 4 & 3 & 3 & 2 \end{pmatrix} $
QuanTA-4	$\begin{pmatrix} 3 & 2 & 1 & 1 \\ 4 & 3 & 2 & 4 \end{pmatrix}$
QuanTA-2	$\begin{pmatrix} 3 & 1 \\ 4 & 2 \end{pmatrix}$

B.3 IMPLEMENTATION

In our experiments on PEFT for LLMs, training is carried out on two NVIDIA RTX A6000 GPUs (each with 48 GB of memory). Similar to synthetic data, an Intel Xeon E5-2690 CPU node manages results collection and TNO generation.

B.4 Additional experimental results on LLM fine-tuning

Fig. 7 depicts the dynamic information during the Phase II on five datasets. (a)-(b) report the average accuracy degradation ($\times 100\%$) on **Test** and **Valid split**, respectively. The curves with light color in the background denote the accuracy at each iteration, and the curves with dark color foreground show the average accuracy smoothed with window size 20. (c)-(d) show the total number of core tensors and parameters at each iteration. It should be noted that the accuracy degradation on **Valid split** is bounded due to the settings of tolerance parameter η_2 . As can be seen, for all datasets except SIQA, the total number of core tensors and parameters reduce quickly at the first 150 iterations and then slow down due to the setting of the tolerance η_2 .

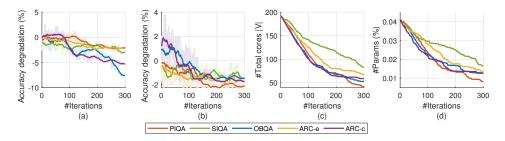


Figure 7: Accuracy degradation on test data and validation data over iterations (a,b). Total number of core tensors and #Params over iterations (c, d);

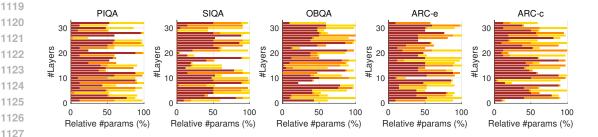


Figure 8: #Params percentage relative to QuanTA-6 per layer of Llama2-7B on five datasets. Yellow, orange, and brown bars indicate #Params of 0.031%, 0.024%, and 0.017%, respectively.

Fig. 8 shows the number of parameters on each transformer layer compared with QuanTA-6. The yellow, orange, and brown bars correspond to the three settings with the total number of parameters 0.031%, 0.024%, and 0.017%, respectively. Further, the VI matrices of corresponding TNOs for

8/16/24/32 layers with the total number of parameters 0.024% are shown in Table 13. The results imply the asymmetric property of TNO across different layers and different datasets.

Table 13: TNO structures (in VI matrix form) of different transformer layers with #Params 0.024%

Layer	PIQA	SIQA	OBQA	ARC-e	ARC-c
8	$ \begin{pmatrix} 3 & 2 & 1 & 1 & 1 \\ 4 & 4 & 4 & 3 & 2 \end{pmatrix} $	$ \begin{pmatrix} 3 & 2 & 1 & 1 & 1 \\ 4 & 4 & 4 & 3 & 2 \end{pmatrix} $	$\begin{pmatrix} 3 & 1 & 2 \\ 4 & 4 & 3 \end{pmatrix}$	$\begin{pmatrix} 1 & 1 \\ 4 & 2 \end{pmatrix}$	$ \begin{pmatrix} 3 & 2 & 2 \\ 4 & 4 & 3 \end{pmatrix} $
16	$ \begin{pmatrix} 3 & 2 & 1 & 2 & 1 \\ 4 & 4 & 4 & 3 & 3 \end{pmatrix} $	$\begin{pmatrix} 1 & 2 \\ 4 & 3 \end{pmatrix}$	$\begin{pmatrix} 1 & 1 & 1 \\ 4 & 3 & 2 \end{pmatrix}$	$\begin{pmatrix} 2 & 1 & 1 \\ 4 & 4 & 3 \end{pmatrix}$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$
24	$\begin{pmatrix} 1 & 2 & 1 & 1 \\ 4 & 3 & 3 & 2 \end{pmatrix}$	$\begin{pmatrix} 2 & 1 & 1 \\ 3 & 3 & 2 \end{pmatrix}$	$ \begin{pmatrix} 1 & 2 & 1 & 1 \\ 4 & 3 & 3 & 2 \end{pmatrix} $	$\begin{pmatrix} 1 & 2 & 1 \\ 4 & 3 & 2 \end{pmatrix}$	$\begin{pmatrix} 3 & 2 & 1 \\ 4 & 3 & 2 \end{pmatrix}$
32	$ \begin{pmatrix} 3 & 2 & 1 & 1 \\ 4 & 4 & 4 & 3 \end{pmatrix} $	$ \begin{pmatrix} 3 & 2 & 1 & 1 \\ 4 & 3 & 3 & 2 \end{pmatrix} $	$\begin{pmatrix} 2 \\ 3 \end{pmatrix}$	$\begin{pmatrix} 1 & 2 & 1 \\ 4 & 3 & 3 \end{pmatrix}$	$ \begin{pmatrix} 2 & 1 & 2 \\ 4 & 4 & 3 \end{pmatrix} $

C TECHNICAL APPENDICES WITH ADDITIONAL RESULTS FOR REPRESENTATION OF QUANTUM CIRCUIT OPTIMIZATION

C.1 PROBLEM SETTINGS AND SOLUTIONS

The goal of quantum circuit optimization is to approximate a quantum operator U by a quantum circuit composed of several quantum gates. Following the common setting in quantum computing, we assume that each two-qubit gate operates on adjacent qubits in a 1D chain. As shown in Fig. 9, a 4 qubits quantum operator U (a) can be represented as a quantum circuit with five two-qubit gates (b). In practice, finding a quantum circuit like (b) is challenging. An efficient way that avoid circuit structure searching is directly using the quantum circuits with a brick-wall structure with 3 blocks, as shown in Fig. 9 (c). However, the quantum circuits constructed using brick-wall structures are not always compact due to the redundancy of quantum gates (e.g., the gates in the red-dotted boxes in Fig. 9 (c)). In this work, we aim to efficiently search for a more compact representation with fewer quantum gates compared to the brick-wall structure.

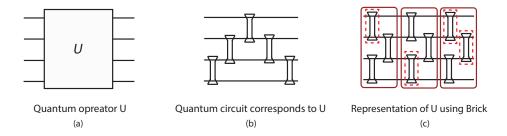


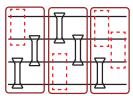
Figure 9: Illustration of the representation of a quantum operator.

C.2 GENERATION OF SYNTHETIC QUANTUM OPERATORS

In our experiment, we consider synthetic quantum operators with dimension I=2 and qubit counts Q=4,8,12. For each qubit count, a quantum circuit using a brick-wall structure with M blocks is constructed. Then, a quantum circuit is generated by randomly masking a number of quantum gates, as shown in Fig.10 and Table.14. Finally, for each quantum circuit, 20 quantum operators are obtained by sampling quantum gates from the Haar-Gaussian distribution and contracting them according to the generated quantum circuit structure.

C.3 GENERATION OF QUANTUM FOURIER TRANSFORM (QFT) OPERATORS

The QFT is to represent the discrete Fourier transform (DFT) using a quantum circuit. As analyzed in (Chen et al., 2023; Camps et al., 2021), the DFT matrix can be efficiently represented by QFT using



quantum circuit (4 qubits)

Figure 10: Example of the quantum circuits generated for the synthetic quantum operator experiment (Q=4).

Table 14: TNO structure for generation of synthetic quantum operators.

#qubits	#block	Mask locations for each block						
Q	M	Block 1	Block 2	Block 3	Block 4	Block 5		
4	3	1	2	1, 3	/			
8	5	3, 4, 6	2, 3, 6, 7	1, 2, 4, 5, 7	1, 3, 5, 6	2, 4, 5, 6		
12	3	2, 3, 4, 6, 8, 10	1, 3, 5, 6, 7, 9, 11	1, 2, 4, 5, 7, 8, 9, 10, 11	/	/		

matrix (tensor) decomposition due to its special structure. In our experiment, we consider the DFT matrix F of size $2^Q \times 2^Q$ with Q = 4, 6, and the QFT operator is set as U = F.

C.4 ALGORITHM FOR QUANTUM CIRCUIT OPTIMIZATION

From the view of TNO, the quantum circuit can be seen as a TNO graph, and the quantum gates can be represented using the core tensors. Thus, representing a quantum operator U is equivalent to approximating it using a TNO graph with some core tensors. Additionally, according to the property of quantum computing, we should constrain the core tensors to be unitary and relax the value to be complex-valued.

Concentrating on the goal of compact representation, we assume the M for each generated synthetic quantum operator U is known and directly deploy Phase II to the brick-wall TNO structure with M blocks. While for each QFT operator, M is set to a relatively large value to achieve a satisfactory high fidelity. At each evaluation, given a TNO and the quantum operator U, the density-matrix renormalization group (DMRG) method (Schollwöck, 2005) is applied to optimize the core tensors, with each core tensor initialized as an identity tensor. The loss function $\pi_D(\cdot)$ in Eq. (2) is defined as

$$\pi_D(X) = 1 - |tr\left(U^{\dagger}X\right)|/I^Q. \tag{6}$$

where $|tr(U^{\dagger}X)|/I^Q$ is the normalized fidelity quantifying the 'closeness' between two operators U and X. The η_2 is set to 10^{-3} . For each case of representing synthetic quantum operators, the fidelity is averaged over the generated 20 unitary operators.

C.5 IMPLEMENTATION

Due to limited resources, we implement the quantum experiment just in the proof-of-concept scale using a laptop MacBook Air 13-inch with M3 chip and memory of 24 GB.

D Proofs for Proposition 3.1

In this section, we present the proof of Proposition 3.1 from the main manuscript. Recall the proposition as follows.

Proposition D.1 (Perturbation analysis). Let $\mathcal{Z}^M = \{Z_1, \dots, Z_M\}$ be a system of TNOs with $\mathcal{Z}^M \in \arg\min_{\mathcal{Y}^M \in ts(G_1, \dots, G_M; Q, K, I)} \pi_D(\mathcal{Y}^M)$, where π_D is a task loss defined jointly on multiple

operators. Suppose one operator Z_j contains a maskable core \mathcal{G} with $\|\mathcal{G}\|_F = 1$ and operator–Schmidt rank S. Let $\mathcal{Z}^{M,\star} = \{Z_1, \ldots, Y_j^{\star}, \ldots, Z_M\}$ be the system obtained by replacing Z_j with its best masked version

$$Y_j^{\star} \in \arg\min_{Y \in tno(G_{\text{mask}}; Q, K, I)} \pi_D(\{Z_1, \dots, Y, \dots, Z_M\}),$$

where G_{mask} is obtained by masking \mathcal{G} . If π_D is L-Lipschitz in $\|\cdot\|_F$ with respect to each operator, then

$$|\pi_D(\mathcal{Z}^{M,\star}) - \pi_D(\mathcal{Z}^M)| \leq L C \sqrt{S-1},$$

where C depends only on the contracted environment of G. In particular, if S=1, masking G does not increase the loss.

Proof. Let $\mathcal{Z}^M = \{Z_1, \dots, Z_M\} \in ts(G_1, \dots, G_M; Q, K, I)$ be system-optimal:

$$\mathcal{Z}^M \in \arg\min_{\mathcal{Y}^M \in ts(G_1,...,G_M)} \pi_D(\mathcal{Y}^M).$$

Fix an index j and a maskable core $\mathcal G$ inside Z_j with $\|\mathcal G\|_F=1$ and operator–Schmidt decomposition $\mathcal G=\sum_{r=1}^S \lambda_r\,\mathcal U_r\otimes \mathcal V_r$, where $\lambda_1\geq \cdots \geq \lambda_S\geq 0$ and $\|\mathcal U_r\|_F=\|\mathcal V_r\|_F=1$. By the bipartition assumption, contracting the entire network except $\mathcal G$ defines a linear map

$$\mathcal{T}: \mathbb{F}^{I^K \times I^K} \longrightarrow \mathbb{F}^{I^Q \times I^Q}, \quad \text{such that } Z_i = \mathcal{T}(\mathcal{G}).$$

Since \mathcal{T} is linear, write

$$Z_j = \mathcal{T}(\mathcal{G}) = \sum_{r=1}^S \lambda_r \, \mathcal{T}(\mathcal{U}_r \otimes \mathcal{V}_r).$$

A realizable masked candidate. Masking $\mathcal G$ removes that core and reconnects its neighbors. Because the top Schmidt term $\mathcal U_1\otimes\mathcal V_1$ factorizes across the bipartition, it can be absorbed into the adjacent cores on each side of the cut, hence is realizable in the masked class. Therefore, there exists $Y_j^{(0)}\in tno(G_{\mathrm{mask}};Q,K,I)$ such that

$$Y_i^{(0)} = \mathcal{T}(\mathcal{U}_1 \otimes \mathcal{V}_1).$$

Define the system candidate $\mathcal{Z}^{M,(0)} := \{Z_1,\ldots,Y_j^{(0)},\ldots,Z_M\}.$

Environment norm bound. Let $C := \|\mathcal{T}\|_{2 \to F} := \sup_{X \neq 0} \frac{\|\mathcal{T}(X)\|_F}{\|X\|_F}$, which depends only on the (fixed) contracted environment surrounding \mathcal{G} . Then

$$||Z_j - Y_j^{(0)}||_F = \left\| \sum_{r=2}^S \lambda_r \, \mathcal{T}(\mathcal{U}_r \otimes \mathcal{V}_r) \right\|_F \le \sum_{r=2}^S \lambda_r \, ||\mathcal{T}(\mathcal{U}_r \otimes \mathcal{V}_r)||_F \le C \sum_{r=2}^S \lambda_r.$$

Since $\|\mathcal{G}\|_F^2 = \sum_{r=1}^S \lambda_r^2 = 1$, Cauchy–Schwarz yields $\sum_{r=2}^S \lambda_r \leq \sqrt{S-1} (\sum_{r=2}^S \lambda_r^2)^{1/2} \leq \sqrt{S-1}$, hence

$$||Z_j - Y_j^{(0)}||_F \le C\sqrt{S-1}.$$

Lipschitz transfer to task loss. By the *L*-Lipschitz property of π_D w.r.t. each operator in Frobenius norm,

$$|\pi_D(\mathcal{Z}^{M,(0)}) - \pi_D(\mathcal{Z}^M)| \le L \|Z_j - Y_j^{(0)}\|_F \le L C \sqrt{S - 1}.$$

Optimal masked replacement. By definition,

$$Y_j^{\star} \in \arg\min_{Y \in tno(G_{\text{mask}}; Q, K, I)} \pi_D(\{Z_1, \dots, Y, \dots, Z_M\}),$$

so $\pi_D(\mathcal{Z}^{M,\star}) \leq \pi_D(\mathcal{Z}^{M,(0)})$, where $\mathcal{Z}^{M,\star} := \{Z_1, \dots, Y_i^{\star}, \dots, Z_M\}$. Therefore,

$$|\pi_D(\mathcal{Z}^{M,\star}) - \pi_D(\mathcal{Z}^M)| \le |\pi_D(\mathcal{Z}^{M,(0)}) - \pi_D(\mathcal{Z}^M)| \le LC\sqrt{S-1}.$$

Finally, if S=1 then $Z_j=\mathcal{T}(\mathcal{U}_1\otimes\mathcal{V}_1)$ is itself realizable post-masking, so we may choose $Y_j^{(0)}=Z_j$, giving zero loss change. This proves the claim.

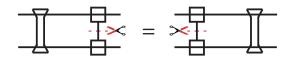


Figure 11: A simple example demonstrating the equivalence of masking different core tensors in a QTN.

We emphasize that a small Schmidt rank, as discussed in Proposition 3.1, is *not* a necessary condition for masking a core tensor to have minimal impact on approximation error. This is demonstrated by a simple yet insightful example in Figure 11, where two core tensors occupy a commutative position within a TNO, and each individually possesses full model representation. In this case, masking either core tensor yields an equivalent result, regardless of its Schmidt number. More generally, this condition holds when the Zariski closure (Landsberg et al., 2012) of the remaining TNO spans the entire ambient space $\mathbb{F}^{I^K \times I^K}$. Such universal representation properties are well-studied in the context of quantum computing Barenco et al. (1995); Möttönen et al. (2004).

E STATEMENT OF THE USE OF LARGE LANGUAGE MODELS (LLMS)

During the writing of this paper, language polishing and grammar checking were partially assisted by Large Language Models (LLMs). The LLMs were used to improve the accuracy and fluency of the text, with all modifications reviewed and approved by the author to ensure originality and academic integrity.