

# EFFICIENT CAUSAL STRUCTURE LEARNING VIA MODULAR SUBGRAPH INTEGRATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Learning causal structures from observational data remains a fundamental yet computationally intensive task, particularly in high-dimensional settings where existing methods face challenges such as the super-exponential growth of the search space and increasing computational demands. To address this, we introduce VISTA (Voting-based Integration of Subgraph Topologies for Acyclicity), a modular framework that decomposes the global causal structure learning problem into local subgraphs based on Markov Blankets. The global integration is achieved through a weighted voting mechanism that penalizes low-support edges via exponential decay, filters unreliable ones with an adaptive threshold, and ensures acyclicity using a Feedback Arc Set (FAS) algorithm. The framework is model-agnostic, imposing no assumptions on the inductive biases of base learners, is compatible with arbitrary data settings without requiring specific structural forms, and fully supports parallelization. We also theoretically establish finite-sample error bounds for VISTA, and prove its asymptotic consistency under mild conditions. Extensive experiments on both synthetic and real datasets consistently demonstrate the effectiveness of VISTA, yielding notable improvements in both accuracy and efficiency over a wide range of base learners.

## 1 INTRODUCTION

Understanding causal relationships from observational data Pearl (2009) is critical across numerous fields such as biology Petersen et al. (2024), economics Hünernund & Bareinboim (2023), and healthcare Sanchez et al. (2022b). Identifying causal structures enables reliable interventions and scientific insights. A common modeling framework represents the system as a causal graph—a Directed Acyclic Graph (DAG) where nodes are variables and directed edges denote causal links Spirtes et al. (2000). While identifiability of the true DAG generally requires additional structural assumptions, our VISTA framework inherits whatever identifiability guarantees each base learner provides. In practice, large-scale observational datasets further complicate structure recovery, as most existing algorithms struggle to scale efficiently. Constraint-based pipelines Spirtes et al. (2000); Meek (2013) must search over large conditioning sets while the number of CI tests grows combinatorially with the size of graph, and finite-sample CI tests become unreliable in high dimensions, so early mistakes can easily propagate to later steps. Score-based learners Chickering (2002); Loh & Bühlmann (2014) optimize over a super-exponential DAG space; practical solvers still require heavy global searches or acyclicity constraints with repeated dense updates, driving time and memory up sharply. These disadvantages make them difficult to perform well in large-scale datasets.

Given the challenges of learning large-scale causal structures, divide-and-conquer strategies have emerged as a natural solution. By decomposing the global graph into smaller, tractable subgraphs, these methods significantly reduce computational complexity, particularly in sparse settings, and facilitate parallel or distributed computation. In addition, aggregating local structures often enhances robustness relative to learning the full graph in a single pass. Early approaches expand neighborhoods from a random node Gao et al. (2017) or apply hierarchical clustering Gu & Zhou (2020). More recent work often partition the variable set into local neighborhoods, such as Markov Blankets, before aggregating them Dong et al. (2024); Mokhtarian et al. (2021); Tsamardinos et al. (2003); Wu et al. (2023; 2022). However, the majority of these “conquer” steps rely on fixed heuristics for merging, such as voting thresholds, edge overlap rules, or manual conflict resolution. While simple, such rule-based schemes lack adaptability to noise and offer limited theoretical guarantees for global

consistency. DCILP Dong et al. (2024) formulates the merging process as an Integer Linear Program (ILP) and introduces solver-based reconciliation. Although this approach benefits from advances in ILP solvers and distributed optimization, it remains NP-hard and often incurs substantial solver overhead. In practice, even moderate-sized subproblems can lead to high memory usage and long runtimes. Alternatively, recent methods like Shah et al. (2024) retain heuristic-based fusion steps, which are efficient but similarly sensitive to noise and lack theoretical support.

In this paper, we propose VISTA (Voting-based Integration of Subgraph Topologies for Acyclicity), a novel modular framework for large-scale causal discovery. The method proceeds in three main stages. First, for each variable we identify its Markov Blanket, thereby reducing the global problem into tractable local neighborhoods. A base learner is then applied to each neighborhood using the data restricted to that subset of variables, producing local subgraphs. Second, these local subgraphs are aggregated through an adaptive voting mechanism that down-weights low-support edges, suppressing statistical noise and inconsistencies. Finally, the aggregated graph is post-processed with an efficient approximation algorithm that enforces acyclicity while preserving as many high-confidence orientations as possible. We also establish a theoretical result showing that the overall error rate of the procedure is bounded above by that of the subgraph-level aggregation, ensuring soundness of the divide-and-conquer strategy.

Crucially, VISTA is strictly model-agnostic and highly efficient. It makes no assumptions about the internal design or inductive biases of the base learners, places no restrictions on the choice of Markov Blanket identification algorithm, and imposes no conditions on the underlying data distribution beyond standard faithfulness assumptions. It operates purely on the edge-level outputs of local subgraphs and requires only a one-time  $\mathcal{O}(|V|^2)$  aggregation without any additional solver or training overhead. This lightweight design makes VISTA framework readily compatible with any causal discovery method while enabling broad applicability across baselines and full parallelism in the divide phase.

Our key contributions include:

- We propose VISTA, a model-agnostic and modular framework that decomposes global DAG learning into node-centered Markov Blanket subgraphs. It is fully plug-and-play with respect to MB identification and local learners, requiring no identifiability or distributional assumptions on the chosen base learners.
- Our aggregation is lightweight, efficient, and edge-level, performing a one-pass weighted voting instead of relying on expensive global searches or solver-based optimization. We derive finite-sample error bounds and an asymptotic consistency guarantee for this aggregation, which explicitly calibrates errors from imperfect base learners.
- Extensive experiments across diverse graphs and a wide range of base learners demonstrate that VISTA remedies the typical performance drop of base learners, consistently improving robustness and scalability over standalone baselines.

## 2 PRELIMINARIES

**Setup and notation.** Let  $\mathbf{V} = \{V_1, \dots, V_n\}$  be random variables generated by a structural causal model with mutually independent noises  $\epsilon_i$ :

$$V_i = f_i(\text{Pa}(V_i), \epsilon_i), \quad \epsilon_i \perp\!\!\!\perp \text{Pa}(V_i).$$

This induces a directed acyclic graph (DAG)  $\mathcal{G} = (\mathbf{V}, \mathbf{E})$  where  $V_i \rightarrow V_j \in \mathbf{E}$  iff  $V_i$  appears in  $f_j$ , and the observational distribution factorizes as  $\mathbb{P}(\mathbf{V}) = \prod_{i=1}^n \mathbb{P}(V_i \mid \text{Pa}(V_i))$ .

**Markov Blanket locality.** Assuming causal sufficiency for exposition, the *Markov Blanket*  $\text{MB}(V)$  of a node  $V$  is the minimal set that renders  $V$  independent of all others given  $\text{MB}(V)$ ; it consists of parents, children, and *spouses* (other parents of the children). Equivalently,  $\text{MB}(V)$   $d$ -separates  $V$  from  $\mathbf{V} \setminus (\{V\} \cup \text{MB}(V))$ . This locality motivates our divide-conquer design: by learning  $\text{MB}(V)$ , causal discovery can be restricted to the induced subgraph  $\mathcal{G}[\{V\} \cup \text{MB}(V)]$ , substantially reducing search complexity while preserving relevant adjacencies for  $V$ .

**Existing Modular Causal Discovery Paradigms.** For large-scale causal discovery, several local-to-global or fusion-style schemes decompose a graph and then merge the pieces: a top-down CI-driven partition with set-based stitching Xie & Geng (2008), global fusion over multiple full Bayesian networks Puerta et al. (2021), a separation–reunion pipeline that repeatedly searches the structure Liu et al. (2017), a PC-style progressive skeleton requiring iterative bootstraps Guo et al. (2024), and DCILP, which formulates reconciliation as an ILP Dong et al. (2024). However, these methods are typically algorithm-specific rather than modular frameworks; they either assume correct inputs at merging time, depend on heavy global search or solver-based optimization, or perform essentially uncalibrated frequency-based stitching. [There also exists a SADA-based or extended model Cai et al. \(2013; 2018\); Rahman et al. \(2021\), but it is limited to LiNGAM and lacks a calibration process during merging.](#) By contrast, our framework provides a lightweight, calibrated weighted-voting aggregation that down-weights low-support directions and remains compatible with arbitrary base learners. A more detailed related work discussion appears in Appendix B.

### 3 METHODOLOGY

We introduce VISTA (Voting-based Integration of Subgraph Topologies for Acyclicity), a novel modular framework for large-scale DAG learning that is both model-agnostic and efficient. Instead of searching the full graph, VISTA focuses on edge-level evidence: for each node  $V$ , we form the subgraph induced by  $\{V\} \cup \text{MB}(V)$  and run any off-the-shelf local learner, regardless of its parametric form, identifiability assumptions, or internal design. The resulting local predictions are reconciled by a lightweight weighted voting on each ordered pair  $(X, Y)$ , which calibrates errors from imperfect base learners, and acyclicity is then enforced by a Feedback Arc Set heuristic Eades et al. (1993). This modular design makes VISTA fully plug-and-play: MB identification and local learning can be tailored to the data regime, while aggregation and acyclicity remain fixed, scalable, and consistent.

**Proposition 3.1** (Coverage of a DAG by Markov-Blanket Subgraphs). *Let  $\mathcal{G} = (V, E)$  be a DAG. For each  $V \in V$ , define*

$$\mathcal{G}' = \bigcup_{V_i \in V} \mathcal{G}[\{V_i\} \cup \text{MB}(V_i)]. \quad (1)$$

*Then every edge of  $\mathcal{G}$  is present in  $\mathcal{G}'$ , i.e.,  $E \subseteq E(\mathcal{G}')$ .*

*Proof.* Take any edge  $(X, Y) \in E$ . If  $X \rightarrow Y$ , then  $Y$  is a child of  $X$  and  $X$  is a parent of  $Y$ , hence  $Y \in \text{MB}(X)$  and  $X \in \text{MB}(Y)$ . Therefore  $(X, Y)$  appears in  $\mathcal{G}[\{X\} \cup \text{MB}(X)]$  and in  $\mathcal{G}[\{Y\} \cup \text{MB}(Y)]$ , and thus in the union  $\mathcal{G}'$ .  $\square$

This coverage property is the foundation of VISTA: once MBs and their local subgraphs are correctly identified, no true edge is lost in the decomposition. Importantly, our framework remains *agnostic* to the specific MB estimator or local learner, that any method suitable for the data distribution can be plugged in. All subsequent aggregation and acyclicity enforcement operate purely at the edge level and rely only on this coverage guarantee. Besides, as shown in Figure 1, the accuracy of MB identification remains relatively stable as the number of nodes increases, whereas the performance of base learners degrades more sharply. This empirical observation is consistent with our theoretical analysis in Section 3.2, where we prove that the proposed merging scheme converges to the correct edge orientations. Furthermore, across different graph sizes, the VISTA-enhanced versions consistently outperform their corresponding baselines, demonstrating the robustness of our framework.

Moreover, since our framework is agnostic to the choice of MB identification methods, we also provide a flexible interface in our implementation that allows practitioners to plug in any suitable MB estimator depending on the specific

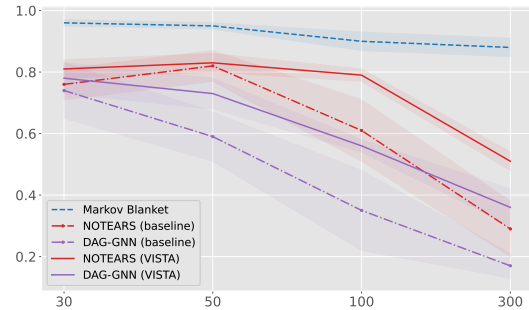


Figure 1: F1 score comparison as the number of nodes increases.

data distribution. Notably, we assume that each base learner outputs directed edges on local subgraphs throughout this work. If an undirected adjacency  $X - Y$  is returned, it is treated as providing no directional vote in the aggregation.

### 3.1 VISTA: VOTING-BASED INTEGRATION OF SUBGRAPH TOPOLOGIES FOR ACYCLICITY

**Naive Voting (NV)** To merge estimated subgraphs into a globally causal graph, we first consider a naive voting strategy. For each pair of nodes  $X$  and  $Y$ , let  $A$  denote the number of times the directed edge  $X \rightarrow Y$  appears across all subgraphs, and  $B$  denote the number of times  $Y \rightarrow X$  appears. The directional support ratio for each orientation is computed as:

$$r_{X \rightarrow Y} = \frac{A}{A+B}, \quad r_{Y \rightarrow X} = \frac{B}{A+B}.$$

This NV rule serves to demonstrate an important property of our divide-and-conquer framework. By Theorem 3.1, every ground-truth causal edge must appear in the union of MB subgraphs. Therefore, even this unweighted scheme, which simply aggregates raw directional votes, already ensures that all true edges are included in the candidate pool. In other words, NV validates that our subgraph decomposition does not lose any causal edges, providing an essential guarantee for the global reconstruction stage.

However, while NV does not distinguish between strong and weak statistical support. Edges appearing rarely across subgraphs receive the same confidence as frequently supported ones, and directional conflicts cannot be resolved in a principled manner. These issues motivate the introduction of our weighted voting formulation, which incorporates frequency-based confidence to produce more reliable global orientation decisions.

**Weighted Voting (WV)** For each pair of nodes  $X$  and  $Y$ , let  $A$  and  $B$  denote the number of times  $X \rightarrow Y$  and  $Y \rightarrow X$  appear across all subgraphs, respectively, and let  $m = A + B$  be the total occurrence. We define the confidence-adjusted score as:

$$s(X \rightarrow Y) = (1 - e^{-\lambda m}) \frac{A}{m}, \quad (2)$$

where  $\lambda > 0$  is a tunable weighting parameter. An edge  $X \rightarrow Y$  is retained if  $s(X \rightarrow Y) \geq t$ , where  $t \in (0, 1)$  is a global decision threshold.

Here, the weighting term  $(1 - e^{-\lambda m})$  serves as a soft confidence modulator that adapts to the reliability of directional evidence. It plays a role analogous to smoothing priors in Bayesian estimation, where rare events are regularized toward lower confidence. The details are illustrated in Appendix D.1. The inclusion threshold  $t$  determines the minimum score required to retain an edge.

Compared to naive voting, which treats all local decisions equally, the weighted scheme jointly calibrates confidence and sparsity. Specifically, the parameter  $\lambda$  penalizes edges with weak support, while the threshold  $t$  determines the final inclusion criterion. Together, the two parameters govern the precision-recall trade-off, since a larger  $\lambda$  tends to preserve edges with limited but consistent evidence and thus improves recall, while a higher  $t$  enforces stricter acceptance and thereby improves precision. This mechanism is particularly beneficial in sparse graphs, where many candidate edges receive only minimal support; the exponential weighting amplifies even small differences in frequency, effectively suppressing unreliable edges. As a result, the aggregation remains robust without relying on strong parametric assumptions, and it provides a tunable handle for balancing false discoveries and missed edges. Beyond the divide-and-conquer efficiency of VISTA, the weighted voting strategy itself enhances the performance of base learners,

```
def VISTA(nodes, base_learner, ...,
          MB_solver, lam, t):
    local_graphs = []

    for v in nodes:
        MB_v = MB_solver(v)
        G_v = base_learner(MB_v ∪ v)
        local_graphs.append(G_v)

    G_merged = WV(local_graphs, lam, t)
    G_final = post_prune(G_merged)
    return G_final
```

Figure 2: Pseudocode of VISTA framework

yielding substantial gains in recall while tightening theoretical error bounds. A detailed analysis of these effects is provided in Section 3.2 and Appendices D - E.

**Acyclicity guarantee** While the weighted voting improves robustness, the resulting merged graph may still contain cycles. To ensure that the final output is a valid DAG, it is necessary to explicitly break loops introduced during the merging process. So we explicitly enforce acyclicity by solving a Feedback Arc Set (FAS) problem Simpson et al. (2016). As FAS is NP-hard, we adopt a fast GreedyFAS heuristic Eades et al. (1993) adapted to weighted edges; the implementation is detailed in Algorithm 2 in Appendix C.

Notably, an important implementation detail involves the ordering between GreedyFAS and threshold-based filtering. In VISTA, cycles are first removed using GreedyFAS, after which edges with weights below a global threshold  $t$  are filtered out. This ordering avoids forcing the cycle removal step to act on already sparse graphs, where eliminating a cycle may require discarding high-confidence edges. In contrast, applying filtering before GreedyFAS can lead to unnecessary precision loss, as the remaining cycles must be resolved by removing stronger edges that would otherwise have been preserved. Besides, taking a subset of nodes from a causal graph introduces unobserved confounding, which will lead to additional edges in the subgraph; the post-processing step here can remove part of these redundant edges.

In general, our VISTA offers several key advantages that make it particularly suited for large-scale causal discovery. It operates purely on aggregated edge counts and requires only matrix-level operations, with no reliance on optimization solvers or iterative training. Importantly, it is model-agnostic, i.e., the aggregation is independent of the internal structure of base learners and can be applied to any method that outputs directed subgraphs. This modularity allows seamless integration with a broad class of causal discovery algorithms and supports parallel execution in the divide stage. The complete procedure is implemented as a simple and modular pipeline, summarized in Figure 3.

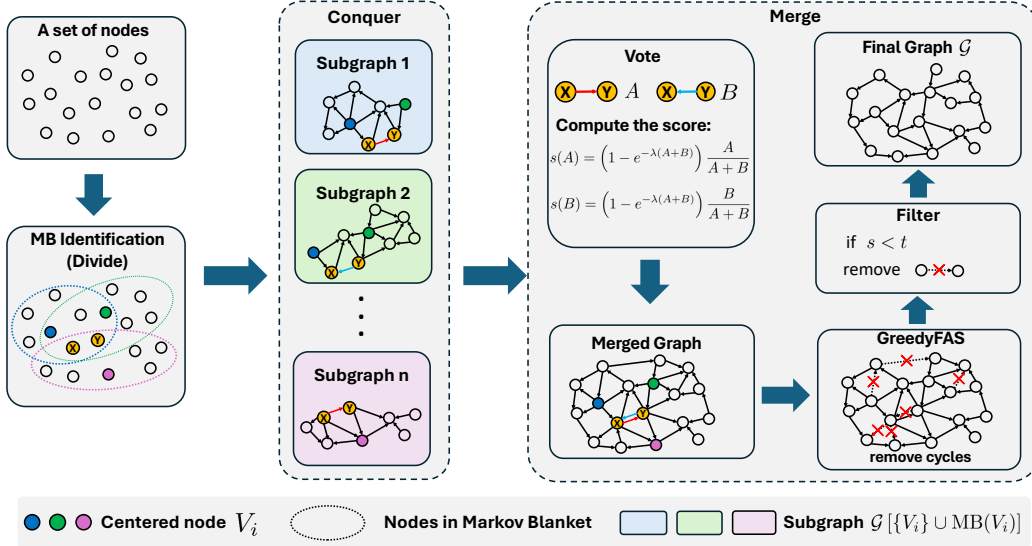


Figure 3: Overview of VISTA, a modular framework for causal discovery: (1) dividing via Markov Blankets identification, (2) parallel subgraph structure identification using a base learner, and (3) global aggregation through weighted voting. The framework then applies cycle resolution (GreedyFAS) and weight-based filtering to produce the final DAG.

**Theoretical guarantees for Weighted Voting** To ensure the reliability of our edge orientation decisions based on the weighted voting mechanism described above, we provide theoretical guarantees derived from concentration inequalities. The core idea is to determine the minimum number of votes (subgraphs)  $m$  required to achieve a desired level of confidence  $1 - \epsilon$  in our decision.

**Theorem 3.2** (Sufficient Condition for Weighted Voting Accuracy). *Let  $A \sim \text{Binomial}(m, p)$  represent the number of successful votes in  $m$  independent subgraphs for the edge direction  $X \rightarrow Y$ , where each subgraph supports this direction independently with probability  $p \in (0, 1)$ , a decision threshold  $t \in (0, 1)$  and the weight function  $w(m) = 1 - e^{-\lambda m}$ ,  $\lambda > 0$ . Assume the effective*

threshold for accepting the edge direction  $X \rightarrow Y$  is  $r(m) = \frac{t}{1-e^{-\lambda m}} < p$ , i.e., the true support rate  $p$  is above the effective threshold. Then, if

$$\frac{mp}{2} \left( 1 - \frac{t}{p(1-e^{-\lambda m})} \right)^2 \geq \log \frac{1}{\epsilon}, \quad (3)$$

it follows that  $P(s(A) \geq t) \geq 1 - \epsilon$ .

This theorem guarantees that if  $m$  is large enough to satisfy the given inequality, the weighted voting procedure will correctly identify the edge direction with high probability. The condition highlights that the required  $m$  depends on the squared relative difference between the true probability  $p$  and the effective threshold  $r(m)$ . Note that  $r(m)$  itself depends on  $m$  and  $\lambda$ . As  $m$  increases or  $\lambda$  increases,  $1 - e^{-\lambda m}$  approaches 1, and  $r(m)$  approaches  $t$ . The inequality requires larger  $m$  and becomes more difficult to satisfy when  $p$  is close to  $r(m)$  or when higher confidence is desired. This trade-off illustrates the role of  $\lambda$  in controlling the conservativeness of the decision rule, which we will analyze further in later sections. In practice, the true value of  $p$  is unknown, but we can empirically validate the trend predicted by this condition using observed vote frequencies and measured recovery accuracy across different values of  $\lambda$  and  $t$ .

Notably, Theorem 3.2 is stated under an idealized assumption that the votes from different local subgraphs are independent. In practice, subgraphs learned from the same dataset can induce correlations among votes, so the bound should be interpreted as a qualitative guide, and we expect the same monotone trend to hold more effectively independent votes still reduce error and the gap between  $p$  and the effective threshold continues to govern sample complexity. Extending the theory to low-correlation weakly dependent votes will be an interesting future direction.

**Corollary 3.3** (Lower bound on node in subgraphs). *Let  $\lambda > 0$ ,  $t \in (0, 1)$ , and  $\epsilon \in (0, 1)$  be fixed. For a candidate edge  $(X, Y)$ , denote by  $m$  the number of local subgraphs whose Markov Blankets contain both endpoints. Under the setting of Theorem 3.2, the sufficient condition (3) can be converted into an explicit bound*

$$m \geq \frac{2 \log(1/\epsilon)}{p((1-t/p)^2 - 2(t/p)(1-t/p)e^{-\lambda})}. \quad (4)$$

Generally, a lower error rate  $\epsilon$  leads to a larger  $\log(1/\epsilon)$  term, which increases the required size of  $m$ . When  $p$  is much greater than  $t$ , it results in a small required  $m$ . This aligns with intuition: if the true voting rate  $p$  is far from the threshold  $t$ , the distinction is easier, and fewer votes are needed for reliable decisions. Similarly, when the gap  $p - t$  is small, it will result in a significantly larger required  $m$ . A large lower bound on  $m$  primarily indicates that the current setting yields a very small gap between  $p$  and  $t$ , which, in turn, implies that the decision task has intrinsically high sample complexity.

### 3.2 ERROR BOUND ANALYSIS

We analyze the edge-level errors of the weighted voting rule to understand how the weighting parameter  $\lambda$  and the threshold  $t$  affect false positives and false negatives. We first characterize a sufficient condition that converts  $t$  into a probability threshold and yields a feasible range for  $\lambda$ , and then show that under this regime, weighted voting achieves asymptotic consistency as the graph size grows.

**Theorem 3.4** (Practical choice of  $\lambda$ ). *Fix a vote count  $m \geq 1$ , a decision threshold  $t \in (0, 1)$ , and a target error level  $\epsilon \in (0, 1)$ . If  $\lambda$  satisfies*

$$-\frac{1}{m} \ln(1-t) < \lambda \leq -\frac{1}{m} \ln \epsilon, \quad (5)$$

*then the weighted-vote rule achieves the prescribed error control under the union bound.*

Theorem 3.4 establishes a feasible interval for  $\lambda$  that guarantees uniform control of edge-level errors. While the confidence weight  $1 - e^{-\lambda m}$  down-weights low-support orientations at a fixed  $t$ , the smaller  $\lambda$  values impose stricter thresholds  $r_\lambda(m)$  to suppress low-support edges, while larger values retain weaker true edges and improve recall. The proof of the theorem, as well as detailed discussions, is in Appendix E.1. In practice, we adopt the relatively large admissible  $\lambda$  in (5), which lowers the



effective threshold and reduces false negatives at the cost of more false positives. This choice is well suited to sparse graphs since false positives typically dominate. The empirical behavior of varying  $\lambda$  is further examined in Section 4.1. Notably, as  $\lambda \rightarrow 0$ , the rule reduces to naive voting with a fixed threshold  $t$ . Building on the finite-sample guarantees above, we next analyze the asymptotic behavior of the weighted voting rule as the number of variables grows. Similarly to  $p$ , let  $q \in (0, 1)$  denote the probability that a false edge is erroneously included. In practice, both  $p$  and  $q$  can be empirically estimated.

**Theorem 3.5** (Asymptotic Consistency). *Fix a threshold  $t \in (0, 1)$  and let  $\delta_p = p - t$  and  $\delta_q = t - q$  denote the positive margins between  $t$  and the inclusion probabilities  $p, q$  of true and false edges respectively. Assume  $\delta_p, \delta_q > 0$  and that  $\lambda$  satisfies the conditions in Theorem 3.4. If the number of local subgraphs per candidate edge is  $m = C \log n$  with  $C > \frac{2}{\min\{\delta_p^2, \delta_q^2\}}$ , then we have*

$$\Pr(\text{global error}) = o(1), \quad \text{as } n \rightarrow \infty. \quad (6)$$

Since most base solvers are reliable and can correctly identify a substantial fraction of true edges, our assumptions are quite mild and practically easy to satisfy. Theorem 3.5 establishes that weighted voting is asymptotically consistent: as the number of subgraph samples increases, the probability of edge-level misclassification vanishes. Notably, the required number of independent subgraphs per edge grows only logarithmically with the graph size, i.e.,  $\mathcal{O}(\log n)$ , making the approach efficient. From a computational perspective, the global merging procedure involves only one pass of edge counting and scoring, with an overall complexity  $\mathcal{O}(n^2)$  regardless of the base learner. These guarantees jointly ensure that the method remains scalable and reliable for large-scale structure discovery. The proof of the theorem is provided in Appendix E.3.

## 4 EXPERIMENTS

### 4.1 SYNTHETIC DATA

We empirically evaluate the performance of the proposed VISTA framework on a range of graph structures and sizes, as well as diverse base learners. To demonstrate the improvement and effectiveness of VISTA, we report representative results that highlight the structural recovery performance of VISTA, its runtime benefits from our modular strategy, and the precision-recall trade-offs induced by different values of  $\lambda$ . All experiments are conducted on a machine equipped with 13th Gen Intel(R) Core(TM) i9-13900HX CPU (24 cores) and NVIDIA A30 GPU (24GB).

**Baselines** We benchmark VISTA against recent typical state-of-the-art causal discovery algorithms, including CAM Bühlmann & Peters (2016), NOTEARS Zheng et al. (2018), DAG-GNN Yu et al. (2019), and GOLEM Ng et al. (2020) for the linear setting, which we modeled as linear Structural Equation Model (SEM) with Gaussian noise, as well as SCORE Rolland et al. (2022) and GraN-DAG Lachapelle et al. (2020) for the nonlinear setting, defined as quadratic SEM. Each baseline is evaluated both in isolation and when integrated with our modular framework VISTA. Additionally, in Appendix F.2, we provide a comparison between VISTA and DCILP Dong et al. (2024), a recent distributed framework for causal structure learning, where we also implemented the MB solver used in that work.

We evaluate the accuracy of our VISTA framework under the Naive Voting (NV) and the Weighted Voting (WV) aggregation schemes. Each base learner is tested standalone and with both VISTA variants. We evaluate the proposed method on synthetic datasets generated from Erdős-Rényi (ER) and scale-free (SF) graphs, with average out-degree  $h \in \{3, 5\}$  and number of nodes  $n \in \{30, 50, 100, 300\}$ . Performance is assessed using False Discovery Rate (FDR), True Positive Rate (TPR), Structural Hamming Distance (SHD), and F1 score, as well as runtime metrics. Experiments are conducted under multiple simulation settings, and we report the average performance, with the  $\pm$  values indicating the corresponding standard deviations.

**Results** Table 1 shows two complementary roles of our aggregation. The NV variant already lifts recall by pooling evidence from overlapping neighborhoods, recovering more true edges. Building on this, WV acts as a principled edge-level filter. By down-weighting orientations with small or inconsistent support and applying a single global threshold, it removes noisy connections and yields substantially cleaner structures. Quantitatively, WV reduces FDR by 50 ~ 80% relative to the original baselines and by 40 ~ 70% compared to NV, while generally keeping TPR no less than 0.70.

Table 1: Results with linear and nonlinear synthetic datasets ( $n = 100, h = 5$ ).

| Method    | ER5                |                    |                        |                    | SF5                |                    |                        |                    |
|-----------|--------------------|--------------------|------------------------|--------------------|--------------------|--------------------|------------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                   | F1↑                | FDR↓               | TPR↑               | SHD↓                   | F1↑                |
| NOTEARS   | 0.21 ± 0.21        | 0.74 ± 0.26        | 208.80 ± 199.71        | 0.76 ± 0.24        | 0.37 ± 0.15        | 0.60 ± 0.14        | 352.60 ± 125.39        | 0.61 ± 0.14        |
| +VISTA-NV | 0.87 ± 0.01        | <b>0.97 ± 0.01</b> | 3171.80 ± 174.02       | 0.23 ± 0.01        | 0.84 ± 0.01        | <b>0.97 ± 0.01</b> | 2443.60 ± 143.74       | 0.27 ± 0.01        |
| +VISTA-WV | <b>0.08 ± 0.03</b> | 0.68 ± 0.01        | <b>182.40 ± 16.03</b>  | <b>0.79 ± 0.02</b> | <b>0.18 ± 0.07</b> | 0.68 ± 0.03        | <b>233.00 ± 34.76</b>  | <b>0.74 ± 0.03</b> |
| GOLEM     | 0.61 ± 0.16        | 0.35 ± 0.17        | 567.00 ± 129.77        | 0.35 ± 0.15        | 0.70 ± 0.15        | 0.29 ± 0.19        | 610.10 ± 118.00        | 0.29 ± 0.17        |
| +VISTA-NV | 0.87 ± 0.01        | <b>0.91 ± 0.04</b> | 2891.00 ± 224.42       | 0.23 ± 0.01        | 0.86 ± 0.01        | <b>0.90 ± 0.02</b> | 2589.00 ± 270.09       | 0.25 ± 0.02        |
| +VISTA-WV | <b>0.23 ± 0.12</b> | 0.50 ± 0.13        | <b>306.70 ± 87.75</b>  | <b>0.60 ± 0.14</b> | <b>0.33 ± 0.15</b> | 0.40 ± 0.12        | <b>371.10 ± 88.21</b>  | <b>0.50 ± 0.13</b> |
| DAG-GNN   | 0.66 ± 0.15        | 0.42 ± 0.23        | 739.20 ± 323.34        | 0.35 ± 0.17        | 0.64 ± 0.15        | 0.47 ± 0.22        | 731.40 ± 303.38        | 0.38 ± 0.17        |
| +VISTA-NV | 0.87 ± 0.01        | <b>0.95 ± 0.01</b> | 3065.00 ± 136.49       | 0.23 ± 0.01        | 0.85 ± 0.01        | <b>0.95 ± 0.00</b> | 2480.00 ± 203.65       | 0.27 ± 0.01        |
| +VISTA-WV | <b>0.36 ± 0.03</b> | 0.56 ± 0.05        | <b>377.00 ± 26.06</b>  | <b>0.59 ± 0.02</b> | <b>0.35 ± 0.10</b> | 0.49 ± 0.08        | <b>363.00 ± 41.10</b>  | <b>0.56 ± 0.09</b> |
| GraN-DAG  | 0.92 ± 0.04        | 0.05 ± 0.03        | 715.00 ± 70.14         | 0.06 ± 0.04        | 0.94 ± 0.02        | 0.05 ± 0.03        | 1088.60 ± 31.49        | 0.05 ± 0.02        |
| +VISTA-NV | 0.86 ± 0.04        | <b>0.18 ± 0.06</b> | 656.60 ± 83.30         | 0.16 ± 0.03        | 0.89 ± 0.02        | <b>0.20 ± 0.04</b> | 947.20 ± 53.33         | 0.14 ± 0.02        |
| +VISTA-WV | <b>0.43 ± 0.06</b> | 0.10 ± 0.02        | <b>503.40 ± 46.68</b>  | <b>0.17 ± 0.03</b> | <b>0.54 ± 0.05</b> | 0.11 ± 0.02        | <b>545.80 ± 65.54</b>  | <b>0.18 ± 0.03</b> |
| SCORE     | 0.92 ± 0.10        | 0.58 ± 0.03        | 4039.60 ± 123.3        | 0.14 ± 0.15        | 0.91 ± 0.03        | 0.62 ± 0.05        | 3166.40 ± 258.7        | 0.16 ± 0.05        |
| +VISTA-NV | 0.95 ± 0.08        | <b>0.76 ± 0.02</b> | 3464.20 ± 215.6        | 0.09 ± 0.14        | 0.95 ± 0.04        | <b>0.76 ± 0.05</b> | 2978.00 ± 367.3        | 0.08 ± 0.07        |
| +VISTA-WV | <b>0.80 ± 0.06</b> | 0.65 ± 0.07        | <b>838.00 ± 364.78</b> | <b>0.31 ± 0.09</b> | <b>0.81 ± 0.05</b> | 0.63 ± 0.04        | <b>892.60 ± 345.58</b> | <b>0.29 ± 0.06</b> |

The trend holds for both differentiable and combinatorial base learners, indicating that the gains stem from the aggregation rule rather than any particular estimator.

Crucially,  $\lambda$  appears only in the final aggregation, so sweeping it is retraining-free: we reuse cached votes, recompute  $r_\lambda(m)$ , and rerun the DAG projection to obtain the full curves. To avoid per-dataset hyperparameter tuning and cherry-picking, all VISTA results in the main tables use a single, fixed operating point:  $\lambda = 0.5$  and  $t = 0.7$ . This choice lies within (5) and serves as a stable compromise between precision and recall across settings. We report the full precision–recall curves for transparency, but no post-hoc selection is performed for the tabulated results.

The observed improvement in WV cases against NV aligns with Theorem 3.4. Edges with limited empirical support are selectively pruned while strongly supported ones are preserved, which is exactly the filtering behavior reflected in Table 1. This validates our weighted voting scheme as an effective, model-agnostic mechanism for stabilizing global structures. To further substantiate this model-agnostic property, we next examine the impact of data standardization as it is known to influence baseline performance Reisach et al. (2021).

Table 2: Results with normalized linear and nonlinear synthetic datasets ( $n = 50, h = 5$ ).

| Method    | ER5                |                    |                       |                    | SF5                |                    |                       |                    |
|-----------|--------------------|--------------------|-----------------------|--------------------|--------------------|--------------------|-----------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                  | F1↑                | FDR↓               | TPR↑               | SHD↓                  | F1↑                |
| NOTEARS   | <b>0.04 ± 0.02</b> | 0.39 ± 0.01        | 140.00 ± 4.90         | 0.56 ± 0.01        | <b>0.02 ± 0.02</b> | 0.38 ± 0.04        | 138.50 ± 9.87         | 0.55 ± 0.05        |
| +VISTA-NV | 0.27 ± 0.05        | <b>0.61 ± 0.03</b> | 135.20 ± 6.16         | 0.66 ± 0.02        | 0.35 ± 0.04        | <b>0.62 ± 0.04</b> | 132.80 ± 18.82        | 0.63 ± 0.03        |
| +VISTA-WV | 0.19 ± 0.05        | 0.58 ± 0.03        | <b>122.90 ± 7.54</b>  | <b>0.68 ± 0.02</b> | 0.08 ± 0.04        | 0.54 ± 0.06        | <b>109.10 ± 19.91</b> | <b>0.68 ± 0.05</b> |
| GOLEM     | 0.40 ± 0.03        | 0.22 ± 0.04        | 182.00 ± 15.51        | 0.32 ± 0.05        | 0.44 ± 0.07        | 0.20 ± 0.04        | 183.60 ± 6.55         | 0.29 ± 0.05        |
| +VISTA-NV | 0.31 ± 0.03        | <b>0.75 ± 0.03</b> | 129.50 ± 4.97         | 0.72 ± 0.02        | 0.29 ± 0.05        | <b>0.70 ± 0.05</b> | 122.80 ± 19.87        | 0.70 ± 0.04        |
| +VISTA-WV | <b>0.06 ± 0.03</b> | 0.62 ± 0.04        | <b>95.30 ± 9.88</b>   | <b>0.75 ± 0.02</b> | <b>0.10 ± 0.04</b> | 0.60 ± 0.06        | <b>100.20 ± 15.69</b> | <b>0.72 ± 0.05</b> |
| DAG-GNN   | 0.16 ± 0.03        | 0.41 ± 0.05        | 160.80 ± 53.55        | 0.55 ± 0.05        | 0.19 ± 0.05        | 0.48 ± 0.04        | 183.60 ± 45.37        | 0.60 ± 0.03        |
| +VISTA-NV | 0.85 ± 0.09        | <b>0.74 ± 0.14</b> | 609.80 ± 72.70        | 0.25 ± 0.12        | 0.79 ± 0.04        | <b>0.72 ± 0.09</b> | 538.40 ± 25.55        | 0.33 ± 0.05        |
| +VISTA-WV | <b>0.14 ± 0.05</b> | 0.50 ± 0.09        | <b>93.50 ± 29.12</b>  | <b>0.63 ± 0.07</b> | <b>0.13 ± 0.08</b> | 0.56 ± 0.06        | <b>87.80 ± 16.56</b>  | <b>0.68 ± 0.05</b> |
| GraN-DAG  | 0.82 ± 0.01        | 0.06 ± 0.01        | 275.00 ± 18.50        | 0.09 ± 0.01        | 0.92 ± 0.02        | 0.02 ± 0.02        | 269.80 ± 45.50        | 0.03 ± 0.02        |
| +VISTA-NV | 0.66 ± 0.15        | <b>0.26 ± 0.06</b> | 219.20 ± 46.41        | 0.29 ± 0.07        | 0.68 ± 0.05        | <b>0.17 ± 0.04</b> | 223.00 ± 26.25        | 0.22 ± 0.04        |
| +VISTA-WV | <b>0.15 ± 0.06</b> | 0.18 ± 0.05        | <b>199.20 ± 13.64</b> | <b>0.32 ± 0.07</b> | <b>0.33 ± 0.03</b> | 0.13 ± 0.03        | <b>205.40 ± 59.15</b> | <b>0.23 ± 0.04</b> |
| SCORE     | 0.71 ± 0.05        | 0.50 ± 0.05        | 386.80 ± 67.99        | 0.37 ± 0.04        | 0.65 ± 0.13        | 0.52 ± 0.15        | 340.40 ± 81.08        | 0.38 ± 0.05        |
| +VISTA-NV | 0.79 ± 0.03        | <b>0.60 ± 0.14</b> | 489.70 ± 123.82       | 0.31 ± 0.04        | 0.77 ± 0.03        | <b>0.56 ± 0.05</b> | 471.10 ± 16.68        | 0.33 ± 0.03        |
| +VISTA-WV | <b>0.64 ± 0.09</b> | 0.42 ± 0.11        | <b>305.80 ± 49.93</b> | <b>0.39 ± 0.07</b> | <b>0.57 ± 0.04</b> | 0.36 ± 0.06        | <b>244.20 ± 53.35</b> | <b>0.39 ± 0.04</b> |

The results show that, regardless of fluctuations in the performance of individual base learners, the improvements brought by VISTA remain consistent. This stability further supports our claim that VISTA does not rely on any inductive bias of the base learner or data distribution. Rather, the edge-level aggregation mechanism provides robustness across settings. These findings further highlight the model-agnostic nature of our framework. Additional experiments under alternative parameter settings are provided in Appendix F.4.

**Time efficiency** To assess the scalability of our framework, we report the total computation time for different base learners in Table 3. All results are presented as mean  $\pm$  standard deviation over repeated runs. Across all tested graph sizes, integrating VISTA consistently yields substantial runtime

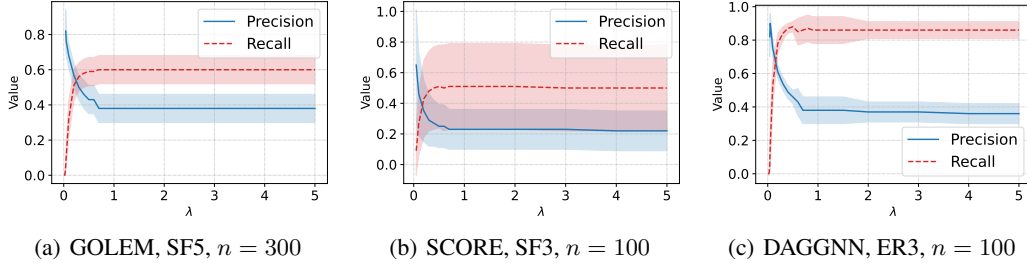


reductions compared to the original methods. These improvements are not due to algorithm-specific acceleration but result directly from our divide-and-conquer design: since each local subgraph is processed independently, the learning procedure naturally supports parallel execution. This decomposition effectively reduces the per-task computational load and alleviates memory bottlenecks, enabling scalable causal discovery even with large node counts. Further results for other settings are included in Appendix F.3.

Table 3: Comparison of total computing time (s) under ER3 setting.

| Method   | $n = 50$                             | $n = 100$                             | $n = 300$                               |
|----------|--------------------------------------|---------------------------------------|---|
| NOTEARS  | 494.40 $\pm$ 98.24                   | 1473.69 $\pm$ 395.59                  | 12515.63 $\pm$ 1599.06                  |
| +VISTA   | <b>189.15 <math>\pm</math> 65.37</b> | <b>339.90 <math>\pm</math> 158.75</b> | <b>2136.72 <math>\pm</math> 708.15</b>  |
| GOLEM    | 72.65 $\pm$ 15.41                    | 108.82 $\pm$ 70.56                    | 261.84 $\pm$ 30.44                      |
| +VISTA   | <b>21.93 <math>\pm</math> 0.81</b>   | <b>26.16 <math>\pm</math> 2.68</b>    | <b>43.40 <math>\pm</math> 3.21</b>      |
| DAG-GNN  | 628.63 $\pm$ 55.29                   | 2192.97 $\pm$ 323.59                  | 17713.84 $\pm$ 2861.06                  |
| +VISTA   | <b>201.31 <math>\pm</math> 43.36</b> | <b>371.25 <math>\pm</math> 199.91</b> | <b>1960.43 <math>\pm</math> 794.02</b>  |
| GraN-DAG | 730.42 $\pm$ 89.95                   | 3035.76 $\pm$ 481.85                  | 25205.64 $\pm$ 2098.85                  |
| +VISTA   | <b>238.53 <math>\pm</math> 51.36</b> | <b>472.30 <math>\pm</math> 172.77</b> | <b>2336.32 <math>\pm</math> 1028.04</b> |
| SCORE    | 426.63 $\pm$ 61.15                   | 10040.65 $\pm$ 209.31                 | —                                       |
| +VISTA   | <b>105.64 <math>\pm</math> 39.65</b> | <b>198.82 <math>\pm</math> 34.12</b>  | <b>225.16 <math>\pm</math> 11.45</b>    |

**Sensitivity study of  $\lambda$**  We sweep  $\lambda$  and plot precision/recall in Figure 4. By the conclusion of Theorem 3.4 and Appendix E.1, larger  $\lambda$  shifts the method toward higher recall and lower precision by relaxing the penalty on low-support edges. Within the theoretical range, this precision–recall trade-off is smooth and yields informative voting thresholds  $r_\lambda(m)$ . The figure also substantiate this point, Small  $\lambda$  strongly discounts low-support edges, yielding high precision and low recall. Similarly, as  $\lambda$  increases, recall rises while precision falls. Beyond the upper end of (5) we have  $(1 - e^{-\lambda m}) \approx 1$  and thus  $s(X \rightarrow Y) \approx A/m$ , so the curves plateau and further increases of  $\lambda$  have negligible effect. Therefore, to balance precision and recall in practice, a moderate value of the hyperparameter could be fixed within the theoretical range, which serves as a stable operating point.

Figure 4: Precision–recall trade-off under varying  $\lambda$ , where threshold  $t = 0.5$ .

## 4.2 REAL DATA

We further evaluate all methods on the well-known Sachs protein signaling network based on expression levels of proteins and phospholipids Sachs et al. (2005). This benchmark is widely used in causal discovery research, and the ground-truth graph with 11 nodes and 17 directed edges is consistently accepted by the community.

Here we trained normalized data with 853 samples and reported the results in Table 4. Incorporating VISTA consistently reduces false discoveries and improves structural accuracy, measured by SHD and SID Peters & Bühlmann (2015) across different baselines. This highlights that VISTA is a plug-and-play module that can reliably enhance the performance of arbitrary causal discovery algorithms.

Table 4: Results on the Sachs protein-signaling network.

| Method        | FDR↓        | TPR↑        | SHD↓      | SID↓      |
|---------------|-------------|-------------|-----------|-----------|
| GOLEM         | 0.80        | 0.26        | 16        | 50        |
| <b>+VISTA</b> | <b>0.57</b> | <b>0.18</b> | <b>16</b> | <b>48</b> |
| SCORE         | 0.81        | 0.18        | 18        | 57        |
| <b>+VISTA</b> | <b>0.60</b> | <b>0.12</b> | <b>15</b> | <b>53</b> |
| DAG-GNN       | 0.50        | 0.12        | 15        | 54        |
| <b>+VISTA</b> | <b>0.25</b> | <b>0.18</b> | <b>14</b> | <b>52</b> |
| GraN-DAG      | 0.82        | 0.53        | 16        | 48        |
| <b>+VISTA</b> | <b>0.00</b> | <b>0.29</b> | <b>12</b> | <b>45</b> |

## 5 CONCLUSION

In this paper, we introduced VISTA, a scalable and model-agnostic framework for causal discovery that decomposes global structure learning into Markov Blanket neighborhoods, aggregates them via a weighted voting scheme, and enforces acyclicity through FAS post-processing. The design is fully parallelizable, and the aggregation step operates only at the edge level, enabling efficient exploration of operating points regardless of the base learner. Theoretically, we establish finite-sample error guarantees and asymptotic consistency under mild conditions. Empirically, across diverse graph families and base learners, VISTA improves accuracy and runtime efficiency, typically increasing precision without sacrificing recall.

Despite the favorable performance of VISTA, the framework has several limitations. First, when aggregating local graphs, latent confounding introduced by restricting the learner to subsets may produce high-confidence redundant edges. In some cases these edges do not necessarily participate in cycles and our current framework can only mitigate them through the combination of GreedyFAS and threshold-based filtering. Moreover, although the FAS projection guarantees acyclicity, it may also prune edges that are weakly supported yet correct, which can negatively affect downstream tasks that are sensitive to edge directions. Future work includes incorporating interventional data to improve orientation accuracy and extending the VISTA framework to online settings for large-scale applications.

## REPRODUCIBILITY STATEMENT

We provide the code in the supplementary material, together with a README file that allows experimental results to be reproduced.

## THE USE OF LLM

We used LLM to polish the writing and correct grammar in some paragraphs, but it did not contribute to ideas or conceptual content.

## REFERENCES

- Steen A Andersson, David Madigan, and Michael D Perlman. A characterization of markov equivalence classes for acyclic digraphs. *The Annals of Statistics*, 25(2):505–541, 1997.
- Taiyu Ban, Lyuzhou Chen, Xiangyu Wang, Xin Wang, Derui Lyu, and Huanhuan Chen. Differentiable structure learning with partial orders. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Kevin Bello, Bryon Aragam, and Pradeep Ravikumar. Dagma: Learning dags via m-matrices and a log-determinant acyclicity characterization. *Advances in Neural Information Processing Systems*, 35:8226–8239, 2022.
- Peter Bühlmann and Jonas Peters. CAM: Causal additive models, high-dimensional order search and penalized regression. *The Annals of Statistics*, 44(1):243–274, 2016.
- Ruichu Cai, Zhenjie Zhang, and Zhifeng Hao. Sada: A general framework to support robust causation discovery. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 208–216, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- Ruichu Cai, Zhenjie Zhang, Zhifeng Hao, and Marianne Winslett. Sophisticated merging over random partitions: A scalable and robust causal discovery approach. *IEEE Transactions on Neural Networks and Learning Systems*, 29(8):3623–3635, 2018. doi: 10.1109/TNNLS.2017.2734804.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.

- David Maxwell Chickering, David Heckerman, and Christopher Meek. Large-sample learning of bayesian networks is np-hard. *Journal of Machine Learning Research*, 5:1287–1330, 2004.
- Davin Choo, Kirankumar Shiragur, and Arnab Bhattacharyya. Verification and search algorithms for causal dags. *Advances in Neural Information Processing Systems*, 35:12787–12799, 2022.
- Diego Colombo, Marloes H Maathuis, Markus Kalisch, and Thomas S Richardson. Learning high-dimensional directed acyclic graphs with latent and selection variables. *The Annals of Statistics*, pp. 294–321, 2012.
- Shuyu Dong, Michèle Sebag, Kento Uemura, Akito Fujii, Shuang Chang, Yusuke Koyanagi, and Koji Maruhashi. DCDILP: a distributed learning method for large-scale causal structure learning. *arXiv preprint arXiv:2406.10481*, 2024.
- Peter Eades, Xuemin Lin, and William F Smyth. A fast and effective heuristic for the feedback arc set problem. *Information Processing Letters*, 47(6):319–323, 1993. doi: 10.1016/0020-0190(93)90079-O.
- Zhuangyan Fang, Shengyu Zhu, Jiji Zhang, Yue Liu, Zhitang Chen, and Yangbo He. On low-rank directed acyclic graphs and causal structure learning. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4):4924–4937, 2023.
- Tian Gao, Kshitij Fadnis, and Murray Campbell. Local-to-global bayesian network structure learning. In *International Conference on Machine Learning*, pp. 1193–1202. PMLR, 2017.
- Jiaying Gu and Qing Zhou. Learning big Gaussian Bayesian networks: Partition, estimation and fusion. *Journal of machine learning research*, 21(158):1–31, 2020.
- Xianjie Guo, Kui Yu, Lin Liu, Jiuyong Li, Jiye Liang, Fuyuan Cao, and Xindong Wu. Progressive skeleton learning for effective local-to-global causal structure learning. *IEEE Transactions on Knowledge and Data Engineering*, 36(12):9065–9079, 2024.
- Alain Hauser and Peter Bühlmann. Characterization and greedy learning of interventional markov equivalence classes of directed acyclic graphs. *The Journal of Machine Learning Research*, 13(1): 2409–2464, 2012.
- Yang-Bo He and Zhi Geng. Active learning of causal networks with intervention experiments and optimal designs. *Journal of Machine Learning Research*, 9(Nov):2523–2547, 2008.
- Biwei Huang, Charles Jia Han Low, Feng Xie, Clark Glymour, and Kun Zhang. Latent hierarchical causal structure discovery with rank constraints. *Advances in neural information processing systems*, 35:5549–5561, 2022.
- Paul Hünermund and Elias Bareinboim. Causal inference and data fusion in econometrics. *The Econometrics Journal*, pp. utad008, 2023.
- Maximilian Kaiser, Stefan Bauer, and Bernhard Schölkopf. Bootstrap aggregation and confidence measures for time-series causal discovery. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, 2024.
- Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural dag learning. In *International Conference on Learning Representations*, 2020.
- Thuc Duy Le, Tao Hoang, Jiuyong Li, Lin Liu, Huawen Liu, and Shu Hu. A fast PC algorithm for high dimensional causal discovery with multi-core pcs. *IEEE/ACM transactions on computational biology and bioinformatics*, 16(5):1483–1495, 2016.
- Phillip Lippe, Taco Cohen, and Efstratios Gavves. Efficient neural causal discovery without acyclicity constraints. *arXiv preprint arXiv:2107.10483*, 2021.
- Hui Liu, Shuigeng Zhou, Wai Lam, and Jihong Guan. A new hybrid method for learning bayesian networks: Separation and reunion. *Knowledge-Based Systems*, 121:185–197, 2017. ISSN 0950-7051.

- Po-Ling Loh and Peter Bühlmann. High-dimensional learning of linear causal networks via inverse covariance estimation. *Journal of Machine Learning Research*, 15(140):3065–3105, 2014.
- Sindy Löwe, David Madras, Richard Zemel, and Max Welling. Amortized causal discovery: Learning to infer causal graphs from time-series data. In *Conference on Causal Learning and Reasoning*, pp. 509–525. PMLR, 2022.
- Christopher Meek. Causal inference and causal explanation with background knowledge. *arXiv preprint arXiv:1302.4972*, 2013.
- Ehsan Mokhtarian, Sina Akbari, AmirEmad Ghassami, and Negar Kiyavash. A recursive markov boundary-based approach to causal structure learning. In *The KDD’21 Workshop on Causal Discovery*, pp. 26–54. PMLR, 2021.
- Francesco Montagna, Nicoletta Noceti, Lorenzo Rosasco, Kun Zhang, and Francesco Locatello. Causal discovery with score matching on additive models with arbitrary noise. In *Conference on Causal Learning and Reasoning*, pp. 726–751. PMLR, 2023a.
- Francesco Montagna, Nicoletta Noceti, Lorenzo Rosasco, Kun Zhang, and Francesco Locatello. Scalable causal discovery with score matching. In *Conference on Causal Learning and Reasoning*, pp. 752–771. PMLR, 2023b.
- Ivan Ng, Xun Zheng, and Bryon Aragam. Learning sparse causal models is not np-hard. In *Advances in Neural Information Processing Systems*, volume 33, pp. 16888–16900, 2020.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Jonas Peters and Peter Bühlmann. Structural intervention distance for evaluating causal graphs. *Neural Computation*, 27(3):771–799, 2015. doi: 10.1162/NECO\_a.00708.
- Jonas Peters, Joris M Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery with continuous additive noise models. *The Journal of Machine Learning Research*, 15(1):2009–2053, 2014.
- Anne Helby Petersen, Claus Thorn Ekstrøm, Peter Spirtes, and Merete Osler. Causal discovery and epidemiology: A potential for synergy. *American Journal of Epidemiology*, pp. kwae101, 2024.
- José M. Puerta, Juan A. Aledo, José A. Gámez, and Jorge D. Laborda. Efficient and accurate structural fusion of bayesian networks. *Information Fusion*, 66:155–169, 2021. ISSN 1566-2535.
- Md Musfiquir Rahman, Ayman Rasheed, Md Mosaddek Khan, Mohammad Ali Javidian, Pooyan Jamshidi, and Md Mamun-Or-Rashid. Accelerating recursive partition-based causal structure learning. *arXiv preprint arXiv:2102.11545*, 2021.
- Joseph Ramsey, Madelyn Glymour, Ruben Sanchez-Romero, and Clark Glymour. A million variables and more: the fast greedy equivalence search algorithm for learning high-dimensional graphical causal models, with an application to functional magnetic resonance images. *International Journal of Data Science and Analytics*, 3(2):121–129, 2017. doi: 10.1007/s41060-016-0032-z.
- Alexander Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag! causal discovery benchmarks may be easy to game. *Advances in Neural Information Processing Systems*, 34:27772–27784, 2021.
- Paul Rolland, Volkan Cevher, Matthäus Kleindessner, Chris Russell, Dominik Janzing, Bernhard Schölkopf, and Francesco Locatello. Score matching enables causal discovery of nonlinear additive noise models. In *International Conference on Machine Learning*, pp. 18741–18753. PMLR, 2022.
- Karen Sachs, Omar Perez, Dana Pe’er, Douglas A. Lauffenburger, and Garry P. Nolan. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721): 523–529, 2005. doi: 10.1126/science.1105809.
- Pedro Sanchez, Xiao Liu, Alison Q O’Neil, and Sotirios A Tsafaris. Diffusion models for causal discovery via topological ordering. *arXiv preprint arXiv:2210.06201*, 2022a.

- Pedro Sanchez, Jeremy P Voisey, Tian Xia, Hannah I Watson, Alison Q O’Neil, and Sotirios A Tsaftaris. Causal machine learning for healthcare and precision medicine. *Royal Society Open Science*, 9(8):220638, 2022b.
- Ashka Shah, Adela DePavia, Nathaniel Hudson, Ian Foster, and Rick Stevens. Causal discovery over high-dimensional structured hypothesis spaces with causal graph partitioning. *arXiv preprint arXiv:2406.06348*, 2024.
- Karthikeyan Shanmugam, Murat Kocaoglu, Alexandros G Dimakis, and Sriram Vishwanath. Learning causal graphs with small interventions. *Advances in Neural Information Processing Systems*, 28, 2015.
- Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara, Takashi Washio, Patrik O Hoyer, Kenneth Bollen, and Patrik Hoyer. Directlingam: A direct method for learning a linear non-Gaussian structural equation model. *Journal of Machine Learning Research-JMLR*, 12(Apr):1225–1248, 2011.
- Michael Simpson, Venkatesh Srinivasan, and Alex Thomo. Efficient computation of feedback arc set at web-scale. *Proceedings of the VLDB Endowment*, 10(3):133–144, 2016.
- Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT press, 2nd edition, 2000.
- Chandler Squires, Sara Magliacane, Kristjan Greenewald, Dmitriy Katz, Murat Kocaoglu, and Karthikeyan Shanmugam. Active structure learning of causal dags via directed clique trees. *Advances in Neural Information Processing Systems*, 33:21500–21511, 2020.
- Ioannis Tsamardinos, Constantin F Aliferis, Alexander R Statnikov, and Er Statnikov. Algorithms for large scale markov blanket discovery. In *FLAIRS*, volume 2, pp. 376–81, 2003.
- Ioannis Tsamardinos, Laura E. Brown, and Constantin F. Aliferis. The max-min hill-climbing bayesian network structure learning algorithm. *Machine Learning*, 65(1):31–78, 2006. doi: 10.1007/s10994-006-6889-7.
- Thomas S Verma and Judea Pearl. Equivalence and synthesis of causal models. In *Probabilistic and causal inference: The works of Judea Pearl*, pp. 221–236. 2022.
- Xiaoqiang Wang, Yali Du, Shengyu Zhu, Liangjun Ke, Zhitang Chen, Jianye Hao, and Jun Wang. Ordering-based causal discovery with reinforcement learning. *arXiv preprint arXiv:2105.06631*, 2021.
- Xingyu Wu, Bingbing Jiang, Yan Zhong, and Huanhuan Chen. Multi-target markov boundary discovery: Theory, algorithm, and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4964–4980, 2022.
- Xingyu Wu, Bingbing Jiang, Tianhao Wu, and Huanhuan Chen. Practical markov boundary learning without strong assumptions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 10388–10398, 2023.
- Xianchao Xie and Zhi Geng. A recursive method for structural learning of directed acyclic graphs. *Journal of Machine Learning Research*, 9(14):459–483, 2008.
- Yuhuai Yu, Xun Zheng, Yichong Cheng, Tengyu Zhao, Oluwasanmi Koyejo, Bo Huang, and Yifan Wang. DAG-GNN: DAG structure learning with graph neural networks. In *International Conference on Machine Learning*, pp. 7154–7163. PMLR, 2019.
- Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P Xing. DAGs with NO TEARS: Continuous optimization for structure learning. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Zihan Zhou, Muhammad Qasim Elahi, and Murat Kocaoglu. Sample efficient bayesian learning of causal graphs from interventions. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Shengyu Zhu, Ignavier Ng, and Zhitang Chen. Causal discovery with reinforcement learning. *arXiv preprint arXiv:1906.04477*, 2019.

## A ALGORITHM OF VISTA

---

### Algorithm 1 VISTA-Weighted Voting

---

**Require:** A set of local subgraphs  $\{\mathcal{G}_V : V \in \mathbf{V}\}$ , each induced by the Markov Blanket of node  $V$ ; hyperparameters  $\lambda$  and threshold  $t$ .

- 1: Initialize zero matrix `EdgeCount` to record counts for each edge  $V_i \rightarrow V_j$  where  $V_i, V_j \in \mathbf{V}$ .
- 2: **for** each local subgraph  $\mathcal{G}_V$  **do**
- 3:   **for** each directed edge  $V_i \rightarrow V_j$  in  $\mathcal{G}_V$  **with**  $i \neq j$  **do**
- 4:     Increment `EdgeCount`[ $V_i, V_j$ ] by 1.
- 5:   **end for**
- 6: **end for**
- 7: Compute the Occurrence matrix as `Occurrence`  $\leftarrow$  `EdgeCount` + `EdgeCount`<sup>T</sup>
- 8: Compute the coefficient matrix elementwise: `Coef`  $\leftarrow$   $1 - \exp(-\lambda \cdot \text{Occurrence})$ .
- 9: Compute the merged weighted directed graph  $\mathcal{G}_1 = \text{Coef} \odot \text{EdgeCount} / \text{Occurrence}$ .
- 10: Use Algorithm 2 to break cycles in  $\mathcal{G}_1$  and obtain a DAG  $\mathcal{G}_2$ .
- 11: Remove edges in  $\mathcal{G}_2$  whose weights are less than threshold  $t$  to obtain the final DAG  $\mathcal{G}$ .
- 12: **return** the global causal graph  $\mathcal{G}$ .

---

## B DETAILED RELATED WORKS

**General Causal Discovery Methods:** Classical algorithms recover Directed Acyclic Graphs (DAGs) by either testing conditional independencies or maximizing a score on the discrete space of graphs. Constraint-based methods such as PC and FCI Colombo et al. (2012); Spirtes et al. (2000) iteratively remove edges whose endpoints become independent given bounded-size conditioning sets. Assuming faithfulness, with only observational data, a common result in causal discovery shows that one can only recover the causal graph up to its Markov Equivalence Class (MEC) Andersson et al. (1997); Verma & Pearl (2022). Therefore, interventional data is usually required to fully recover the graph. Many works propose algorithms that aim to learn the graph with minimal interventional data Choo et al. (2022); Hauser & Bühlmann (2012); He & Geng (2008); Shanmugam et al. (2015); Squires et al. (2020); Zhou et al. (2024). Score-based searches, e.g., GES Chickering (2002) and exact DP-based optimizers Chickering et al. (2004), evaluate a decomposable metric (BIC, MDL) while heuristically exploring the super-exponential DAG space. Hybrid strategies typified by MMHC Tsamardinos et al. (2006) first identify each variable’s Markov Blanket and then run a restricted greedy search. Although provably sound under the causal Markov and faithfulness assumptions, all three lines are NP-hard and their run time or memory grows super-polynomially with node count, limiting practical use to  $\lesssim 10^2$  variables.

Ordering-based methods constitute a distinct and increasingly influential category. These approaches first attempt to infer a topological ordering of variables and then determine parent sets accordingly. Early examples such as DirectLiNGAM Shimizu et al. (2011) and RESIT Peters et al. (2014) exploit non-Gaussianity or additive-noise assumptions to infer edge directions from regression residuals. CAM Bühlmann & Peters (2016) extends this idea to nonlinear settings via generalized additive models and greedy order search. More recently, SCORE Rolland et al. (2022) proposes to identify causal ordering by minimizing the variance of the score function, which has inspired several scalable extensions leveraging score-matching or diffusion-based estimation Montagna et al. (2023a;b); Sanchez et al. (2022a). These methods achieve promising empirical results on graphs with thousands of nodes, but typically rely on strong functional assumptions and remain sensitive to latent confounders.

Besides, recent years have seen a growing emphasis on continuous and differentiable formulations in causal structure learning, aiming to overcome the combinatorial challenges associated with discrete DAG optimization. NOTEARS Zheng et al. (2018), DAG-GNN Yu et al. (2019), GraN-DAG Lachapelle et al. (2020), and their low-rank or log-det variants Bello et al. (2022); Fang et al. (2023) convert acyclicity into a smooth penalty and learn graphs via gradient descent. Reinforcement learning and meta-learning schemes Wang et al. (2021); Zhu et al. (2019); Lippe et al. (2021) treat node ordering as a policy and bypass explicit acyclicity constraints. These methods alleviate combinatorial search but still entail an  $O(d^2)$  adjacency parameterization or an  $O(d^3)$  matrix exponential, so GPU



memory becomes a bottleneck beyond a few hundred nodes. In summary, although continuous optimization and ordering-based heuristics mitigate the need for discrete search, general-purpose methods typically incur  $\mathcal{O}(d^2)$  memory overhead or rely on restrictive assumptions, which constrains their applicability to graphs of moderate size.

**Large-Scale Causal Discovery:** To push causal discovery into the high-dimensional regime, researchers have explored sparsity-aware and parallel variants of the above paradigms. Fast Greedy Search (FGS) Ramsey et al. (2017) and parallel-PC Le et al. (2016) cache CI tests and distribute computations over multi-core CPUs, handling tens of thousands of genes. In the continuous camp, DAGMA Bello et al. (2022) and NOTEARS-LowRank Fang et al. (2023) reduce memory usage by factorizing the weight matrix, achieving 5k–10k nodes on a single GPU, while Amortized Causal Discovery Löwe et al. (2022) shares a latent decoder across samples to scale to massive time-series. Bootstrap and bagging strategies aggregate multiple weak graphs to improve stability without increasing per-run complexity Wu et al. (2023); Kaiser et al. (2024). Despite these advances, most scalable algorithms either rely on heavy solvers (e.g., SDP/MILP), strong sparsity assumptions, or lack finite-sample guarantees, motivating alternative divide-and-conquer solutions. As a complementary approach, our proposed VISTA framework addresses these challenges through modular subgraph decomposition and lightweight aggregation, while providing finite-sample error control and scalability to graphs with a large scale of nodes.

**Scalable or Modular Structure Learning:** Partition-based approaches decompose the global graph into overlapping neighbourhoods, learn local substructures, and then reconcile conflicts. Early local-to-global techniques grow random neighbourhoods until conditional independence saturates Gao et al. (2017). Gu & Zhou (2020); Huang et al. (2022) apply hierarchical clustering before local search, whereas Shah et al. (2024) first estimates a coarse skeleton and then partitions it to learn subgraphs in parallel. DCILP Dong et al. (2024) formulates the fusion step as an integer program that guarantees optimal conflict resolution but suffers from MILP infeasibility on dense regions. Recent ensemble methods perform Markov-Blanket bootstrap with majority or confidence-weighted voting Wu et al. (2023); Ban et al. (2024), yet provide limited theoretical analysis of the aggregated error.

Our method VISTA follows the divide-and-conquer paradigm but departs from prior work by integrating a frequency-aware weighted voting mechanism that admits closed-form error analysis, and by enforcing global acyclicity through a lightweight GreedyFAS post-processing step instead of solving large-scale ILPs. These design choices lead to near-linear memory usage, full parallelizability, and theoretical consistency guarantees, enabling scalable causal discovery on graphs with thousands of nodes.

## C PSEUDOCODE OF FEEDBACK ARC SET

After obtaining a directed graph with weighted edges from the voting stage, the final step is to enforce acyclicity, formulated as a *feedback arc set (FAS)* problem. Since exact FAS is NP-hard, we adopt a greedy approximation based on node degree imbalance.

For each node  $V_i \in \mathbf{V}$ , let  $d^o(V_i)$  and  $d^i(V_i)$  be its out- and in-degrees, and define imbalance  $\delta(V_i) = d^o(V_i) - d^i(V_i)$ . At each iteration, we remove one node: sources are appended to a sequence  $s_1$ , sinks are prepended to a sequence  $s_2$ , and if neither exists, we select the node with the largest absolute imbalance  $|\delta(V_i)|$ . This process continues until all nodes are removed, yielding a topological order  $s = s_1 // s_2$ .

Given this order, any edge  $(V_i, V_j) \in \mathbf{E}$  that points from a later node to an earlier node in  $s$  is marked as a backward edge. These are sorted by weight and the lightest ones are iteratively removed until the graph becomes acyclic. Algorithm 2 summarizes the procedure.

## D STATISTICAL ACCURACY ANALYSIS OF WEIGHTED VOTING

This section provides a theoretical analysis of the statistical behavior of the weighted voting mechanism introduced in Section 3.1. The goal is to characterize the conditions under which a candidate edge is correctly retained or excluded based on its empirical directional support. The analysis builds on a probabilistic interpretation of the weighted score as a posterior expectation, and derives sufficient conditions for edge-level accuracy using concentration inequalities.

**Algorithm 2** Solve FAS to guarantee acyclicity on the weighted directed graph

---

**Require:** A weighted directed graph  $\mathcal{G} = (\mathbf{V}, \mathbf{E})$ , where  $E_{UV}$  denotes the edge from  $U \in \mathbf{V}$  to  $V \in \mathbf{V}$  and  $w_{UV} > 0$  is its weight.

- 1: For the ease of description,  $\mathcal{G}'$  is a copy of input graph  $\mathcal{G}$ .
- 2: Initialize two empty sequences  $s_1 \leftarrow \emptyset$ ,  $s_2 \leftarrow \emptyset$ , and a backward edge set  $b \leftarrow \emptyset$ .
- 3: **while**  $\mathcal{G} \neq \emptyset$  **do**
- 4:   **if**  $\mathcal{G}$  contains a source **then**
- 5:     choose the sink  $u$  with maximum  $\delta(U)$
- 6:      $s_1 \leftarrow s_1 // U$
- 7:      $\mathbf{V} \leftarrow \mathbf{V} \setminus U$ ;  $\mathbf{E} \leftarrow \mathbf{E} \setminus \{E_{UV}, E_{VU}\}, \forall U \in \mathbf{V}$
- 8:      $\mathcal{G} = (\mathbf{V}, \mathbf{E})$
- 9:   **end if**
- 10:   **if**  $\mathcal{G}$  contains a sink **then**
- 11:     choose the sink  $U$  with minimum  $\delta(U)$
- 12:      $s_2 \leftarrow U // s_2$
- 13:      $\mathbf{V} \leftarrow \mathbf{V} \setminus U$ ;  $\mathbf{E} \leftarrow \mathbf{E} \setminus \{E_{UV}, E_{VU}\}, \forall U \in \mathbf{V}$
- 14:      $\mathcal{G} = (\mathbf{V}, \mathbf{E})$
- 15:   **end if**
- 16: **end while**
- 17: The topological ordering is  $s = s_1 // s_2$
- 18: **for**  $E_{UV}$  in the input graph  $\mathcal{G}'$  **do**
- 19:   **if**  $U$  is after  $V$  in  $s$  **then**
- 20:      $b \leftarrow b // E_{UV}$
- 21:   **end if**
- 22: **end for**
- 23: Sort  $b$  in ascending order according to  $w_{UV}$
- 24: **while**  $\mathcal{G}'$  is not a DAG **do**
- 25:   remove the edge with smallest  $w_{UV}$  from  $\mathcal{G}'$
- 26: **end while**
- 27: **return** The directed acyclic causal graph  $\mathcal{G}'$ .

---

We begin by examining the relationship between the weighting parameter  $\lambda$ , the empirical support rate, and the effective threshold. We then establish a general bound on the probability of edge-level error under the weighted voting rule, and provide sufficient conditions under specific support distributions that guarantee accurate recovery.

#### D.1 BAYESIAN MOTIVATION FOR THE WEIGHTED VOTING RULE

Specifically, we show that the score can be viewed as the posterior mean under a Beta prior whose influence diminishes as the number of supporting subgraphs increases. We first consider each edge direction  $X \rightarrow Y$  as a binary decision problem. Suppose each local subgraph that includes both  $X$  and  $Y$  independently votes for one of the two directions:  $X \rightarrow Y$  or  $Y \rightarrow X$ . Let  $A$  and  $B$  denote the number of times each direction appears, and let  $m = A + B$  be the total number of subgraphs providing directional evidence.

A natural approach is to model the true support probability  $p = \Pr(X \rightarrow Y)$  using a Beta prior:

$$p \sim \text{Beta}(\alpha, \beta),$$

so that the posterior mean becomes

$$\mathbb{E}[p \mid A] = \frac{A + \alpha}{m + \alpha + \beta}. \quad (7)$$

In classical Laplace smoothing, a fixed prior such as  $\text{Beta}(1, 1)$  adds uniform pseudo-counts regardless of sample size. However, in our setting, most candidate edges are supported by very few subgraphs. The fixed priors are therefore either too weak to suppress noise or too strong to allow learning when evidence grows.

We therefore introduce a data-dependent pseudo-count that decreases with  $m$ . Specifically, we set  $\alpha = 0$ , and define an effective prior strength:

$$\beta = \kappa(m) := \frac{me^{-\lambda m}}{1 - e^{-\lambda m}}, \quad (8)$$

where  $\lambda > 0$  is a tunable parameter. This yields the posterior mean:

$$\mathbb{E}[p \mid A] = \frac{A}{m + \kappa(m)} = (1 - e^{-\lambda m}) \cdot \frac{A}{m}. \quad (9)$$

Thus, our weighted score function  $s(X \rightarrow Y)$  can be viewed as the posterior mean under a Beta prior whose strength vanishes exponentially as the number of supporting subgraphs increases. When  $m$  is small, the exponential decay is slow, and the prior contributes a significant regularization, effectively suppressing low-support edges. As  $m$  grows, the prior influence rapidly vanishes, and the score approaches the empirical frequency  $A/m$ , recovering naive voting.

The hyperparameter  $\lambda$  controls how quickly the prior decays. A larger  $\lambda$  yields more aggressive penalization for rare edges, while a smaller  $\lambda$  allows quicker adaptation to the empirical signal. This dynamic pseudo-count interpretation explains the design of our exponential weight  $1 - e^{-\lambda m}$  and its effectiveness in controlling false positives in sparse and noisy settings.

## D.2 PROOF OF THEOREM 3.2

**Theorem 3.2** (Sufficient Condition for Weighted Voting Accuracy) *Let  $A \sim \text{Binomial}(m, p)$  represent the number of successful votes in  $m$  independent subgraphs for the edge direction  $X_1 \rightarrow X_2$ , where each subgraph supports this direction independently with probability  $p \in (0, 1)$ , decision threshold  $t \in (0, 1)$  and the weight function  $w(m) = 1 - e^{-\lambda m}$ ,  $\lambda > 0$ . Assume the effective threshold for accept the edge direction  $X_1 \rightarrow X_2$  is  $r(m) = \frac{t}{1 - e^{-\lambda m}} < p$ , i.e., the true support rate  $p$  is above the effective threshold. Then, if*

$$\frac{mp}{2} \left( 1 - \frac{t}{p(1 - e^{-\lambda m})} \right)^2 \geq \log \frac{1}{\epsilon},$$

it follows that  $P(s(A) \geq t) \geq 1 - \epsilon$ .

*Proof.* Our goal is to show that

$$P\left(s = [1 - \exp(-\lambda m)] \cdot \frac{A}{m} \geq t\right) \geq 1 - \epsilon, \quad (10)$$

where  $A \sim \text{Binomial}(m, p)$  and  $m$  is the number of (independent) subgraphs or subsamples considered. Rewriting  $s \geq t$  gives

$$[1 - \exp(-\lambda m)] \cdot \frac{A}{m} \geq t \iff \frac{A}{m} \geq \frac{t}{1 - \exp(-\lambda m)}.$$

For notational simplicity, we define  $r = \frac{t}{1 - \exp(-\lambda m)}$ . Hence, our goal becomes ensuring  $P(A \geq mr) \geq 1 - \epsilon$ . Since  $A$  is a binomial random variable  $A \sim \text{Binomial}(m, p)$ ,  $\mathbb{E}[A] = mp$ , we therefore have the Chernoff bound, states that,

$$P\left(A \leq (1 - \delta)mp\right) \leq \exp\left(-\frac{\delta^2}{2}mp\right), \quad (11)$$

for any  $0 < \delta < 1$ . Subsequently, we set  $(1 - \delta)mp = mr$ , i.e.,  $\delta = 1 - \frac{r}{p}$ . Note that for this  $\delta$  to be positive (so that the Chernoff bound form applies), we need  $r < p$ . In other words,

$$\frac{t}{1 - \exp(-\lambda m)} = r < p,$$

which is the intuitive condition that the true probability  $p$  exceeds the effective threshold  $r$ . With the above definition of  $\delta$ ,

$$P\left(A < mr\right) = P\left(A \leq (1 - \delta)mp\right) \leq \exp\left(-\frac{mp}{2}\left(1 - \frac{r}{p}\right)^2\right).$$

Hence

$$P\left(A \geq mr\right) \geq 1 - \exp\left(-\frac{mp}{2}\left(1 - \frac{r}{p}\right)^2\right). \quad (12)$$

To ensure this probability is at least  $1 - \epsilon$ , we impose

$$\exp\left(-\frac{mp}{2}\left(1 - \frac{r}{p}\right)^2\right) \leq \epsilon.$$

Since  $r = \frac{t}{1 - e^{-\lambda m}}$ , this condition explicitly becomes

$$\frac{mp}{2}\left(1 - \frac{t}{p(1 - e^{-\lambda m})}\right)^2 \geq \log \frac{1}{\epsilon}. \quad (13)$$

Therefore, whenever (3) and  $r < p$  is satisfied, we have

$$P\left(\left[1 - e^{-\lambda m}\right] \cdot \frac{A}{m} \geq t\right) = P\left(A \geq mr\right) \geq 1 - \epsilon. \quad (14)$$

Hence the theorem follows.  $\square$

### D.3 PROOF OF COROLLARY 3.3

**Corollary 3.3** (Upper bound of node in subgraphs) *Let  $\lambda > 0$ ,  $t \in (0, 1)$ , and  $\epsilon \in (0, 1)$  be fixed. For a candidate edge  $(X, Y)$ , denote by  $m$  the number of local subgraphs whose Markov Blankets contain both endpoints. Under the setting of Theorem 3.2, the sufficient condition (3) can be converted into the explicit bound*

$$m \geq \frac{2 \log(1/\epsilon)}{p((1 - t/p)^2 - 2(t/p)(1 - t/p)e^{-\lambda})},$$

*Proof.* We first define  $y = \exp(-\lambda m)$ . Then, by the conclusion of Theorem 4.4, we obtain

$$-\frac{p}{2\lambda} \log y \left(1 - \frac{t}{p(1 - y)}\right)^2 \geq \log \frac{1}{\epsilon}. \quad (15)$$

Next, we consider the first-order Taylor expansion:

$$\begin{aligned} \left(1 - \frac{t}{p} \frac{1}{1 - y}\right)^2 &= \left[1 - \frac{t}{p} - \frac{t}{p}y + O(y^2)\right]^2 \\ &= [\gamma - \theta y + O(y^2)]^2 \\ &= \gamma^2 - 2\theta\gamma y + O(y^2), \end{aligned} \quad (16)$$

where we set  $\theta = \frac{t}{p}$  and  $\gamma = 1 - \frac{t}{p}$ . Therefore, (15) becomes

$$\log y [\gamma^2 - 2\theta\gamma y + O(y^2)] \leq -\frac{2\lambda}{p} \log \frac{1}{\epsilon}. \quad (17)$$

Therefore, by substituting  $\log y = -\lambda m$  and dropping the  $O(y^2)$  term (since  $y^2$  is small enough), we get an approximate condition:

$$m(\gamma^2 - 2\theta\gamma e^{-\lambda m}) \geq \frac{2}{p} \log \frac{1}{\epsilon}. \quad (18)$$

This is an implicit condition on  $m$ . To derive an explicit and sufficient lower bound, we strengthen the left-hand side. Since  $m \geq 1$ , we have  $e^{-\lambda m} \leq e^{-\lambda}$ , therefore

$$\gamma^2 - 2\theta\gamma e^{-\lambda m} \geq \gamma^2 - 2\theta\gamma e^{-\lambda}.$$

Let  $K_\lambda = \gamma^2 - 2\theta\gamma e^{-\lambda}$ . To ensure the lower bound is positive, we require  $\gamma^2 > 2\theta\gamma e^{-\lambda}$ , or equivalently  $\gamma > 2\theta e^{-\lambda}$  (since  $\gamma = 1 - t/p > 0$ ). This condition simplifies to  $1 - t/p > 2(t/p)e^{-\lambda}$ .

By our Theorem 3.4, since  $\lambda > \frac{1}{m} \log(1-t)$ , we have

$$\frac{t}{p} (2e^{-\lambda} + 1) < \frac{t}{p} \left( 2(1-t)^{\frac{1}{m}} + 1 \right) < 1.$$

This inequality is easily satisfied and thus represents a very mild condition. Consequently, we can regard the lower bound  $K_\lambda > 0$  as being established. Then the inequality (18) is satisfied under a stronger condition:

$$m \cdot K_\lambda \geq \frac{2}{p} \log \frac{1}{\epsilon}.$$

Solving for  $m$  gives an explicit lower bound:

$$m \geq \frac{2 \log(1/\epsilon)}{p K_\lambda} = \frac{2 \log(1/\epsilon)}{p(\gamma^2 - 2\theta\gamma e^{-\lambda})}.$$

Substituting the definitions of  $\gamma$  and  $\theta$ , we obtain:

$$m \geq \frac{2 \log(1/\epsilon)}{p((1-t/p)^2 - 2(t/p)(1-t/p)e^{-\lambda})}.$$

□

## E DISCUSSION OF THE STRUCTURE-AWARE ERROR BOUND

The weighted voting procedure serves as the core mechanism for aggregating local subgraph estimates into a global DAG. While this method adjusts edge confidence based on empirical support, its effectiveness ultimately depends on the ability to balance false positives and false negatives across the merged graph. To better understand this behavior, we analyze the global error induced by the weighted voting rule and how it interacts with the sparsity of the graph, the choice of voting threshold, and the distribution of subgraph overlaps.

This section formalizes that analysis. We first derive a decomposition of the total error into false positive and false negative components, followed by a structure-aware upper bound based on the union bound. The role of the weighting parameter  $\lambda$  is then examined in detail, culminating in formal proofs of Theorem 3.4 and Theorem 3.5, which establish a feasible range for  $\lambda$  and the asymptotic vanishing of global error, respectively. These bounds are further instantiated under Erdős-Rényi (ER) and scale-free (SF) graph models to characterize how graph topology influences the merging accuracy.

To begin with, we formalize the decomposition of the global error into *false negatives* (FN) and *false positives* (FP), and derive a structure-aware upper bound based on the union bound. We summarized it into the following lemma:

**Lemma E.1** (Structure-aware global error bound). *Each candidate directed edge  $(V_i, V_j)$  is evaluated in  $m_{ij}$  independent local sub-graphs whose Markov Blankets contain both endpoints.*

- For a **true** edge, the vote count obeys  $A_{ij} \sim \text{Binomial}(m_{ij}, p)$ .
- For a **false** edge,  $A_{ij} \sim \text{Binomial}(m_{ij}, q)$  with  $p > q$ .

Using the weighted rule

$$s_{ij} = [1 - e^{-\lambda m_{ij}}] \frac{A_{ij}}{m_{ij}} \geq t, \quad r_\lambda(m_{ij}) = \frac{t}{1 - e^{-\lambda m_{ij}}},$$

assume  $p > r_\lambda(m_{ij})$  and  $q < r_\lambda(m_{ij})$  for every edge. Then

$$\Pr(\text{global error}) \leq \underbrace{\sum_{(i,j) \in E^*} e^{-2m_{ij}[p-r_\lambda(m_{ij})]^2}}_{\text{FN contribution}} + \underbrace{\sum_{(i,j) \notin E^*} e^{-2m_{ij}[r_\lambda(m_{ij})-q]^2}}_{\text{FP contribution}}, \quad (19)$$

where  $E^*$  denotes the ground-truth edge set.

*Proof.* For  $(V_i, V_j) \in \mathbf{E}^*$ , we have

$$\Pr(\text{FN on } (V_i, V_j)) = \Pr(A_{ij}/m_{ij} < r_\lambda(m_{ij})) \leq e^{-2m_{ij}[p-r_\lambda(m_{ij})]^2} \quad (20)$$

by Hoeffding's inequality. A symmetric argument gives the FP term for  $(V_i, V_j) \notin \mathbf{E}^*$ . Finally, the union bound over all edges yields the claimed inequality.  $\square$

**Corollary E.2** (Worst-case simplification). *If  $m_{ij} \geq m_{\min}$  for all edges, then*

$$\Pr(\text{global error}) \leq N_{\text{FN}} e^{-2m_{\min}[p-r_\lambda(m_{\min})]^2} + N_{\text{FP}} e^{-2m_{\min}[r_\lambda(m_{\min})-q]^2}, \quad (21)$$

where  $N_{\text{FN}} = |\mathbf{E}^*|$  and  $N_{\text{FP}} = \binom{n}{2} - N_{\text{FN}}$  for a graph with  $n$  nodes.

The error bound derived above depends on the effective threshold  $r_\lambda(m)$ , which is controlled by the weighting parameter  $\lambda$ . To understand the role of this parameter, it is instructive to consider the limiting case  $\lambda = 0$ , which corresponds to the naive voting scheme. In this case, the weight term disappears, and the edge inclusion rule reduces to comparing the raw directional frequency  $A/m$  against the fixed threshold  $t$ .

**Remark E.3** (Naive voting baseline). If we drop the weight and decide solely on the unweighted fraction  $\frac{A_{ij}}{m_{ij}} \geq t$ , Lemma E.1 specialises to

$$\Pr(\text{global error}) \leq \sum_{(i,j) \in \mathbf{E}^*} e^{-2m_{ij}(p-t)^2} + \sum_{(i,j) \notin \mathbf{E}^*} e^{-2m_{ij}(t-q)^2}. \quad (22)$$

In sparse graphs, where the number of candidate false positive edges vastly exceeds the number of true positives (i.e.,  $N_{\text{FP}} \gg N_{\text{FN}}$ ), the overall error is typically dominated by the first summation term. Therefore, a moderate increase in  $\lambda$  can lead to a significant reduction in total error by aggressively penalizing low-support spurious edges, even if it slightly increases the false negative rate. This trade-off is particularly favorable in high-dimensional settings, where controlling the false discovery rate is often more critical than maximizing recall. These insights align with the empirical results reported in Section 4.1, where the weighted voting scheme consistently improves FDR without severely compromising TPR across a wide range of base learners.

## E.1 INFLUENCE AND PRACTICAL RANGE OF THE WEIGHT PARAMETER $\lambda$

To ensure that the weighted voting mechanism achieves a reliable trade-off between false positives and false negatives, it is necessary to understand how the choice of the weighting parameter  $\lambda$  affects the acceptance threshold and the overall error bound. The following derivation provides a characterization of the feasible range of  $\lambda$  that satisfies the conditions used in the theoretical analysis of edge decisions. This directly supports the proof of Theorem 3.4 in the main text.

**Theorem 3.4** (Practical choice of  $\lambda$ ) *Fix a vote count  $m \geq 1$ , a decision threshold  $t \in (0, 1)$ , and a target error level  $\epsilon \in (0, 1)$ . If  $\lambda$  satisfies*

$$-\frac{1}{m} \ln(1-t) < \lambda \leq -\frac{1}{m} \ln \epsilon,$$

*then the weighted-vote rule achieves the prescribed error control under the union bound.*

*Proof.* Define the Hoeffding-based global error upper bound  $\mathcal{L}(\lambda) = N_{\text{FN}} e^{-2m_{\min}(p-r_\lambda)^2} + N_{\text{FP}} e^{-2m_{\min}(r_\lambda-q)^2}$ , where  $N_{\text{FN}}$  ( $N_{\text{FP}}$ ) is the number of true (false) candidate edges rescaled by their respective cost coefficients. For notational simplicity, we omit the subscripts, and use  $m$  to represent  $m_{\min}$  in our later proof. We first differentiate  $\mathcal{L}$  w.r.t.  $\lambda$ :

$$\begin{aligned} \frac{\partial r_\lambda}{\partial \lambda} &= \frac{tme^{-\lambda m}}{(1-e^{-\lambda m})^2} = r_\lambda \frac{me^{-\lambda m}}{1-e^{-\lambda m}} > 0, \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= 2m \frac{te^{-\lambda m}}{(1-e^{-\lambda m})^2} \left[ N_{\text{FN}} \delta_p e^{-2m\delta_p^2} - N_{\text{FP}} \delta_q e^{-2m\delta_q^2} \right], \end{aligned} \quad (23)$$

with  $\delta_p = p - r_\lambda > 0$  and  $\delta_q = r_\lambda - q > 0$ .



Because  $\delta_p$  increases and  $\delta_q$  decreases as  $\lambda$  grows, a larger  $\lambda$  lowers the false-negative term (higher **recall**) but raises the false-positive term (lower **precision**). For sparse causal graphs we typically have  $N_{FP} \gg N_{FN}$ , making the second term dominant and hence  $\partial\mathcal{L}/\partial\lambda < 0$  until the exponential weight saturates. Consequently, increasing  $\lambda$  is beneficial *only* inside a finite interval.

**Upper bound for  $\lambda$ .** In the worst-case scenario where all candidate edges are consistently supported in the same direction, the voting scores for both true and false edges become uniformly close to  $1 - e^{-\lambda m}$ . If  $1 - e^{-\lambda m} \geq 1 - \epsilon$  ( $0 < \epsilon \ll 1$ ), then even false edges can exceed the decision threshold  $t$ , leading to a large number of number of false positives. Therefore To avoid such indiscriminate acceptance,  $\lambda$  must be chosen to ensure that  $1 - e^{-\lambda m}$  remains sufficiently below 1. Solving  $e^{-\lambda m} = \epsilon$  gives

$$0 < \lambda \leq \lambda_{\max}(\epsilon) = -\frac{1}{m} \ln \epsilon \quad (\text{e.g., } \epsilon = 0.01 \Rightarrow \lambda_{\max} \approx 4.6/m). \quad (24)$$

**Lower bound for  $\lambda$ .** The effective threshold  $r_\lambda(m) = \frac{t}{1 - e^{-\lambda m}}$  must satisfy  $0 < r_\lambda < 1$ ; otherwise the acceptance condition  $A/m \geq r_\lambda$  can never be met because  $A/m \leq 1$  by definition. Solving the inequality  $r_\lambda < 1$  yields

$$\frac{t}{1 - e^{-\lambda m}} < 1 \iff e^{-\lambda m} < 1 - t \iff \lambda > \lambda_{\min} := -\frac{1}{m} \ln(1 - t). \quad (25)$$

Intuitively, when  $\lambda$  falls below this bound the exponential weight is so close to 1 that the prefactor  $1 - e^{-\lambda m}$  becomes *smaller* than  $t$ , inflating  $r_\lambda$  beyond 1 and blocking every candidate edge, including true ones. Hence  $\lambda_{\min}$  is the *viability threshold*: only for  $\lambda > \lambda_{\min}$  does the weighted-voting rule retain a non-zero recall. Therefore a practical search range is

$$\lambda \in \left[ -\frac{1}{m} \ln(1 - t), \frac{1}{m} \ln \epsilon \right], \quad (26)$$

within which cross-validation or the closed-form condition  $\partial\mathcal{L}/\partial\lambda = 0$  can be used to pinpoint an optimal  $\lambda^*$ .  $\square$

**Exponentially vanishing reversal error.** For any  $\lambda$  in this range and any true edge with support probability  $p > r_\lambda$ , the probability of being accepted in the reverse direction is  $\Pr(\text{reverse}) \leq \exp[-2m(p - r_\lambda)^2]$ , which decays exponentially with the number of independent subgraphs  $m$ . This guarantees that the weighted-voting merger remains statistically consistent as data grow, while a properly chosen  $\lambda$  suppresses spurious edges in finite-sample regimes.

The general error bound depends on the number of subgraphs in which each edge appears. This quantity is influenced by the underlying graph topology. In the following, the behavior of the bound is examined under two commonly used random graph models: Erdős–Rényi and scale-free graphs. The analysis characterizes typical support counts and their implications for the error terms derived in Lemma E.1.

## E.2 ERDŐS–RÉNYI AND SCALE-FREE GRAPHS

The error bounds derived in the previous section depend not only on the weighting parameter  $\lambda$ , but also on the empirical support count  $m_{ij}$  which measures the number of subgraphs in which each edge appear. This quantity is influenced by the underlying graph topology and the statistical properties of the Markov Blanket construction.

To understand how  $m_{ij}$  behaves in practice, we analyze two representative random graph models: Erdős–Rényi (ER) and scale-free (SF) networks. These models differ significantly in their degree distributions, which in turn affect the overlap patterns among Markov Blankets and the expected frequency with which edges are covered by local subgraphs. The analysis below characterizes typical support rates under each model, providing context for interpreting the global error bounds and informing the expected sample complexity of reliable aggregation.

**Theorem E.4** (ER- $h$  graph). *Let  $G \sim ER(n, \theta)$  with edge probability  $\theta = h/(n-1)$ , and assign directions by a random topological order so that the expected out-degree is  $h$ . Denote*

$$\delta_p := p - r_\lambda(2), \quad \delta_q := r_\lambda(2) - q \quad (\delta_p, \delta_q > 0).$$

*Then, with probability at least  $1 - O(\theta^2)$  over the graph draw,*

$$\Pr(\text{global error}) \leq \frac{nh}{2} e^{-4\delta_p^2} + \frac{n(n-1) - nh}{2} e^{-4\delta_q^2} + O(\theta^2). \quad (27)$$

*The  $O(\theta^2)$  term covers the negligible fraction of edges whose vote count  $m_{ij} > 2$ .*

*Proof.* In a directed ER graph each vertex has  $\deg^{\text{in}}, \deg^{\text{out}} \sim \text{Pois}(\theta/2)$ , so  $\mathbb{E}[|\text{MB}(v)|] = \mathbb{E}[\deg^{\text{in}} + \deg^{\text{out}} + \text{spouses}] \approx 2h$ , where the “spouse” term (*co-parents*) shares the same mean as  $\deg^{\text{out}}$ .

For an oriented edge  $(i, j)$ , it appears in *both*  $\text{MB}(V_i)$  and  $\text{MB}(V_j)$ , giving a baseline  $m_{ij} \geq 2$ . Additionally, it appears in  $\text{MB}(V_k)$  for every common child  $V_k$  of  $V_i$  and  $V_j$ . For fixed  $V_k$ , the events “ $V_i \rightarrow V_k$ ” and “ $V_j \rightarrow V_k$ ” are independent with probability  $\theta^2$ . Hence the number of common children follows  $\text{Pois}(\lambda_c)$  with  $\lambda_c = (n-2)\theta^2 \approx h^2/n$ .

Thus,

$$m_{ij} = 2 + X, \quad X \sim \text{Pois}(\lambda_c).$$

When  $h = O(1)$ ,  $\lambda_c = O(\theta^2) \ll 1$ , whence

$$\Pr(m_{ij} = 2) = 1 - O(\theta^2), \quad \Pr(m_{ij} \geq 3) = O(\theta^2).$$

For the overwhelming majority of edges ( $m_{ij} = 2$ ), lemma E.1 gives:

$$\Pr(\text{FN on } (i, j)) \leq e^{-4\delta_p^2}, \quad \Pr(\text{FP on } (i, j)) \leq e^{-4\delta_q^2}.$$

Counting edges:

$$N_{\text{FN}} \approx \frac{nh}{2}, \quad N_{\text{FP}} = \binom{n}{2} - N_{\text{FN}}.$$

Summing the two contributions yields

$$\frac{nh}{2} e^{-4\delta_p^2} + \frac{n(n-1) - nh}{2} e^{-4\delta_q^2}. \quad (28)$$

□

We can obtain similar results from the SF graph.

**Theorem E.5** (SF- $h$  graph). *Let  $\mathcal{G}$  be a directed scale-free graph on  $n$  nodes, obtained by sampling an undirected Chung–Lu (or Barabási–Albert) graph whose degree sequence  $(d_1, \dots, d_n)$  satisfies*

$$\Pr(d \geq k) \leq C_\alpha k^{1-\alpha}, \quad 2 < \alpha < 3, \quad (29)$$

*and whose mean degree is  $h$ ; and orienting edges according to a random topological order. Then, for a universal constant  $C_\alpha$  that depends only on  $\alpha$ ,*

$$\Pr(\text{global error}) \leq \frac{nh}{2} e^{-4\delta_p^2} + \frac{n(n-1)}{2} - \frac{nh}{2} e^{-4\delta_q^2} + \frac{n(n-1)}{2} C_\alpha n^{-(\alpha-2)}. \quad (30)$$

*Proof.* For an oriented edge  $(V_i, V_j)$  let  $d_i, d_j$  be its endpoint degrees. Exactly as in the ER case each edge appears at least twice; additional occurrences come from every *common child*  $k$  with probability  $(d_i/n)(d_j/n)$ . Hence

$$m_{ij} = 2 + X, \quad X \sim \text{Pois}(\lambda_{ij}), \quad \lambda_{ij} := \frac{d_i d_j}{n}.$$

For any fixed  $\lambda$  and  $\delta \in \{\delta_p, \delta_q\}$

$$\begin{aligned}\mathbb{E}[e^{-2m_{ij}\delta^2} \mid \lambda_{ij}] &= e^{-4\lambda_{ij}^2} \mathbb{E}[e^{-2X\delta^2}] \\ &= e^{-4\delta^2} e^{\lambda_{ij}(e^{-2\delta^2}-1)}.\end{aligned}\tag{31}$$

Since the value of  $e^{-2\delta^2} - 1$  varies, we splitted the expectation (31) into two regimes:

$$\mathbb{E}[e^{-2m_{ij}\delta^2}] = \mathbb{E}[e^{-2m_{ij}\delta^2} \mathbf{1}_{\{\lambda_{ij} \leq 1\}}] + \mathbb{E}[e^{-2m_{ij}\delta^2} \mathbf{1}_{\{\lambda_{ij} > 1\}}].\tag{32}$$

- Non-hub regime  $\lambda_{ij} \leq 1$ :

$$\mathbb{E}[e^{-2m_{ij}\delta^2} \mid \lambda_{ij}] \leq e^{-4\delta^2}.\tag{33}$$

- Hub regime  $\lambda_{ij} > 1$ . By (31) the conditional term is  $\leq e^{-4\delta^2} e^{-\lambda_{ij}/2} \leq 1$ , but the probability of this event can be bounded with the degree tail:

$$\mathbb{E}[e^{-2m_{ij}\delta^2} \mathbf{1}_{\{\lambda_{ij} > 1\}}] \leq \Pr(\lambda_{ij} > 1) = \Pr\left(\frac{d_i d_j}{n} > 1\right) \leq C_\alpha n^{-(\alpha-2)},$$

where we apply the union bound to decompose the event  $d_i d_j > n$  into two simpler events,  $d_i > n^{1/2}$  or  $d_j > n^{1/2}$ , control each using the degree tail bound  $\Pr(d \geq k) \leq C_\alpha k^{1-\alpha}$ , and then combine the two estimates.

Therefore, for either  $\delta = \delta_p$  or  $\delta_q$ ,

$$\mathbb{E}[e^{-2m_{ij}\delta^2}] \leq e^{-4\delta^2} + C_\alpha n^{-(\alpha-2)}.\tag{34}$$

There are  $N_{\text{FN}} \approx nh/2$  true and  $N_{\text{FP}} = \binom{n}{2} - N_{\text{FN}}$  false edges on average. Multiplying the expectation (34) by these counts and plugging into Lemma E.1 yields inequality (30).  $\square$

Theorem E.5 completes the structure-aware error analysis by characterizing the influence of heterogeneous degree distributions on the residual error bound. While the dominant exponential terms governing false positive and false negative rates are structurally similar to those in Theorem E.4, the residual term exhibits a slower decay due to the presence of high-degree nodes. These hub-related structures lead to greater variability in the support count  $m_{ij}$  across candidate edges.

This variability has practical implications. In networks where edge supports are highly non-uniform, the weighted voting mechanism implicitly induces a form of confidence calibration: high-support edges, typically associated with structurally central nodes, retain larger weights and are more likely to be preserved. In contrast, low-support edges often arising from sparse or weakly connected regions, will be heavily penalized by the exponential weighting term. This differential treatment improves robustness to statistical noise and helps suppress false positives without uniformly raising the threshold for all decisions.

As a result, the error reduction effect of the weighting scheme is not solely determined by the average support level, but also by the variance in subgraph overlap. Networks with broader support distributions provide more opportunities for selective edge retention, which enhances the overall effectiveness of the aggregation procedure. This observation complements the earlier asymptotic result, and offers a finer-grained explanation of the empirical precision gains observed in our experiments.

To complete the analysis, we examine how the global error behaves asymptotically under increasing graph size.

### E.3 ASYMPTOTIC ANALYSIS

**Theorem 3.5** (Asymptotic Consistency) *Fix a threshold  $t \in (0, 1)$  and let  $\delta_p = p - t$  and  $\delta_q = t - q$  denote the positive margins between  $t$  and the inclusion probabilities  $p, q$  of true and false edges*

respectively. Assume  $\delta_p, \delta_q > 0$  and that  $\lambda$  satisfies the conditions in Theorem 3.4. If the number of local subgraphs per candidate edge is  $m = C \log n$  with  $C > \frac{2}{\min\{\delta_p^2, \delta_q^2\}}$ , then we have

$$\Pr(\text{global error}) = o(1), \quad \text{as } n \rightarrow \infty. \quad (35)$$

*Proof.* By the conclusion of Lemma E.1, the global error probability is bounded by

$$\Pr(\text{global error}) \leq \sum_{(i,j) \in E^*} e^{-2m_{ij}(p-t)^2} + \sum_{(i,j) \notin E^*} e^{-2m_{ij}(t-q)^2}. \quad (36)$$

Since the number of true edges satisfies  $N_{\text{FN}} = |\mathbf{E}^*| = \mathcal{O}(n)$ , and the number of false edges is  $N_{\text{FP}} = \binom{n}{2} - N_{\text{FN}} = \mathcal{O}(n^2)$ , we can simplify the above bound by letting  $m_{ij} \equiv m$  for all edges:

$$\Pr(\text{global error}) \leq N_{\text{FN}} e^{-2m\delta_p^2} + N_{\text{FP}} e^{-2m\delta_q^2},$$

where we denote  $\delta_p = p - t > 0$  and  $\delta_q = t - q > 0$ .

To ensure that both terms remain bounded by a constant, we require

$$e^{-2m\delta_p^2} \leq n^{-1} \quad \Rightarrow \quad m \geq \frac{1}{2\delta_p^2} \log n,$$

and

$$e^{-2m\delta_q^2} \leq n^{-2} \quad \Rightarrow \quad m \geq \frac{1}{\delta_q^2} \log n.$$

Therefore, it suffices to set

$$m = C \log n, \quad C > \max \left\{ \frac{1}{2\delta_p^2}, \frac{1}{\delta_q^2} \right\},$$

which guarantees that

$$\Pr(\text{global error}) \leq \underbrace{\mathcal{O}(n \cdot n^{-1})}_{=\mathcal{O}(1)} + \underbrace{\mathcal{O}(n^2 \cdot n^{-2})}_{=\mathcal{O}(1)} = \mathcal{O}(1).$$

In fact, choosing a slightly larger constant  $C$  makes both terms decay to zero, which establishes asymptotic consistency as  $n \rightarrow \infty$ .  $\square$

**Complexity.** Finally, we analyze the computational complexity, which consists of two parts:

- The local structure learning phase takes  $\mathcal{O}(m^3)$  per node, and there are  $n$  nodes, resulting in  $\mathcal{O}(nm^3)$  total cost.
- The voting and merging phase requires computing pairwise edge counts and resolving cycles over  $\mathcal{O}(n^2)$  edge pairs, leading to an additional  $\mathcal{O}(n^2)$  term.

Substituting  $m = \mathcal{O}(\log n)$ , the total runtime becomes

$$\mathcal{O}(n(\log n)^3 + n^2) = \tilde{\mathcal{O}}(n^2),$$

where the soft- $\mathcal{O}$  notation hides polylogarithmic factors. Thus, the proposed divide-and-conquer method achieves both statistical consistency and near-quadratic scalability.

## F IMPLEMENTATION DETAILS

Our code is based on two open-source packages: `gcastle`, which provides implementations of score-based and continuous causal discovery methods such as NOTEARS, GOLEM, GraN-DAG and DAG-GNN, and `dodiscover`, which implements ordering-based methods. These packages form the backbone of our experimental framework. On top of them, we implement our own modules for subgraph construction, weighted voting aggregation, and cycle removal. The full pipeline with configuration scripts and reproducibility controls is described in detail below. Subsequent subsections provide additional implementation details for baseline configuration, extended experimental results, runtime breakdown, and comparison against DCILP.

## F.1 BASELINES

All baseline methods are implemented using publicly available code and configured with recommended hyperparameters. For methods involving continuous optimization, the primary computational bottleneck lies in gradient-based acyclicity constraints, which require  $\mathcal{O}(d^3)$  time and  $\mathcal{O}(d^2)$  memory due to matrix operations over the full graph. Discrete search-based methods such as SCORE and CAM incur combinatorial overhead when handling larger node counts. In all cases, integrating these methods into the VISTA framework significantly reduces both runtime and memory usage, as the local subgraphs are orders of magnitude smaller and can be processed independently.

**NOTEARS** This method reformulates the combinatorial problem of DAG structure learning into a purely continuous optimization problem. It introduces a novel, smooth, and exact characterization of acyclicity using a matrix exponential function  $h(W) = \text{tr}(W \circ W) - d = 0$ . This transformation allows the problem to be solved efficiently using standard gradient-based optimization techniques, avoiding discrete search over graph structures.

**GOLEM** This work analyzes the role of sparsity and DAG constraints in learning linear DAGs, noting potential optimization issues with hard DAG constraints required by prior methods like NOTEARS. It proposes GOLEM (Gradient-based Optimization of dag-penalized Likelihood for learning linEar dag Models), which uses a likelihood-based score function instead of least squares. The key finding is that applying soft sparsity and DAG penalties to this likelihood objective suffices to recover the ground truth DAG structure asymptotically, resulting in an unconstrained optimization problem that is easier to solve.

**DAG-GNN** This method employs a deep generative model, specifically a Variational Autoencoder (VAE), to learn DAG structures, extending beyond linear models. It parameterizes the VAE’s encoder and decoder using a novel Graph Neural Network (GNN) architecture, designed to capture complex non-linear relationships inherent in data. The approach learns the graph’s weighted adjacency matrix alongside the neural network parameters, enforcing acyclicity through a continuous polynomial constraint, and naturally handles both continuous and discrete variables.

**GraN-DAG** This work extends continuous DAG learning to nonlinear settings by parameterizing each conditional distribution with neural networks and constructing a weighted adjacency matrix from network connectivity. Acyclicity is enforced through a smooth matrix-exponential constraint, enabling gradient-based optimization of the likelihood objective. Post-processing with thresholding and pruning helps recover sparse graphs.

**SCORE** This method recovers causal graphs for non-linear additive noise models by utilizing the score function ( $\nabla \log p(x)$ ) of the observational data distribution. It establishes that the Jacobian of the score function reveals information sufficient to identify leaf nodes in the causal graph. By iteratively identifying and removing leaves based on the variance of the score’s Jacobian diagonal elements, a topological ordering is estimated. The SCORE algorithm employs score matching techniques, specifically an extension of Stein’s identity, to approximate the necessary score Jacobian components from data samples.

**CAM** This approach estimates additive SEMs by decoupling the task into order search and edge selection. It first estimates a causal ordering of the variables using (potentially restricted) maximum likelihood, exploiting the identifiability property of additive models. Given the estimated order, sparse additive regression methods are then applied to select relevant parent variables (edges) for each node and estimate the corresponding additive functions. For high-dimensional data, an initial neighborhood selection step can reduce the search space before estimating the order.

In addition to the above baselines, we also include DCILP Dong et al. (2024), a recently proposed divide-and-conquer method that combines Markov Blanket estimation with global structure recovery via integer linear programming. While DCILP shares a similar high-level motivation with VISTA, it suffers from several practical limitations. Most notably, its final merging step relies on solving a large-scale ILP problem, which becomes computationally infeasible as the graph size increases. In many of our experimental settings, DCILP either fails to complete within a reasonable time window or produces no feasible solution at all. These issues highlight the need for a more lightweight and

scalable integration procedure, which motivates the design of VISTA. We provide a direct comparison with DCILP in the following section.

## F.2 COMPARISON WITH DCILP

We provide a detailed comparison between VISTA and DCILP, two methods that share a high-level divide-and-conquer strategy based on Markov Blanket decomposition. Although both approaches follow a similar decomposition principle, they differ notably in how they perform the aggregation step and enforce global acyclicity.

DCILP formulates the merging process as an integer linear program that guarantees the removal of 2-cycles, but relies on iterative post-processing to eliminate larger cycles. This procedure can be computationally intensive and may not always yield globally consistent solutions without additional refinement. In contrast, VISTA enforces acyclicity using a feedback arc set-based heuristic, which is algorithmically simpler and ensures a valid DAG by construction. Another distinction lies in how the two frameworks handle local estimation errors: DCILP applies aggressive pruning to Markov Blanket outputs before global optimization, which may propagate early-stage errors. VISTA instead retains a broader set of subgraph information and applies confidence-aware filtering during aggregation, providing more flexibility and robustness to local variability.

For empirical evaluation, we followed DCILP’s implementation baseline by using DAGMA Bello et al. (2022) as the phase-2 solver in both frameworks. This matched setup enables a controlled comparison under consistent base learners and subgraph configurations.

Table 5: Comparison of DCILP and VISTA under DAGMA baseline.

| Scenario      | Model    | FDR↓                              | TPR↑                              | SHD↓                                | F1↑                               |
|---------------|----------|-----------------------------------|-----------------------------------|-------------------------------------|-----------------------------------|
| ER5, $n = 30$ | DCILP    | $0.74 \pm 0.04$                   | $0.52 \pm 0.06$                   | $227.00 \pm 27.17$                  | $0.35 \pm 0.04$                   |
|               | VISTA-NV | $0.63 \pm 0.02$                   | <b><math>0.98 \pm 0.01</math></b> | $236.80 \pm 14.86$                  | $0.54 \pm 0.02$                   |
|               | VISTA-WV | <b><math>0.09 \pm 0.07</math></b> | $0.75 \pm 0.11$                   | <b><math>45.80 \pm 23.57</math></b> | <b><math>0.82 \pm 0.09</math></b> |
| SF5, $n = 30$ | DCILP    | $0.81 \pm 0.04$                   | $0.49 \pm 0.07$                   | $309.70 \pm 60.87$                  | $0.27 \pm 0.04$                   |
|               | VISTA-NV | $0.63 \pm 0.01$                   | <b><math>0.97 \pm 0.01</math></b> | $208.20 \pm 7.98$                   | $0.54 \pm 0.01$                   |
|               | VISTA-WV | <b><math>0.13 \pm 0.09</math></b> | $0.85 \pm 0.06$                   | <b><math>35.00 \pm 16.97</math></b> | <b><math>0.86 \pm 0.07</math></b> |
| ER3, $n = 50$ | DCILP    | $0.79 \pm 0.02$                   | $0.49 \pm 0.05$                   | $282.50 \pm 26.23$                  | $0.29 \pm 0.02$                   |
|               | VISTA-NV | $0.74 \pm 0.02$                   | <b><math>0.97 \pm 0.02</math></b> | $397.20 \pm 41.38$                  | $0.40 \pm 0.03$                   |
|               | VISTA-WV | <b><math>0.06 \pm 0.01</math></b> | $0.76 \pm 0.04$                   | <b><math>39.00 \pm 3.22</math></b>  | <b><math>0.84 \pm 0.02</math></b> |
| SF3, $n = 50$ | DCILP    | $0.91 \pm 0.01$                   | $0.52 \pm 0.03$                   | $820.40 \pm 110.23$                 | $0.15 \pm 0.01$                   |
|               | VISTA-NV | $0.71 \pm 0.05$                   | <b><math>0.95 \pm 0.02</math></b> | $345.00 \pm 71.05$                  | $0.44 \pm 0.05$                   |
|               | VISTA-WV | <b><math>0.14 \pm 0.04</math></b> | $0.84 \pm 0.06$                   | <b><math>40.80 \pm 12.38</math></b> | <b><math>0.81 \pm 0.08</math></b> |
| ER5, $n = 50$ | DCILP    | $0.80 \pm 0.01$                   | $0.52 \pm 0.04$                   | $520.20 \pm 29.51$                  | $0.29 \pm 0.01$                   |
|               | VISTA-NV | $0.76 \pm 0.01$                   | <b><math>0.98 \pm 0.01</math></b> | $730.80 \pm 24.85$                  | $0.38 \pm 0.01$                   |
|               | VISTA-WV | <b><math>0.09 \pm 0.03</math></b> | $0.83 \pm 0.03$                   | <b><math>59.20 \pm 10.32</math></b> | <b><math>0.86 \pm 0.02</math></b> |
| SF5, $n = 50$ | DCILP    | $0.90 \pm 0.01$                   | $0.49 \pm 0.03$                   | $1019.90 \pm 57.87$                 | $0.17 \pm 0.01$                   |
|               | VISTA-NV | $0.75 \pm 0.01$                   | <b><math>0.97 \pm 0.01</math></b> | $665.50 \pm 42.65$                  | $0.40 \pm 0.02$                   |
|               | VISTA-WV | <b><math>0.10 \pm 0.02</math></b> | $0.80 \pm 0.02$                   | <b><math>64.50 \pm 6.50</math></b>  | <b><math>0.85 \pm 0.02</math></b> |

Results in Table 5 show that, under the same configuration using DAGMA as the local structure learner, both VISTA variants (NV and WV) consistently outperform DCILP across all benchmark settings. Even the naive voting variant achieves lower FDR and SHD while maintaining competitive or higher TPR and F1 scores, suggesting that the ILP-based merging step in DCILP may introduce additional overhead without proportional accuracy gains. The weighted voting variant further improves performance by adaptively resolving directional conflicts based on edge support. We also note that as graph size increases such as  $n = 100$ , DCILP occasionally encounters solver infeasibility or produces solutions with substantially higher error rates, likely due to the combinatorial complexity of ILP formulation. In contrast, VISTA maintains stable performance with reduced computational demands. These comparisons underscore the scalability and robustness benefits of our framework in large-graph causal discovery settings.



### F.3 TIME COMPARISON

Since each local structure is learned independently based on a variable’s Markov Blanket, the entire divide phase can be executed in parallel across variables or computing nodes. This distributed strategy significantly reduces total runtime, especially when base learners are computationally intensive, such as neural network based models such as DAG-GNN and GraN-DAG or algorithms involving topological sorting such as SCORE.

Table 6, 7 and 8 confirm that VISTA consistently reduces the total execution time across a variety of settings. In large-scale graphs, where direct application of base methods may be computationally prohibitive, our framework provides a scalable alternative that decomposes the original problem into tractable subproblems. The integration step is lightweight and adds negligible overhead relative to the base learners. These results demonstrate that the benefits of VISTA are not limited to statistical performance, but also extend to practical runtime efficiency, enabling the application of complex causal discovery methods to larger and more realistic graphs.

Table 6: Comparison of total computing time (s) under ER5 setting.

| Method   | $n = 50$                              | $n = 100$                              | $n = 300$                              |
|----------|---------------------------------------|--|--|
| NOTEARS  | 510.73 $\pm$ 84.15                    | 2465.33 $\pm$ 58.02                    | 22407.77 $\pm$ 940.32                  |
| +VISTA   | <b>213.73 <math>\pm</math> 149.68</b> | <b>1096.51 <math>\pm</math> 142.87</b> | <b>3714.30 <math>\pm</math> 908.81</b> |
| GOLEM    | 76.16 $\pm$ 7.59                      | 115.01 $\pm$ 35.82                     | 276.80 $\pm$ 11.03                     |
| +VISTA   | <b>23.25 <math>\pm</math> 0.67</b>    | <b>37.57 <math>\pm</math> 1.36</b>     | <b>46.53 <math>\pm</math> 4.61</b>     |
| DAG-GNN  | 794.42 $\pm$ 72.61                    | 3137.68 $\pm$ 214.75                   | 29801.46 $\pm$ 1105.64                 |
| +VISTA   | <b>311.34 <math>\pm</math> 54.23</b>  | <b>818.52 <math>\pm</math> 501.88</b>  | <b>3313.86 <math>\pm</math> 945.29</b> |
| GraN-DAG | 919.26 $\pm$ 106.65                   | 5613.13 $\pm$ 1068.14                  | 25684.95 $\pm$ 2035.14                 |
| +VISTA   | <b>208.43 <math>\pm</math> 26.62</b>  | <b>934.72 <math>\pm</math> 50.64</b>   | <b>2851.04 <math>\pm</math> 376.84</b> |
| SCORE    | 629.88 $\pm$ 93.72                    | 15876.42 $\pm$ 807.89                  | —                                      |
| +VISTA   | <b>191.84 <math>\pm</math> 33.49</b>  | <b>479.60 <math>\pm</math> 38.19</b>   | <b>945.45 <math>\pm</math> 72.27</b>   |

Table 7: Comparison of total computing time (s) under SF3 setting.

| Method   | $n = 50$                              | $n = 100$                             | $n = 300$                               |
|----------|---------------------------------------|---------------------------------------|---|
| NOTEARS  | 713.30 $\pm$ 58.81                    | 2813.36 $\pm$ 804.27                  | 16631.62 $\pm$ 632.76                   |
| +VISTA   | <b>400.82 <math>\pm</math> 83.65</b>  | <b>652.59 <math>\pm</math> 57.99</b>  | <b>1714.09 <math>\pm</math> 237.32</b>  |
| GOLEM    | 100.73 $\pm$ 45.25                    | 169.20 $\pm$ 16.68                    | 398.58 $\pm$ 45.03                      |
| +VISTA   | <b>23.63 <math>\pm</math> 1.61</b>    | <b>35.47 <math>\pm</math> 2.26</b>    | <b>60.20 <math>\pm</math> 13.66</b>     |
| DAG-GNN  | 697.78 $\pm$ 93.37                    | 3555.95 $\pm$ 2050.91                 | 21242.03 $\pm$ 2178.95                  |
| +VISTA   | <b>282.70 <math>\pm</math> 297.75</b> | <b>645.77 <math>\pm</math> 308.88</b> | <b>2020.79 <math>\pm</math> 811.42</b>  |
| GraN-DAG | 890.96 $\pm$ 135.84                   | 4978.95 $\pm$ 656.25                  | 19372.84 $\pm$ 3037.94                  |
| +VISTA   | <b>319.29 <math>\pm</math> 389.88</b> | <b>817.63 <math>\pm</math> 145.63</b> | <b>2849.62 <math>\pm</math> 1558.40</b> |
| SCORE    | 495.12 $\pm$ 62.44                    | 18643.16 $\pm$ 970.22                 | —                                       |
| +VISTA   | <b>153.65 <math>\pm</math> 46.35</b>  | <b>354.37 <math>\pm</math> 86.45</b>  | <b>5080.02 <math>\pm</math> 3674.36</b> |

### F.4 ADDITIONAL EXPERIMENTS

To assess the effectiveness and scalability of VISTA, we conduct extensive experiments across a diverse set of synthetic graph families, varying both in size and structural complexity. This part is a detailed supplement of our Section 4.1. Specifically, we evaluate performance on 14 different graph configurations derived from ER and SF graphs, each instantiated with average degrees of 3 and 5, and node sizes  $n \in \{30, 50, 100, 300\}$ . This results in a comprehensive benchmark covering both sparse and dense regimes under varying dimensionalities. For each configuration, we benchmark recent representative causal discovery methods. Each method is tested under three settings: the original

Table 8: Comparison of total computing time (s) under SF5 setting.

| Method   | $n = 50$                              | $n = 100$                              | $n = 300$                               |
|----------|---------------------------------------|--|---|
| NOTEARS  | 808.15 $\pm$ 102.23                   | 2842.83 $\pm$ 312.22                   | 18676.80 $\pm$ 6873.83                  |
| +VISTA   | <b>501.96 <math>\pm</math> 62.14</b>  | <b>1200.41 <math>\pm</math> 536.14</b> | <b>3041.62 <math>\pm</math> 1003.68</b> |
| GOLEM    | 77.82 $\pm$ 11.90                     | 217.95 $\pm$ 73.36                     | 446.16 $\pm$ 30.04                      |
| +VISTA   | <b>23.71 <math>\pm</math> 2.37</b>    | <b>99.20 <math>\pm</math> 132.68</b>   | <b>167.89 <math>\pm</math> 214.19</b>   |
| DAG-GNN  | 911.14 $\pm$ 315.63                   | 5762.00 $\pm$ 1714.01                  | 31106.3 $\pm$ 452.12                    |
| +VISTA   | <b>356.46 <math>\pm</math> 101.18</b> | <b>1133.18 <math>\pm</math> 306.12</b> | <b>2641.62 <math>\pm</math> 541.84</b>  |
| GraN-DAG | 853.75 $\pm$ 98.54                    | 4944.93 $\pm$ 2325.58                  | 38163.22 $\pm$ 3919.71                  |
| +VISTA   | <b>313.24 <math>\pm</math> 120.79</b> | <b>934.83 <math>\pm</math> 218.82</b>  | <b>2999.39 <math>\pm</math> 485.66</b>  |
| SCORE    | 637.37 $\pm$ 48.40                    | 18904.31 $\pm$ 344.10                  | —                                       |
| +VISTA   | <b>187.91 <math>\pm</math> 38.43</b>  | <b>2003.45 <math>\pm</math> 882.48</b> | <b>4124.09 <math>\pm</math> 1311.74</b> |

baseline, VISTA with naive voting (+VISTA-NV), and VISTA with weighted voting (+VISTA-WV). Notably, CAM does not scale well with graph size and GraN-DAG fails when  $n$  reaches 300, so we do not report the results here. Due to the increasing computational cost with graph size, we ran each experimental configuration 10 times for  $n = 30$  and  $n = 50$ , 5 times for  $n = 100$ , and 3 times for  $n = 300$ , and report the average and standard deviation across trials.

Although the advantages of VISTA are most pronounced in high-dimensional or structurally complex settings, it is important to note that for some small-scale graphs, particularly relative sparse configurations such as ER3 with low node counts, the original base learners already achieve high accuracy. In these cases, the benefits of decomposition are less clear. Errors introduced during Markov Blanket identification and aggregation, as analyzed in Appendix D, may offset any gains from the divide-and-conquer process. When the true structure is relatively simple and well-recovered by the base model, additional processing may be unnecessary.

By contrast, as the graph size increases, structural coverage from local subgraphs becomes more reliable, and the advantages of localized inference and confidence-aware aggregation become more pronounced. In particular, VISTA consistently improves structural accuracy and reduces false discoveries for base learners that face scalability challenges in large and complex graphs, providing a practical approach to mitigating the curse of dimensionality in causal structure learning.

The remaining experimental results are as follows:

Table 9: Results with linear and nonlinear synthetic datasets ( $n = 30$ ,  $h = 5$ ).

| Method    | ER5                               |                                   |                                      |                                   | SF5                               |                                   |                                      |                                   |
|-----------|-----------------------------------|-----------------------------------|--------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|--------------------------------------|-----------------------------------|
|           | FDR $\downarrow$                  | TPR $\uparrow$                    | SHD $\downarrow$                     | F1 $\uparrow$                     | FDR $\downarrow$                  | TPR $\uparrow$                    | SHD $\downarrow$                     | F1 $\uparrow$                     |
| NOTEARS   | 0.21 $\pm$ 0.08                   | 0.73 $\pm$ 0.07                   | 64.70 $\pm$ 18.22                    | 0.76 $\pm$ 0.05                   | 0.24 $\pm$ 0.13                   | 0.70 $\pm$ 0.10                   | 64.30 $\pm$ 29.24                    | 0.73 $\pm$ 0.08                   |
| +VISTA-NV | 0.63 $\pm$ 0.00                   | <b>0.97 <math>\pm</math> 0.01</b> | 236.60 $\pm$ 7.32                    | 0.53 $\pm$ 0.01                   | 0.56 $\pm$ 0.02                   | <b>0.98 <math>\pm</math> 0.01</b> | 155.20 $\pm$ 11.90                   | 0.61 $\pm$ 0.02                   |
| +VISTA-WV | <b>0.12 <math>\pm</math> 0.04</b> | 0.75 $\pm$ 0.03                   | <b>50.00 <math>\pm</math> 8.64</b>   | <b>0.81 <math>\pm</math> 0.03</b> | <b>0.03 <math>\pm</math> 0.02</b> | 0.86 $\pm$ 0.03                   | <b>20.20 <math>\pm</math> 2.62</b>   | <b>0.84 <math>\pm</math> 0.02</b> |
| GOLEM     | 0.24 $\pm$ 0.08                   | 0.79 $\pm$ 0.07                   | 63.50 $\pm$ 23.61                    | 0.77 $\pm$ 0.05                   | 0.15 $\pm$ 0.11                   | 0.85 $\pm$ 0.09                   | 35.30 $\pm$ 23.39                    | 0.85 $\pm$ 0.07                   |
| +VISTA-NV | 0.65 $\pm$ 0.01                   | <b>0.97 <math>\pm</math> 0.01</b> | 251.00 $\pm$ 2.62                    | 0.52 $\pm$ 0.01                   | 0.56 $\pm$ 0.02                   | <b>0.99 <math>\pm</math> 0.01</b> | 158.40 $\pm$ 3.09                    | 0.61 $\pm$ 0.02                   |
| +VISTA-WV | <b>0.17 <math>\pm</math> 0.03</b> | 0.79 $\pm$ 0.05                   | <b>53.00 <math>\pm</math> 8.73</b>   | <b>0.81 <math>\pm</math> 0.04</b> | <b>0.01 <math>\pm</math> 0.01</b> | 0.89 $\pm$ 0.03                   | <b>15.00 <math>\pm</math> 4.32</b>   | <b>0.94 <math>\pm</math> 0.02</b> |
| DAG-GNN   | 0.29 $\pm$ 0.06                   | 0.77 $\pm$ 0.19                   | 77.50 $\pm$ 20.08                    | 0.74 $\pm$ 0.09                   | 0.29 $\pm$ 0.19                   | 0.72 $\pm$ 0.25                   | 65.50 $\pm$ 36.60                    | 0.72 $\pm$ 0.16                   |
| +VISTA-NV | 0.64 $\pm$ 0.00                   | <b>0.98 <math>\pm</math> 0.01</b> | 250.30 $\pm$ 4.78                    | 0.52 $\pm$ 0.00                   | 0.60 $\pm$ 0.02                   | <b>0.99 <math>\pm</math> 0.01</b> | 189.70 $\pm$ 14.06                   | 0.56 $\pm$ 0.02                   |
| +VISTA-WV | <b>0.27 <math>\pm</math> 0.03</b> | 0.84 $\pm$ 0.07                   | <b>60.00 <math>\pm</math> 7.85</b>   | <b>0.78 <math>\pm</math> 0.05</b> | <b>0.03 <math>\pm</math> 0.02</b> | 0.87 $\pm$ 0.04                   | <b>20.00 <math>\pm</math> 7.26</b>   | <b>0.92 <math>\pm</math> 0.03</b> |
| CAM       | 0.77 $\pm$ 0.04                   | 0.53 $\pm$ 0.08                   | 267.50 $\pm$ 23.32                   | 0.32 $\pm$ 0.04                   | 0.77 $\pm$ 0.06                   | 0.54 $\pm$ 0.11                   | 241.20 $\pm$ 37.51                   | 0.32 $\pm$ 0.06                   |
| +VISTA-NV | 0.78 $\pm$ 0.04                   | <b>0.66 <math>\pm</math> 0.12</b> | 327.00 $\pm$ 24.39                   | 0.33 $\pm$ 0.06                   | 0.82 $\pm$ 0.02                   | 0.57 $\pm$ 0.07                   | 335.60 $\pm$ 9.06                    | 0.27 $\pm$ 0.03                   |
| +VISTA-WV | <b>0.63 <math>\pm</math> 0.08</b> | 0.40 $\pm$ 0.10                   | <b>158.00 <math>\pm</math> 16.87</b> | <b>0.39 <math>\pm</math> 0.09</b> | <b>0.17 <math>\pm</math> 0.04</b> | <b>0.59 <math>\pm</math> 0.05</b> | <b>122.00 <math>\pm</math> 5.35</b>  | <b>0.69 <math>\pm</math> 0.04</b> |
| GraN-DAG  | 0.67 $\pm$ 0.12                   | 0.18 $\pm$ 0.10                   | 159.60 $\pm$ 24.52                   | 0.22 $\pm$ 0.10                   | 0.72 $\pm$ 0.38                   | 0.26 $\pm$ 0.03                   | 187.80 $\pm$ 80.73                   | 0.27 $\pm$ 0.18                   |
| +VISTA-NV | 0.90 $\pm$ 0.08                   | <b>0.39 <math>\pm</math> 0.12</b> | 211.70 $\pm$ 49.91                   | 0.16 $\pm$ 0.10                   | 0.85 $\pm$ 0.29                   | <b>0.41 <math>\pm</math> 0.15</b> | 239.50 $\pm$ 46.61                   | 0.22 $\pm$ 0.31                   |
| +VISTA-WV | <b>0.31 <math>\pm</math> 0.06</b> | 0.14 $\pm$ 0.07                   | <b>134.50 <math>\pm</math> 35.55</b> | <b>0.23 <math>\pm</math> 0.09</b> | <b>0.25 <math>\pm</math> 0.10</b> | 0.18 $\pm$ 0.05                   | <b>128.60 <math>\pm</math> 52.47</b> | <b>0.29 <math>\pm</math> 0.07</b> |
| SCORE     | 0.66 $\pm$ 0.08                   | 0.43 $\pm$ 0.05                   | 117.40 $\pm$ 47.71                   | 0.38 $\pm$ 0.05                   | 0.55 $\pm$ 0.40                   | 0.71 $\pm$ 0.24                   | 153.30 $\pm$ 76.60                   | 0.55 $\pm$ 0.31                   |
| +VISTA-NV | 0.80 $\pm$ 0.06                   | <b>0.83 <math>\pm</math> 0.05</b> | 399.60 $\pm$ 36.65                   | 0.32 $\pm$ 0.08                   | 0.76 $\pm$ 0.25                   | <b>0.88 <math>\pm</math> 0.04</b> | 440.60 $\pm$ 49.79                   | 0.38 $\pm$ 0.31                   |
| +VISTA-WV | <b>0.34 <math>\pm</math> 0.09</b> | 0.56 $\pm$ 0.08                   | <b>95.50 <math>\pm</math> 28.86</b>  | <b>0.61 <math>\pm</math> 0.06</b> | <b>0.36 <math>\pm</math> 0.16</b> | 0.79 $\pm$ 0.05                   | <b>88.80 <math>\pm</math> 11.60</b>  | <b>0.71 <math>\pm</math> 0.10</b> |

Table 10: Results with linear and nonlinear synthetic datasets ( $n = 50, h = 3$ ).

| Method    | ER3                |                    |                       |                    | SF3                |                    |                       |                    |
|-----------|--------------------|--------------------|-----------------------|--------------------|--------------------|--------------------|-----------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                  | F1↑                | FDR↓               | TPR↑               | SHD↓                  | F1↑                |
| NOTEARS   | <b>0.09 ± 0.08</b> | 0.90 ± 0.05        | <b>24.90 ± 17.75</b>  | <b>0.90 ± 0.05</b> | 0.22 ± 0.13        | 0.76 ± 0.11        | 62.80 ± 33.53         | 0.77 ± 0.08        |
| +VISTA-NV | 0.72 ± 0.03        | <b>0.97 ± 0.01</b> | 353.40 ± 57.19        | 0.44 ± 0.04        | 0.68 ± 0.04        | <b>0.97 ± 0.01</b> | 304.40 ± 52.48        | 0.48 ± 0.04        |
| +VISTA-WV | <b>0.09 ± 0.03</b> | 0.86 ± 0.03        | 30.50 ± 5.43          | 0.89 ± 0.02        | <b>0.07 ± 0.03</b> | 0.80 ± 0.07        | <b>36.90 ± 11.73</b>  | <b>0.85 ± 0.05</b> |
| GOLEM     | <b>0.04 ± 0.04</b> | <b>0.97 ± 0.02</b> | <b>8.10 ± 6.98</b>    | <b>0.97 ± 0.03</b> | 0.15 ± 0.06        | 0.84 ± 0.03        | 36.00 ± 11.81         | 0.84 ± 0.03        |
| +VISTA-NV | 0.76 ± 0.03        | 0.93 ± 0.03        | 431.60 ± 63.02        | 0.38 ± 0.04        | 0.75 ± 0.04        | <b>0.93 ± 0.03</b> | 397.70 ± 68.22        | 0.40 ± 0.05        |
| +VISTA-WV | 0.12 ± 0.05        | 0.76 ± 0.04        | 47.90 ± 9.80          | 0.81 ± 0.04        | <b>0.09 ± 0.11</b> | 0.85 ± 0.08        | <b>20.20 ± 16.75</b>  | <b>0.88 ± 0.07</b> |
| DAG-GNN   | 0.14 ± 0.08        | 0.86 ± 0.14        | 38.80 ± 26.22         | 0.86 ± 0.08        | 0.26 ± 0.10        | 0.73 ± 0.11        | 75.40 ± 24.08         | 0.73 ± 0.07        |
| +VISTA-NV | 0.73 ± 0.02        | <b>0.98 ± 0.00</b> | 380.00 ± 48.64        | 0.42 ± 0.03        | 0.73 ± 0.03        | <b>0.98 ± 0.00</b> | 378.00 ± 50.22        | 0.42 ± 0.04        |
| +VISTA-WV | <b>0.07 ± 0.04</b> | 0.84 ± 0.02        | <b>29.30 ± 2.05</b>   | <b>0.89 ± 0.01</b> | <b>0.05 ± 0.03</b> | 0.84 ± 0.05        | <b>41.00 ± 8.01</b>   | <b>0.89 ± 0.03</b> |
| CAM       | —                  | —                  | —                     | —                  | —                  | —                  | —                     | —                  |
| +VISTA-NV | 0.87 ± 0.02        | <b>0.66 ± 0.06</b> | 641.30 ± 67.62        | 0.22 ± 0.03        | 0.86 ± 0.02        | <b>0.71 ± 0.05</b> | 611.20 ± 64.82        | 0.24 ± 0.03        |
| +VISTA-WV | <b>0.66 ± 0.05</b> | 0.51 ± 0.07        | <b>192.00 ± 34.23</b> | <b>0.40 ± 0.05</b> | <b>0.65 ± 0.06</b> | 0.51 ± 0.10        | <b>181.80 ± 21.12</b> | <b>0.41 ± 0.07</b> |
| GraN-DAG  | 0.74 ± 0.32        | 0.09 ± 0.04        | 209.00 ± 54.45        | 0.11 ± 0.05        | 0.34 ± 0.42        | 0.08 ± 0.04        | 166.60 ± 42.38        | 0.12 ± 0.03        |
| +VISTA-NV | 0.75 ± 0.15        | <b>0.31 ± 0.06</b> | 158.80 ± 33.13        | 0.32 ± 0.08        | 0.48 ± 0.30        | <b>0.34 ± 0.09</b> | 195.20 ± 29.46        | <b>0.41 ± 0.11</b> |
| +VISTA-WV | <b>0.29 ± 0.08</b> | 0.26 ± 0.05        | <b>123.40 ± 15.51</b> | <b>0.38 ± 0.05</b> | <b>0.22 ± 0.21</b> | 0.20 ± 0.09        | <b>118.80 ± 23.75</b> | 0.32 ± 0.12        |
| SCORE     | 0.69 ± 0.05        | 0.67 ± 0.08        | 166.20 ± 59.57        | 0.42 ± 0.03        | 0.64 ± 0.06        | 0.64 ± 0.10        | 115.30 ± 31.50        | 0.45 ± 0.02        |
| +VISTA-NV | 0.86 ± 0.06        | <b>0.95 ± 0.03</b> | 980.80 ± 79.12        | 0.24 ± 0.09        | 0.90 ± 0.04        | <b>0.91 ± 0.05</b> | 923.70 ± 146.32       | 0.18 ± 0.07        |
| +VISTA-WV | <b>0.33 ± 0.04</b> | 0.74 ± 0.02        | <b>56.40 ± 19.95</b>  | <b>0.70 ± 0.02</b> | <b>0.49 ± 0.40</b> | 0.80 ± 0.06        | <b>74.60 ± 22.43</b>  | <b>0.62 ± 0.30</b> |

Table 11: Results with linear and nonlinear synthetic datasets ( $n = 50, h = 5$ ).

| Method    | ER5                |                    |                       |                    | SF5                |                    |                       |                    |
|-----------|--------------------|--------------------|-----------------------|--------------------|--------------------|--------------------|-----------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                  | F1↑                | FDR↓               | TPR↑               | SHD↓                  | F1↑                |
| NOTEARS   | 0.16 ± 0.09        | 0.81 ± 0.06        | 81.80 ± 34.23         | 0.82 ± 0.05        | 0.23 ± 0.08        | 0.75 ± 0.04        | 105.90 ± 26.65        | 0.76 ± 0.04        |
| +VISTA-NV | 0.75 ± 0.02        | <b>0.98 ± 0.01</b> | 685.60 ± 68.09        | 0.40 ± 0.02        | 0.72 ± 0.02        | <b>0.98 ± 0.01</b> | 585.30 ± 71.56        | 0.43 ± 0.03        |
| +VISTA-WV | <b>0.08 ± 0.04</b> | 0.76 ± 0.04        | <b>72.90 ± 13.75</b>  | <b>0.83 ± 0.03</b> | <b>0.15 ± 0.06</b> | 0.82 ± 0.04        | <b>71.90 ± 12.54</b>  | <b>0.84 ± 0.02</b> |
| GOLEM     | 0.34 ± 0.18        | 0.73 ± 0.14        | 156.20 ± 88.04        | 0.72 ± 0.14        | 0.30 ± 0.17        | 0.75 ± 0.12        | 130.20 ± 71.76        | 0.72 ± 0.11        |
| +VISTA-NV | 0.75 ± 0.02        | <b>0.96 ± 0.02</b> | 706.90 ± 66.89        | 0.39 ± 0.02        | 0.74 ± 0.03        | <b>0.94 ± 0.01</b> | 618.90 ± 91.19        | 0.41 ± 0.04        |
| +VISTA-WV | <b>0.20 ± 0.11</b> | 0.77 ± 0.08        | <b>99.40 ± 47.64</b>  | <b>0.79 ± 0.10</b> | <b>0.16 ± 0.09</b> | 0.77 ± 0.06        | <b>84.40 ± 34.04</b>  | <b>0.80 ± 0.08</b> |
| DAG-GNN   | 0.29 ± 0.15        | 0.70 ± 0.17        | 141.10 ± 63.30        | 0.71 ± 0.11        | 0.32 ± 0.11        | 0.74 ± 0.07        | 142.40 ± 45.01        | 0.71 ± 0.07        |
| +VISTA-NV | 0.76 ± 0.01        | <b>0.98 ± 0.01</b> | 720.70 ± 37.85        | 0.38 ± 0.02        | 0.74 ± 0.01        | <b>0.98 ± 0.01</b> | 633.00 ± 39.65        | 0.41 ± 0.01        |
| +VISTA-WV | <b>0.22 ± 0.07</b> | 0.79 ± 0.04        | <b>99.00 ± 26.95</b>  | <b>0.79 ± 0.05</b> | <b>0.27 ± 0.05</b> | 0.76 ± 0.03        | <b>116.60 ± 11.84</b> | <b>0.75 ± 0.01</b> |
| CAM       | —                  | —                  | —                     | —                  | —                  | —                  | —                     | —                  |
| +VISTA-NV | 0.86 ± 0.01        | <b>0.69 ± 0.05</b> | 978.00 ± 5.25         | 0.23 ± 0.02        | 0.86 ± 0.01        | <b>0.67 ± 0.06</b> | 941.60 ± 50.21        | 0.23 ± 0.02        |
| +VISTA-WV | <b>0.75 ± 0.02</b> | 0.49 ± 0.08        | <b>426.40 ± 26.82</b> | <b>0.32 ± 0.03</b> | <b>0.75 ± 0.03</b> | 0.47 ± 0.08        | <b>400.80 ± 41.11</b> | <b>0.33 ± 0.05</b> |
| GraN-DAG  | 0.62 ± 0.28        | 0.05 ± 0.03        | 265.80 ± 35.61        | 0.08 ± 0.04        | 0.64 ± 0.33        | 0.07 ± 0.05        | 271.80 ± 29.98        | 0.10 ± 0.06        |
| +VISTA-NV | 0.51 ± 0.08        | <b>0.17 ± 0.09</b> | 213.40 ± 82.48        | <b>0.25 ± 0.10</b> | 0.56 ± 0.15        | <b>0.26 ± 0.05</b> | 229.20 ± 63.68        | <b>0.32 ± 0.06</b> |
| +VISTA-WV | <b>0.36 ± 0.05</b> | 0.13 ± 0.02        | <b>204.00 ± 47.76</b> | 0.22 ± 0.03        | <b>0.27 ± 0.05</b> | 0.18 ± 0.04        | <b>193.00 ± 57.21</b> | 0.29 ± 0.05        |
| SCORE     | 0.73 ± 0.05        | <b>0.61 ± 0.15</b> | 431.60 ± 114.55       | 0.34 ± 0.02        | 0.71 ± 0.04        | 0.42 ± 0.04        | 365.00 ± 68.00        | 0.34 ± 0.03        |
| +VISTA-NV | 0.84 ± 0.01        | 0.47 ± 0.32        | 686.00 ± 62.50        | 0.24 ± 0.04        | 0.81 ± 0.03        | <b>0.54 ± 0.06</b> | 582.50 ± 57.50        | 0.28 ± 0.03        |
| +VISTA-WV | <b>0.64 ± 0.07</b> | 0.38 ± 0.06        | <b>271.00 ± 43.00</b> | <b>0.37 ± 0.05</b> | <b>0.35 ± 0.15</b> | 0.25 ± 0.04        | <b>210.50 ± 11.50</b> | <b>0.36 ± 0.04</b> |

Table 12: Results with linear and nonlinear synthetic datasets ( $n = 100, h = 3$ ).

| Method    | ER3                |                    |                        |                    | SF3                |                    |                        |                    |
|-----------|--------------------|--------------------|------------------------|--------------------|--------------------|--------------------|------------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                   | F1↑                | FDR↓               | TPR↑               | SHD↓                   | F1↑                |
| NOTEARS   | <b>0.09 ± 0.09</b> | 0.91 ± 0.06        | <b>54.60 ± 44.78</b>   | <b>0.91 ± 0.08</b> | 0.15 ± 0.08        | 0.75 ± 0.06        | 108.80 ± 36.84         | 0.80 ± 0.06        |
| +VISTA-NV | 0.81 ± 0.04        | <b>0.95 ± 0.01</b> | 1245.60 ± 349.77       | 0.32 ± 0.05        | 0.75 ± 0.04        | <b>0.95 ± 0.03</b> | 864.00 ± 146.07        | 0.39 ± 0.05        |
| +VISTA-WV | <b>0.09 ± 0.02</b> | 0.73 ± 0.02        | 99.00 ± 3.77           | 0.81 ± 0.01        | <b>0.11 ± 0.03</b> | 0.80 ± 0.05        | <b>92.00 ± 9.67</b>    | <b>0.85 ± 0.03</b> |
| GOLEM     | <b>0.09 ± 0.10</b> | <b>0.95 ± 0.05</b> | <b>39.80 ± 43.52</b>   | <b>0.93 ± 0.08</b> | 0.22 ± 0.05        | 0.72 ± 0.04        | 137.00 ± 22.24         | 0.75 ± 0.04        |
| +VISTA-NV | 0.84 ± 0.01        | 0.91 ± 0.02        | 1373.80 ± 202.28       | 0.28 ± 0.02        | 0.81 ± 0.03        | <b>0.90 ± 0.02</b> | 1180.20 ± 163.44       | 0.31 ± 0.04        |
| +VISTA-WV | 0.22 ± 0.02        | 0.65 ± 0.04        | 147.60 ± 17.55         | 0.71 ± 0.03        | <b>0.18 ± 0.13</b> | 0.78 ± 0.06        | <b>91.80 ± 72.13</b>   | <b>0.80 ± 0.07</b> |
| DAG-GNN   | 0.15 ± 0.11        | 0.71 ± 0.17        | 119.40 ± 63.78         | 0.77 ± 0.15        | 0.31 ± 0.14        | 0.54 ± 0.09        | 215.60 ± 47.23         | 0.59 ± 0.06        |
| +VISTA-NV | 0.63 ± 0.03        | <b>0.95 ± 0.01</b> | 1239.60 ± 131.07       | 0.53 ± 0.03        | 0.78 ± 0.04        | <b>0.94 ± 0.01</b> | 1058.40 ± 259.74       | 0.34 ± 0.06        |
| +VISTA-WV | <b>0.12 ± 0.02</b> | 0.82 ± 0.03        | <b>87.20 ± 15.30</b>   | <b>0.85 ± 0.02</b> | <b>0.23 ± 0.10</b> | 0.70 ± 0.08        | <b>151.20 ± 25.35</b>  | <b>0.73 ± 0.03</b> |
| GraN-DAG  | 0.90 ± 0.03        | 0.04 ± 0.02        | 463.40 ± 22.94         | 0.04 ± 0.02        | 0.83 ± 0.07        | 0.02 ± 0.02        | 366.40 ± 118.26        | 0.05 ± 0.02        |
| +VISTA-NV | 0.88 ± 0.06        | <b>0.25 ± 0.06</b> | 390.80 ± 73.58         | 0.17 ± 0.06        | 0.78 ± 0.06        | <b>0.26 ± 0.08</b> | 308.40 ± 49.60         | 0.24 ± 0.05        |
| +VISTA-WV | <b>0.38 ± 0.05</b> | 0.16 ± 0.03        | <b>250.60 ± 82.64</b>  | <b>0.25 ± 0.04</b> | <b>0.44 ± 0.08</b> | 0.18 ± 0.05        | <b>266.68 ± 67.76</b>  | <b>0.27 ± 0.06</b> |
| SCORE     | 0.91 ± 0.05        | 0.62 ± 0.04        | 2859.40 ± 839.4        | 0.16 ± 0.08        | 0.92 ± 0.03        | 0.66 ± 0.04        | 3131.20 ± 1002         | 0.14 ± 0.05        |
| +VISTA-NV | 0.94 ± 0.04        | <b>0.95 ± 0.12</b> | 2614.80 ± 566.5        | 0.11 ± 0.07        | 0.91 ± 0.05        | <b>0.70 ± 0.05</b> | 2727.60 ± 505.6        | 0.16 ± 0.08        |
| +VISTA-WV | <b>0.53 ± 0.05</b> | 0.75 ± 0.06        | <b>339.00 ± 189.43</b> | <b>0.58 ± 0.04</b> | <b>0.51 ± 0.10</b> | 0.68 ± 0.08        | <b>408.80 ± 205.64</b> | <b>0.57 ± 0.07</b> |

Table 13: Results with linear and nonlinear synthetic datasets ( $n = 300, h = 3$ ).

| Method    | ER3                |                    |                        |                    | SF3                |                    |                        |                    |
|-----------|--------------------|--------------------|------------------------|--------------------|--------------------|--------------------|------------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                   | F1↑                | FDR↓               | TPR↑               | SHD↓                   | F1↑                |
| NOTEARS   | <b>0.13 ± 0.03</b> | 0.89 ± 0.02        | <b>202.33 ± 48.98</b>  | <b>0.88 ± 0.03</b> | <b>0.16 ± 0.06</b> | 0.72 ± 0.04        | 519.33 ± 71.15         | <b>0.78 ± 0.04</b> |
| +VISTA-NV | 0.88 ± 0.01        | <b>0.91 ± 0.00</b> | 6177.00 ± 372.80       | 0.21 ± 0.02        | 0.89 ± 0.00        | <b>0.91 ± 0.01</b> | 6917.67 ± 998.81       | 0.20 ± 0.00        |
| +VISTA-WV | 0.23 ± 0.02        | 0.66 ± 0.03        | 462.00 ± 54.31         | 0.71 ± 0.02        | 0.21 ± 0.03        | 0.55 ± 0.03        | <b>363.00 ± 76.53</b>  | 0.65 ± 0.02        |
| GOLEM     | 0.23 ± 0.15        | 0.76 ± 0.18        | <b>375.67 ± 258.67</b> | <b>0.77 ± 0.16</b> | 0.56 ± 0.19        | 0.34 ± 0.22        | 913.33 ± 304.25        | 0.38 ± 0.21        |
| +VISTA-NV | 0.88 ± 0.00        | <b>0.85 ± 0.01</b> | 5389.67 ± 46.91        | 0.22 ± 0.00        | 0.86 ± 0.03        | <b>0.48 ± 0.28</b> | 3248.67 ± 1304.26      | 0.18 ± 0.08        |
| +VISTA-WV | <b>0.17 ± 0.02</b> | 0.45 ± 0.01        | 628.33 ± 11.90         | 0.50 ± 0.01        | <b>0.21 ± 0.06</b> | 0.44 ± 0.04        | <b>597.00 ± 53.59</b>  | <b>0.56 ± 0.04</b> |
| DAG-GNN   | 0.55 ± 0.34        | 0.19 ± 0.19        | 1288.33 ± 832.49       | 0.21 ± 0.20        | 0.72 ± 0.17        | 0.23 ± 0.22        | 1264.33 ± 484.00       | 0.22 ± 0.19        |
| +VISTA-NV | 0.89 ± 0.01        | <b>0.93 ± 0.00</b> | 6449.00 ± 89.16        | 0.19 ± 0.01        | 0.89 ± 0.01        | <b>0.89 ± 0.02</b> | 6627.00 ± 651.98       | 0.19 ± 0.02        |
| +VISTA-WV | <b>0.18 ± 0.06</b> | 0.57 ± 0.07        | <b>494.67 ± 50.22</b>  | <b>0.66 ± 0.04</b> | <b>0.34 ± 0.06</b> | 0.49 ± 0.07        | <b>633.33 ± 111.52</b> | <b>0.55 ± 0.03</b> |
| SCORE     | —                  | —                  | —                      | —                  | —                  | —                  | —                      | —                  |
| +VISTA-NV | 0.95 ± 0.00        | <b>0.76 ± 0.04</b> | 11064.00 ± 371.63      | 0.09 ± 0.01        | 0.97 ± 0.01        | <b>0.44 ± 0.13</b> | 13057.00 ± 3556.57     | 0.06 ± 0.02        |
| +VISTA-WV | <b>0.19 ± 0.02</b> | 0.32 ± 0.03        | <b>666.67 ± 18.66</b>  | <b>0.46 ± 0.03</b> | <b>0.61 ± 0.29</b> | 0.08 ± 0.04        | <b>970.67 ± 141.74</b> | <b>0.13 ± 0.07</b> |

Table 14: Results with linear and nonlinear synthetic datasets ( $n = 300, h = 5$ ).

| Method    | ER5                |                    |                         |                    | SF5                |                    |                         |                    |
|-----------|--------------------|--------------------|-------------------------|--------------------|--------------------|--------------------|-------------------------|--------------------|
|           | FDR↓               | TPR↑               | SHD↓                    | F1↑                | FDR↓               | TPR↑               | SHD↓                    | F1↑                |
| NOTEARS   | 0.30 ± 0.05        | 0.68 ± 0.12        | 875.33 ± 205.07         | 0.69 ± 0.09        | 0.50 ± 0.09        | 0.22 ± 0.12        | 1402.33 ± 70.59         | 0.29 ± 0.14        |
| +VISTA-NV | 0.93 ± 0.02        | <b>0.94 ± 0.01</b> | 6520.33 ± 1357.12       | 0.10 ± 0.02        | 0.90 ± 0.01        | <b>0.78 ± 0.04</b> | 12180.00 ± 1008.42      | 0.18 ± 0.02        |
| +VISTA-WV | <b>0.15 ± 0.04</b> | 0.67 ± 0.05        | <b>689.67 ± 98.89</b>   | <b>0.75 ± 0.03</b> | <b>0.24 ± 0.02</b> | 0.38 ± 0.03        | <b>890.33 ± 166.61</b>  | <b>0.51 ± 0.03</b> |
| GOLEM     | 0.81 ± 0.05        | 0.10 ± 0.03        | 1921.33 ± 111.21        | 0.13 ± 0.04        | 0.93 ± 0.03        | 0.02 ± 0.01        | 1839.67 ± 139.17        | 0.03 ± 0.02        |
| +VISTA-NV | 0.92 ± 0.01        | <b>0.28 ± 0.14</b> | 5551.00 ± 1310.12       | 0.12 ± 0.03        | 0.92 ± 0.00        | <b>0.77 ± 0.09</b> | 13437.00 ± 1562.43      | 0.14 ± 0.00        |
| +VISTA-WV | <b>0.20 ± 0.23</b> | 0.23 ± 0.02        | <b>1225.00 ± 38.79</b>  | <b>0.36 ± 0.03</b> | <b>0.37 ± 0.11</b> | 0.10 ± 0.06        | <b>1391.00 ± 61.65</b>  | <b>0.17 ± 0.10</b> |
| DAG-GNN   | 0.91 ± 0.06        | 0.33 ± 0.17        | 3858.00 ± 1558.37       | 0.17 ± 0.04        | 0.91 ± 0.04        | 0.15 ± 0.06        | 4617.33 ± 3064.50       | 0.10 ± 0.04        |
| +VISTA-NV | 0.90 ± 0.03        | <b>0.86 ± 0.04</b> | 8988.00 ± 910.33        | 0.18 ± 0.05        | 0.91 ± 0.03        | <b>0.81 ± 0.04</b> | 14578.67 ± 4342.92      | 0.16 ± 0.05        |
| +VISTA-WV | <b>0.37 ± 0.14</b> | 0.25 ± 0.05        | <b>1920.33 ± 809.62</b> | <b>0.36 ± 0.06</b> | <b>0.41 ± 0.03</b> | 0.21 ± 0.06        | <b>2191.33 ± 656.02</b> | <b>0.31 ± 0.09</b> |
| SCORE     | —                  | —                  | —                       | —                  | —                  | —                  | —                       | —                  |
| +VISTA-NV | 0.96 ± 0.00        | <b>0.17 ± 0.07</b> | 18762.67 ± 2501.28      | 0.06 ± 0.01        | 0.98 ± 0.00        | <b>0.13 ± 0.15</b> | 22039.00 ± 2028.89      | 0.03 ± 0.01        |
| +VISTA-WV | <b>0.93 ± 0.00</b> | 0.10 ± 0.03        | <b>1698.33 ± 103.76</b> | <b>0.07 ± 0.01</b> | <b>0.95 ± 0.02</b> | 0.08 ± 0.06        | <b>2582.67 ± 830.34</b> | <b>0.06 ± 0.02</b> |

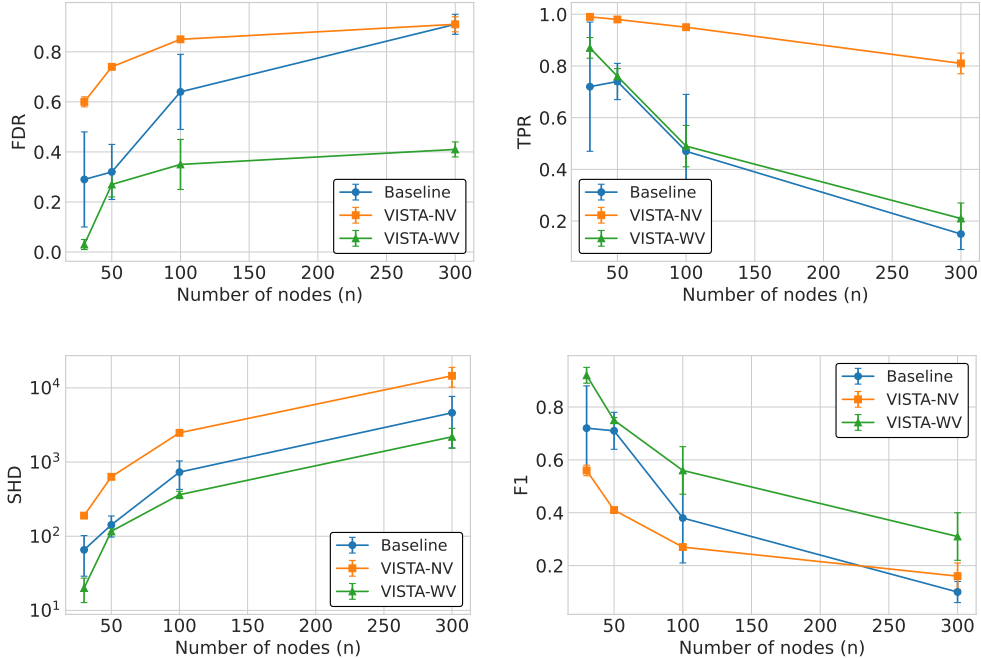


Figure 5: Performance of DAG-GNN on SF5 Graphs.

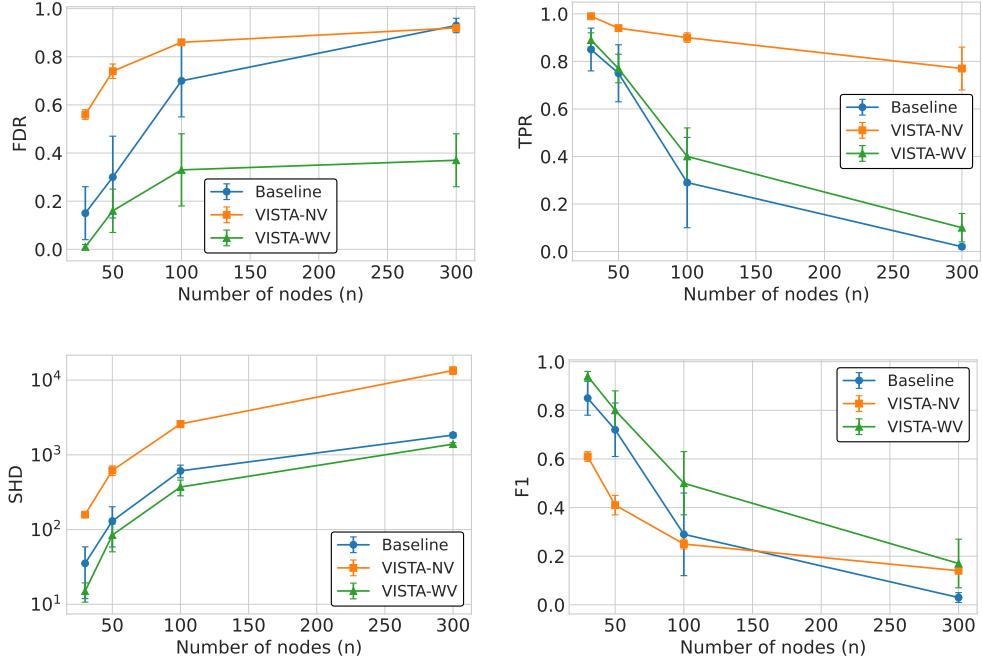


Figure 6: Performance of GOLEM on SF5 Graphs.

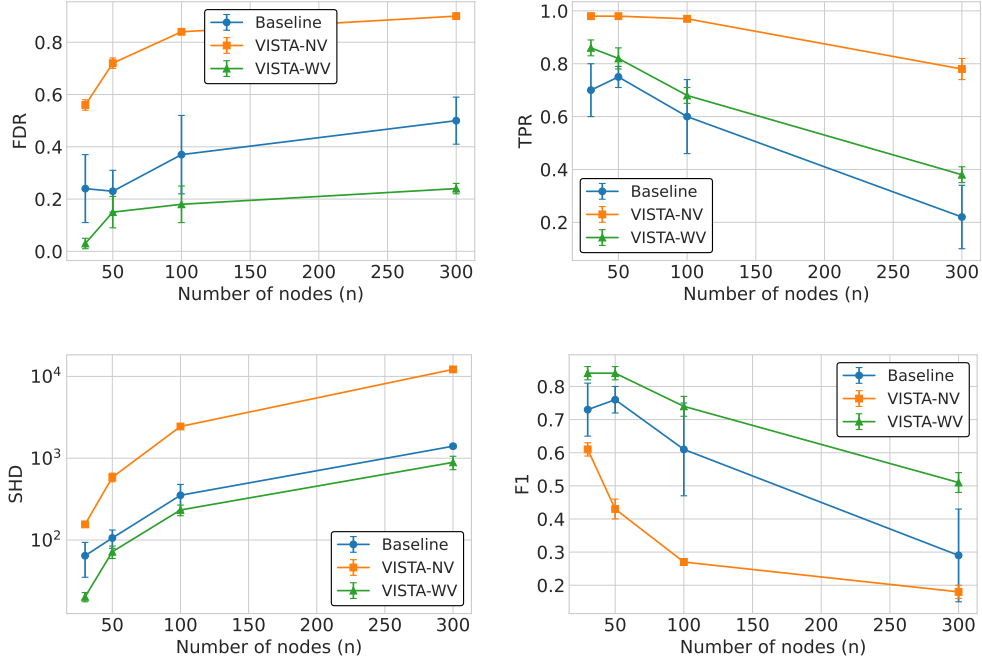


Figure 7: Performance of NOTEARS on SF5 Graphs.

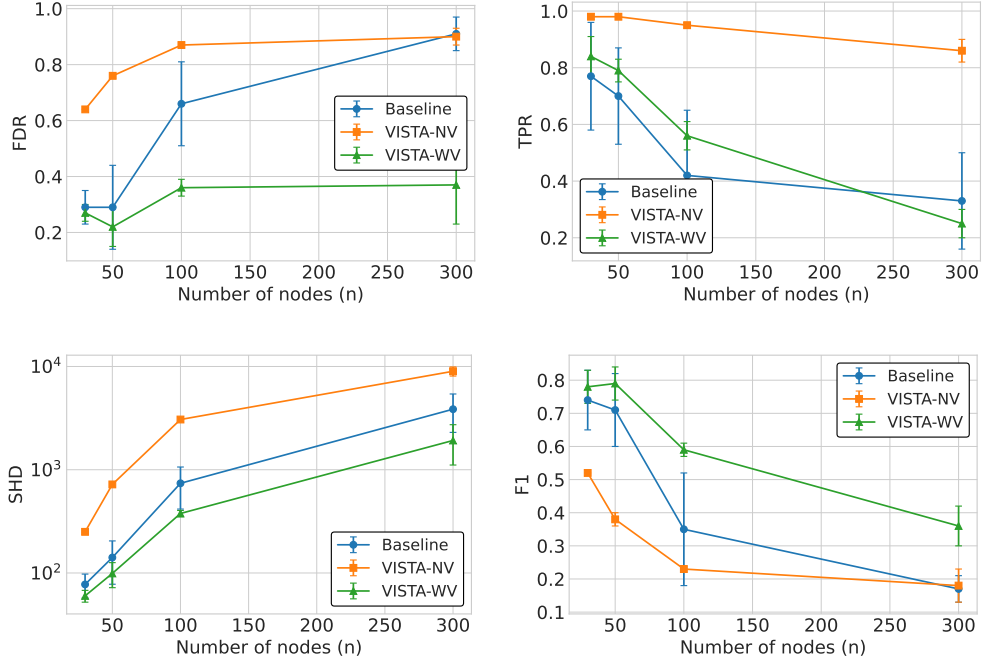


Figure 8: Performance of DAG-GNN on ER5 Graphs.

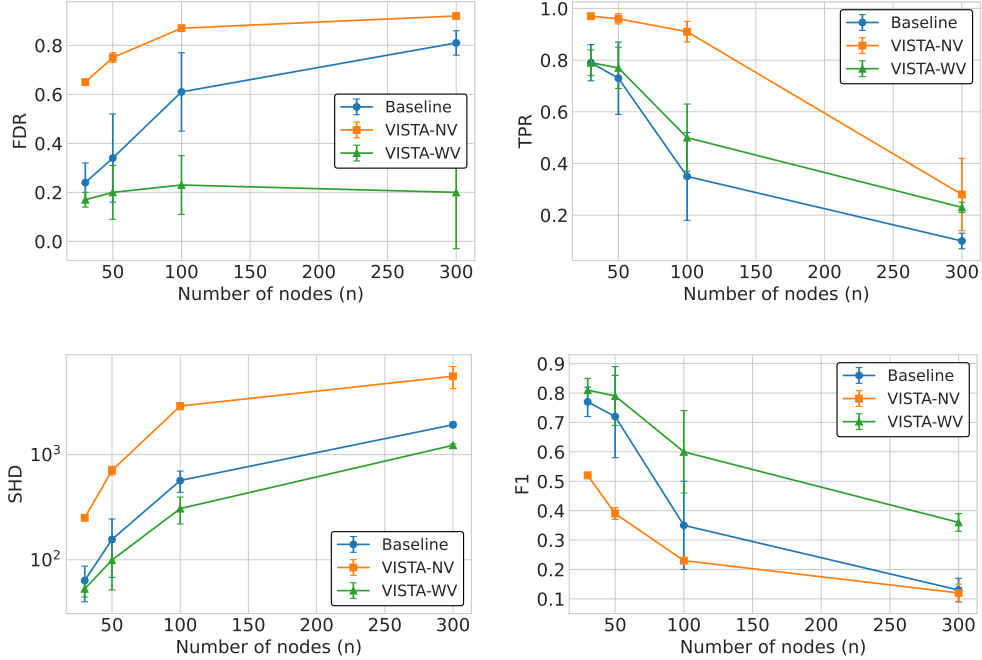


Figure 9: Performance of GOLEM on ER5 Graphs.



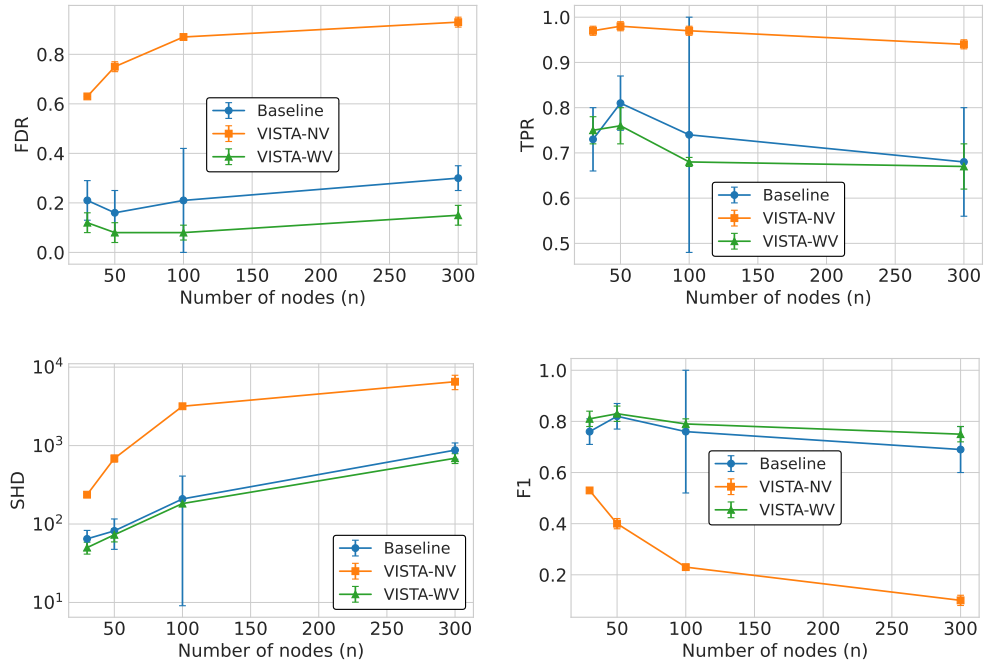


Figure 10: Performance of NOTEARS on ER5 Graphs.