
A Nearly-Linear Time Framework for Graph-Structured Sparsity

Chinmay Hegde
Piotr Indyk
Ludwig Schmidt*

CHINMAY@MIT.EDU
INDYK@MIT.EDU
LUDWIGS@MIT.EDU

Massachusetts Institute of Technology, Cambridge, MA 02139, USA

Abstract

We introduce a framework for sparsity structures defined via graphs. Our approach is flexible and generalizes several previously studied sparsity models. Moreover, we provide efficient projection algorithms for our sparsity model that run in nearly-linear time. In the context of sparse recovery, we show that our framework achieves an information-theoretically optimal sample complexity for a wide range of parameters. We complement our theoretical analysis with experiments demonstrating that our algorithms also improve on prior work in practice.

1. Introduction

Over the past decade, sparsity has emerged as an important tool in several fields including signal processing, statistics, and machine learning. In compressive sensing, sparsity reduces the sample complexity of measuring a signal, and statistics utilizes sparsity for high-dimensional inference tasks. In many settings, sparsity is a useful ingredient because it enables us to model structure in high-dimensional data while still remaining a mathematically tractable concept. For instance, natural images are often sparse when represented in a wavelet basis, and objects in a classification task usually belong to only a small number of classes.

Due to the success of sparsity, a natural question is how we can refine the notion of sparsity in order to capture more complex structures. There are many examples where such an approach is applicable: (i) large wavelet coefficients of natural images tend to form connected *trees*, (ii) active genes can be arranged in functional *groups*, and (iii) approximate point sources in astronomical data often form *clusters*. In such cases, exploiting this additional structure can lead to improved compression ratio for images, bet-

ter multi-label classification, or smaller sample complexity in compressive sensing and statistics. Hence an important question is the following: how can we model such sparsity structures, and how can we make effective use of this additional information in a computationally efficient manner?

There has been a wide range of work addressing these questions, e.g., (Yuan & Lin, 2006; Jacob et al., 2009; He & Carin, 2009; Kim & Xing, 2010; Bi & Kwok, 2011; Huang et al., 2011; Duarte & Eldar, 2011; Bach et al., 2012b; Rao et al., 2012; Negahban et al., 2012; Simon et al., 2013; El Halabi & Cevher, 2015). Usually, the proposed solutions offer a trade-off between the following conflicting goals:

Generality What range of sparsity structures does the approach apply to?

Statistical efficiency What statistical performance improvements does the use of structure enable?

Computational efficiency How fast are the resulting algorithms?

In this paper, we introduce a framework for sparsity models defined through *graphs*, and we show that it achieves a compelling trade-off between the goals outlined above. At a high level, our approach applies to data with an underlying graph structure in which the large coefficients form a small number of connected components (optionally with additional constraints on the edges). Our approach offers three main features: (i) *Generality*: the framework encompasses several previously studied sparsity models, e.g., tree sparsity and cluster sparsity. (ii) *Statistical efficiency*: our sparsity model leads to reduced sample complexity in sparse recovery and achieves the information-theoretic optimum for a wide range of parameters. (iii) *Computational efficiency*: we give a nearly-linear time algorithm for our sparsity model, significantly improving on prior work both in theory and in practice. Due to the growing size of data sets encountered in science and engineering, algorithms with (nearly-)linear running time are becoming increasingly important.

We achieve these goals by connecting our sparsity model to the *prize collecting Steiner tree* (PCST) problem, which has been studied in combinatorial optimization and approximation algorithms. To establish this connection, we introduce a generalized version of the PCST problem and give a nearly-linear time algorithm for our variant. We believe that our sparsity model and the underlying algorithms are useful beyond sparse recovery, and we have already obtained results in this direction. To keep the presentation in this paper coherent, we focus on our results for sparse recovery and briefly mention further applications in Sec. 7.

Before we present our theoretical results in Sections 3 to 5, we give an overview in Section 2. Section 6 complements our theoretical results with an empirical evaluation on both synthetic and real data (a background-subtracted image, an angiogram, and an image of text). We defer proofs and additional details to the supplementary material.

Basic notation Let $[d]$ be the set $\{1, 2, \dots, d\}$. We say that a vector $\beta \in \mathbb{R}^d$ is s -sparse if at most s of its coefficients are nonzero. The support of β contains the indices corresponding to nonzero entries in β , i.e., $\text{supp}(\beta) = \{i \in [d] \mid \beta_i \neq 0\}$. Given a subset $S \subseteq [d]$, we write β_S for the restriction of β to indices in S : we have $(\beta_S)_i = \beta_i$ for $i \in S$ and $(\beta_S)_i = 0$ otherwise. The ℓ_2 -norm of β is $\|\beta\| = \sqrt{\sum_{i \in [d]} \beta_i^2}$.

Sparsity models In some cases, we have more information about a vector than only “standard” s -sparsity. A natural way of encoding such additional structure is via *sparsity models* (Baraniuk et al., 2010): let \mathbb{M} be a family of supports, i.e., $\mathbb{M} = \{S_1, S_2, \dots, S_L\}$ where $S_i \subseteq [d]$. Then the corresponding sparsity model \mathcal{M} is the set of vectors supported on one of the S_i :

$$\mathcal{M} = \{\beta \in \mathbb{R}^d \mid \text{supp}(\beta) \subseteq S \text{ for some } S \in \mathbb{M}\}. \quad (1)$$

2. Our contributions

We state our main contributions in the context of sparse recovery (see Section 7 for further applications). Our goal is to estimate an unknown s -sparse vector $\beta \in \mathbb{R}^d$ from observations of the form

$$y = X\beta + e, \quad (2)$$

where $X \in \mathbb{R}^{n \times d}$ is the design matrix, $y \in \mathbb{R}^n$ are the observations, and $e \in \mathbb{R}^n$ is an observation noise vector. By imposing various assumptions on X and e , sparse recovery encompasses problems such as sparse linear regression and compressive sensing.

2.1. Weighted graph model (WGM)

The core of our framework for structured sparsity is a novel, **general** sparsity model which we call the *weighted*

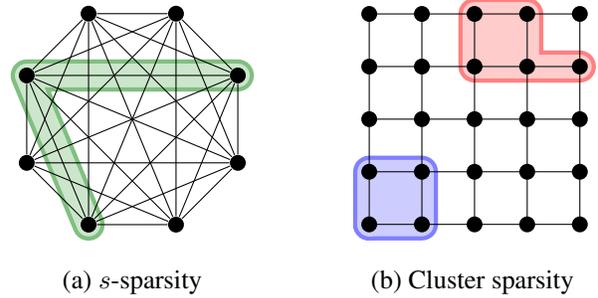


Figure 1. Two examples of the weighted graph model. (a) In a complete graph, any s -sparse support can be mapped to a single tree ($g = 1$). (b) Using a grid graph, we can model a small number of clusters in an image by setting g accordingly. For simplicity, we use unit edge weights and set $B = s - g$ in both examples.

graph model. In the WGM, we use an underlying graph $G = (V, E)$ defined on the coefficients of the unknown vector β , i.e., $V = [d]$. Moreover, the graph is weighted and we denote the edge weights with $w : E \rightarrow \mathbb{N}$. We identify supports $S \subseteq [d]$ with subgraphs in G , in particular *forests* (unions of individual trees). Intuitively, the WGM captures sparsity structures with a small number of connected components in G . In order to control the sparsity patterns, the WGM offers three parameters:

- s , the total sparsity of S .
- g , the maximum number of connected components formed by the forest F corresponding to S .
- B , the bound on the total weight $w(F)$ of edges in the forest F corresponding to S .

More formally, let $\gamma(H)$ be the number of connected components in a graph H . Then we can define the WGM:

Definition 1. The (G, s, g, B) -WGM is the set of supports

$$\mathbb{M} = \{S \subseteq [d] \mid |S| = s \text{ and there is a } F \subseteq G \text{ with } V_F = S, \gamma(F) = g, \text{ and } w(F) \leq B\}. \quad (3)$$

Fig. 1 shows how two sparsity structures can be encoded with the WGM. Since our sparsity model applies to *arbitrary* graphs G , it can describe a wide range of structures. In particular, the model generalizes several previously studied sparsity models, including 1D-clusters, (wavelet) tree hierarchies, the Earth Mover Distance (EMD) model, and the unweighted graph model (see Table 1).

2.2. Recovery of vectors in the WGM

We analyze the **statistical efficiency** of our framework in the context of sparse recovery. In particular, we prove that

the sample complexity of recovering vectors in the WGM is provably smaller than the sample complexity for “standard” s -sparse vectors. To formally state this result, we first introduce a key property of graphs.

Definition 2. Let $G = (V, E)$ be a weighted graph with edge weights $w : E \rightarrow \mathbb{N}$. Then the weight-degree $\rho(v)$ of a node v is the largest number of adjacent nodes connected by edges with the same weight, i.e.,

$$\rho(v) = \max_{b \in \mathbb{N}} |\{(v', v) \in E \mid w(v', v) = b\}|. \quad (4)$$

We define the weight-degree of G to be the maximum weight-degree of any $v \in V$.

Note that for graphs with uniform edge weights, the weight-degree of G is the same as the maximum node degree. Intuitively, the (weight) degree of a graph is an important property for quantifying the sample complexity of the WGM because the degree determines how restrictive the bound on the number of components g is. In the extreme case of a complete graph, any support can be formed with only a single connected component (see Figure 1). Using Definitions 1 and 2, we now state our sparse recovery result (see Theorem 12 in Section 5 for a more general version):

Theorem 3. Let $\beta \in \mathbb{R}^d$ be in the (G, s, g, B) -WGM. Then

$$n = O\left(s \left(\log \rho(G) + \log \frac{B}{s}\right) + g \log \frac{d}{g}\right) \quad (5)$$

i.i.d. Gaussian observations suffice to estimate β . More precisely, let $e \in \mathbb{R}^n$ be an arbitrary noise vector and let $y \in \mathbb{R}^n$ be defined as in Eq. 2 where X is an i.i.d. Gaussian matrix. Then we can efficiently find an estimate $\hat{\beta}$ such that

$$\|\beta - \hat{\beta}\| \leq C \|e\|, \quad (6)$$

where C is a constant independent of all variables above.

Note that in the noiseless case ($e = 0$), we are guaranteed to recover β exactly. Moreover, our estimate $\hat{\beta}$ is in a slightly enlarged WGM for any amount of noise, see Section 5. Our bound (5) can be instantiated to recover previous sample complexity results, e.g., the $n = O(s \log \frac{d}{s})$ bound for “standard” sparse recovery, which is tight (Do Ba et al., 2010).¹ For the image grid graph example in Figure 1, Equation (5) becomes $n = O(s + g \log \frac{d}{g})$, which matches the information-theoretic optimum $n = O(s)$ as long as the number of clusters is not too large, i.e., $g = O(s/\log d)$.²

¹To be precise, encoding s -sparsity with a complete graph as in Figure 1 gives a bound of $n = O(s \log d)$. To match the $\log \frac{d}{s}$ term, we can encode s -sparsity as $g = s$ clusters of size one in a fully disconnected graph with no edges.

²Optimality directly follows from a simple dimensionality argument: even if the s -sparse support of the vector β is known, recovering the unknown coefficients requires solving a linear system with s unknowns uniquely. For this, we need at least s linear equations, i.e., s observations.

2.3. Efficient projection into the WGM

The algorithmic core of our sparsity framework is a **computationally efficient** procedure for projecting arbitrary vectors into the WGM. More precisely, the model-projection problem is the following: given a vector $b \in \mathbb{R}^d$ and a WGM \mathcal{M} , find the best approximation to b in \mathcal{M} , i.e.,

$$P_{\mathcal{M}}(b) = \arg \min_{b' \in \mathcal{M}} \|b - b'\|. \quad (7)$$

If such a model-projection algorithm is available, one can instantiate the framework of (Baraniuk et al., 2010) in order to get an algorithm for sparse recovery with the respective sparsity model.³ However, solving Problem (7) exactly is NP-hard for the WGM due to a reduction from the classical Steiner tree problem (Karp, 1972). To circumvent this hardness result, we use the *approximation-tolerant* framework of (Hegde et al., 2014a). Instead of solving (7) exactly, the framework requires *two* algorithms with the following complementary approximation guarantees.

Tail approximation: Find an $S \in \mathbb{M}$ such that

$$\|b - b_S\| \leq c_T \cdot \min_{S' \in \mathbb{M}} \|b - b_{S'}\|. \quad (8)$$

Head approximation: Find an $S \in \mathbb{M}$ such that

$$\|b_S\| \geq c_H \cdot \max_{S' \in \mathbb{M}} \|b_{S'}\|. \quad (9)$$

Here, $c_T > 1$ and $c_H < 1$ are arbitrary, fixed constants. Note that a head approximation guarantee does not imply a tail guarantee (and vice versa). In fact, stable recovery is not possible with only one type of approximate projection guarantee (Hegde et al., 2014a). We provide two algorithms for solving (8) and (9) (one per guarantee) which both run in *nearly-linear time*.

Our model-projection algorithms are based on a connection to the prize-collecting Steiner tree problem (PCST), which is a generalization of the classical Steiner tree problem. Instead of finding the cheapest way to connect *all* terminal nodes in a given weighted graph, we can instead omit some terminals from the solution and pay a specific price for each omitted node. The goal is to find a subtree with the optimal trade-off between the cost paid for edges used to connect a subset of the nodes and the price of the remaining, unconnected nodes (see Section 3 for a formal definition).

We make the following three main algorithmic contributions. Due to the wide applicability of the PCST problem, we believe that these algorithms can be of independent interest (see Section 7).

³Note that the framework does not supply general projection algorithms. Instead, the model-projection algorithms have to be designed from scratch for each model.

- We introduce a variant of the PCST problem in which the goal is to find a set of g trees instead of a single tree. We call this variant the prize-collecting Steiner forest (PCSF) problem and adapt the algorithm of (Goemans & Williamson, 1995) for this variant.
- We reduce the projection problems (8) and (9) to a small set of adaptively constructed PCSF instances.
- We give a nearly-linear time algorithm for the PCSF problem and hence also the model projection problem.

2.4. Improvements for existing sparsity models

Our results are directly applicable to several previously studied sparsity models that can be encoded with the WGM. Table 1 summarizes these results. In spite of its generality, our approach at least matches the sample complexity of prior work in all cases and actually offers an improvement for the EMD model. Moreover, our running time is always within a polylogarithmic factor of the best algorithm, even in the case of models with specialized solvers such as tree sparsity. For the EMD and cluster models, our algorithm is significantly faster than prior work and improves the time complexity by a polynomial factor. To complement these theoretical results, our experiments in Section 6 show that our algorithm is more than one order of magnitude faster than previous algorithms with provable guarantees and offers a better sample complexity in many cases.

2.5. Comparison to related work

In addition to the “point-solutions” for individual sparsity models outlined above, there has been a wide range of work on general frameworks for utilizing structure in sparse recovery. The approach most similar to ours is (Baraniuk et al., 2010), which gives a framework underlying many of the algorithms in Table 1. However, the framework has one important drawback: it does not come with a full recovery algorithm. Instead, the authors only give a recovery scheme that assumes the existence of a model-projection algorithm satisfying (7). Such an algorithm must be constructed from scratch for each model, and the techniques that have been used for various models so far are quite different. Our contribution can be seen as complementing the framework of (Baraniuk et al., 2010) with a nearly-linear time projection algorithm that is applicable to a wide range of sparsity structures. This answers a question raised by the authors of (Huang et al., 2011), who also give a framework for structured sparsity with a universal and complete recovery algorithm. Their framework is applicable to a wide range of sparsity models, but the corresponding algorithm is significantly slower than ours, both in theory (“Graph clusters” in Table 1) and in practice (see Section 6). Moreover, our recovery algorithm shows more robust performance across different shapes of graph clusters.

Both of the approaches mentioned above use iterative greedy algorithms for sparse recovery. There is also a large body of work on combining M-estimators with convex regularizers that induce structured sparsity, e.g., see the surveys (Bach et al., 2012a) and (Wainwright, 2014). The work closest to ours is (Jacob et al., 2009), which uses an overlapping group Lasso to enforce graph-structured sparsity (graph Lasso). In contrast to their approach, our algorithm gives more fine-grained control over the number of clusters in the graph. Moreover, our algorithm has better computational complexity, and to the best of our knowledge there are no formal results relating the graph structure to the sample complexity of the graph Lasso. Empirically, our algorithm recovers an unknown vector with graph structure faster and from fewer observations than the graph Lasso (see Section A in the supplementary material).

3. The prize-collecting Steiner forest problem

We now establish our connection between prize-collecting Steiner tree (PCST) problems and the weighted graph model. First, we formally define the PCST problem: Let $G = (V, E)$ be an undirected, weighted graph with edge costs $c : E \rightarrow \mathbb{R}_0^+$ and node prizes $\pi : V \rightarrow \mathbb{R}_0^+$. For a subset of edges $E' \subseteq E$, we write $c(E') = \sum_{e \in E'} c(e)$ and adopt the same convention for node subsets. Moreover, for a node subset $V' \subseteq V$, let \bar{V}' be the complement $\bar{V}' = V \setminus V'$. Then the goal of the PCST problem is to find a subtree $T = (V', E')$ such that $c(E') + \pi(\bar{V}')$ is minimized. We sometimes write $c(T)$ and $\pi(\bar{T})$ if the node and edge sets are clear from context.

Similar to the classical Steiner tree problem, PCST is NP-hard. Most algorithms with provable approximation guarantees build on the seminal work of (Goemans & Williamson, 1995) (GW), who gave an efficient primal-dual algorithm with the following guarantee:

$$c(T) + 2\pi(\bar{T}) \leq 2 \min_{T' \text{ is a tree}} c(T') + \pi(\bar{T}'). \quad (10)$$

Note that the PCST problem already captures three important aspects of the WGM: (i) there is an underlying graph G , (ii) edges are weighted, and (iii) nodes have prizes. If we set the prizes to correspond to vector coefficients, i.e., $\pi(i) = b_i^2$, the term $\pi(\bar{T})$ in the PCST objective function becomes $\pi(\bar{T}) = \|b - b_T\|^2$, which matches the objective in the model-projection problems (8) and (9). However, there are two important differences. First, the objective in the PCST problem is to find a *single* tree T , while the WGM can contain supports defined by *multiple* connected components (if $g > 1$). Moreover, the PCST problem optimizes the trade-off $c(T) + \pi(\bar{T})$, but we are interested in minimizing $\|b - b_T\|$ subject to hard constraints on the support cardinality $|T|$ and the support cost $c(T)$ (the parameters s and B , respectively). In this section, we ad-

Table 1. Results of our sparsity framework applied to several sparsity models. In order to simplify the running time bounds, we assume that all coefficients are polynomially bounded in d , and that $s \leq d^{1/2-\mu}$ for some $\mu > 0$. For the graph cluster model, we consider the case of graphs with constant degree. The exponent τ depends on the degree of the graph and is always greater than 1. The parameters w and h are specific to the EMD model, see (Hegde et al., 2014a) for details. We always have $w \cdot h = d$ and $s \geq w$. Our sparsity framework improves on the sample complexity and running time of both the EMD and graph cluster models (bold entries).

Model	Reference	Best previous sample complexity	Our sample complexity	Best previous running time	Our running time
1D-cluster	(Cevher et al., 2009b)	$O(s + g \log \frac{d}{g})$	$O(s + g \log \frac{d}{g})$	$O(d \log^2 d)$	$O(d \log^4 d)$
Trees	(Hegde et al., 2014b)	$O(s)$	$O(s)$	$O(d \log^2 d)$	$O(d \log^4 d)$
EMD	(Hegde et al., 2014a)	$O(s \log \frac{B \log \frac{s}{w}}{s})$	$O(s \log \frac{B}{s})$	$O(sh^2 B \log d)$	$O(wh^2 \log^4 d)$
Graph clusters	(Huang et al., 2011)	$O(s + g \log d)$	$O(s + g \log \frac{d}{g})$	$O(d^\tau)$	$O(d \log^4 d)$

dress the first of these two issues; Section 4 then completes the connection between PCST and the WGM. We begin by defining the following variant of the PCST problem.

Definition 4 (The prize-collecting Steiner forest problem). *Let $g \in \mathbb{N}$ be the target number of connected components. Then the goal of the prize-collecting Steiner forest (PCSF) problem is to find a subgraph $F = (V', E')$ with $\gamma(F) = g$ that minimizes $c(E') + \pi(\overline{V'})$.*

As defined in Section 2.1, $\gamma(F)$ is the number of connected components in the (sub-)graph F . To simplify notation in the rest of the paper, we say that a forest F is a g -forest if $\gamma(F) = g$. There is always an optimal solution for the PCSF problem which consists of g trees because removing edges cannot increase the objective value. This allows us to employ the PCSF problem for finding supports in the WGM that consist of several connected components. In order to give a computationally efficient algorithm for the PCSF variant, we utilize prior work for PCST: (i) To show correctness of our algorithm, we prove that the GW scheme for PCST can be adapted to our PCSF variant. (ii) To achieve a good time complexity, we show how to simulate the GW scheme in nearly-linear time.

3.1. The Goemans-Williamson (GW) scheme for PCSF

A useful view of the GW scheme is the “moat-growing” interpretation of (Jünger & Pulleyblank, 1995), which describes the algorithm as an iterative clustering method that constructs “moats” around every cluster. These moats are essentially the dual variables in the linear program of the GW scheme. Initially, every node forms its own active cluster with a moat of size 0. The moats around each active cluster then grow at a uniform rate until one of the following two events occurs:

Cluster deactivation When the sum of moats in a cluster reaches the sum of node prizes in that cluster, the cluster is deactivated.

Cluster merge When the sum of moats that are intersected by an edge e reaches the cost of e , the clusters at the two endpoints of e are merged and e is added to the current solution.

The moat-growing stage of the algorithm terminates when only a single active cluster remains. After that, the resulting set of edges is pruned in order to achieve a provable approximation ratio. We generalize the proof of (Feofiloff et al., 2010) and show that it is possible to extract more than one tree from the moat-growing phase as long as the trees come from different clusters. Our modification of GW terminates the moat-growing phase when exactly g active clusters remain, and we then apply the GW pruning algorithm to each resulting tree separately. This gives the following result.

Theorem 5. *There is an algorithm for the PCSF problem that returns a g -forest F such that*

$$c(F) + 2\pi(\overline{F}) \leq \min_{F' \subseteq G, \gamma(F') \leq g} 2c(F') + 2\pi(\overline{F'}). \quad (11)$$

For $g = 1$, the theorem recovers the guarantee in (10). We defer the proof to Sec. D.1 of the supplementary material.

3.2. A fast algorithm for Goemans-Williamson

While the modified GW scheme produces good approximate solutions, it is not yet sufficient for a nearly-linear time algorithm: we still need an efficient way of simulating the moat-growing phase. There are two main difficulties: (i) The remaining “slack” amounts on edges can shrink at different rates depending on how many of the edge endpoints are in active clusters. (ii) A single cluster event (merge or deactivation) can change this rate for up to $\Theta(|V|)$ edges. In order to maintain edge events efficiently, we use the dynamic edge splitting approach introduced by (Cole et al., 2001). This technique essentially ensures that every edge always has at most one active endpoint, and hence its slack either shrinks at rate 1 or not

at all. However, edge splitting introduces additional edge events that do not directly lead to a cluster merge. While it is relatively straightforward to show that every such intermediate edge event halves the remaining amount of slack on an edge, we still need an overall bound on the number of intermediate edge events necessary to achieve a given precision. For this, we prove the following new result about the GW moat growing scheme.

Theorem 6. *Let all edge costs $c(e)$ and node prizes $\pi(v)$ be even integers. Then all finished moats produced by the GW scheme have integer sizes.*

In a nutshell, this theorem shows that one additional bit of precision is enough to track all events in the moat-growing phase accurately. We prove the theorem via induction over the events in the GW scheme, see Section D.2.2 for details. On the theoretical side, this result allows us to bound the overall running time of our algorithm for PCSF. Combined with suitable data structures, we can show the following:

Theorem 7. *Let α be the number of bits used to specify a single edge cost or node prize. Then there is an algorithm achieving the PCSF guarantee of Theorem 5 in time $O(\alpha \cdot |E| \log |V|)$.*

On the practical side, we complement Theorem 6 with a new adaptive edge splitting scheme that leads to a small number of intermediate edge events. Our experiments show that our scheme results in less than 3 events per edge on average (see Section D.3 in the supplementary material).

4. Sparse approximation with the WGM

In order to utilize the WGM in sparse recovery, we employ the framework of (Hegde et al., 2014a). As outlined in Section 2.3, the framework requires us to construct two approximation algorithms satisfying the head- and tail-approximation guarantees (8) and (9). We now give two such model-projection algorithms, building on our tools for PCSF developed in the previous section.

4.1. Tail-approximation algorithm

We can connect the PCSF objective to the WGM quantities by setting $\pi(i) = b_i^2$ and $c(e) = w(e) + 1$, which gives:

$$c(F) = w(F) + (|F| - g) \quad \text{and} \quad \pi(\overline{F}) = \|b - b_F\|^2.$$

After multiplying the edge costs with a trade-off parameter λ , the PCSF objective $\lambda \cdot c(F) + \pi(\overline{F})$ essentially becomes a *Lagrangian relaxation* of the model-constrained optimization problem (8). We build our tail-approximation algorithm on top of this connection, starting with an algorithm for the “tail”-variant of the PCSF problem. By performing a binary search over the parameter λ (see Algorithm 1), we get a bicriterion guarantee for the final forest.

Algorithm 1 PCSF-TAIL

```

1: Input:  $G, c, \pi, g, \text{cost-budget } C, \text{ parameters } \nu \text{ and } \delta.$ 
2: We write  $c_\lambda(e) = \lambda \cdot c(e).$ 
3:  $\pi_{\min} \leftarrow \min_{\pi(i) > 0} \pi(i), \quad \lambda_0 \leftarrow \frac{\pi_{\min}}{2C}$ 
4:  $F \leftarrow \text{PCSF-GW}(G, c_{\lambda_0}, \pi, g)$ 
5: if  $c(F) \leq 2C$  and  $\pi(\overline{F}) = 0$  then return  $F$ 
6:  $\lambda_r \leftarrow 0, \quad \lambda_l \leftarrow 3\pi(G), \quad \varepsilon \leftarrow \frac{\pi_{\min} \delta}{C}$ 
7: while  $\lambda_l - \lambda_r > \varepsilon$  do
8:    $\lambda_m \leftarrow (\lambda_l + \lambda_r)/2$ 
9:    $F \leftarrow \text{PCSF-GW}(G, c_{\lambda_m}, \pi, g)$ 
10:  if  $c(F) \geq C$  and  $c(F) \leq \nu C$  then return  $F$ 
11:  if  $c(F) > \gamma C$  then  $\lambda_r \leftarrow \lambda_m$  else  $\lambda_l \leftarrow \lambda_m$ 
12: end while
13: return  $F \leftarrow \text{PCSF-GW}(G, c_{\lambda_l}, \pi, g)$ 

```

Theorem 8. *Let $\nu > 2$ and $\delta > 0$. Then PCSF-TAIL returns a g -forest $F \subseteq G$ such that $c(F) \leq \nu \cdot C$ and*

$$\pi(\overline{F}) \leq \left(1 + \frac{2}{\nu - 2} + \delta\right) \min_{\gamma(F') = g, c(F') \leq C} \pi(\overline{F'}). \quad (12)$$

Theorem 8 does not give $c(F) \leq C$ exactly, but the cost of the resulting forest is still guaranteed to be within a constant factor of C . The framework of (Hegde et al., 2014a) also applies to projections into such slightly larger models. As we will see in Section 5, this increase by a constant factor does not affect the sample complexity.

For the trade-off between support size and support weight, we also make use of approximation. By scalarizing the two constraints carefully, i.e., setting $c(e) = w(e) + \frac{B}{s}$, we get the following result. The proofs of Theorems 8 and 9 can be found in the supplementary material, Section C.1.

Theorem 9. *Let \mathcal{M} be a (G, s, g, B) -WGM, let $b \in \mathbb{R}^d$, and let $\nu > 2$. Then there is an algorithm that returns a support $S \subseteq [d]$ in the $(G, 2\nu \cdot s + g, g, 2\nu \cdot B)$ -WGM satisfying (8) with $c_T = \sqrt{1 + 3/(\nu - 2)}$. Moreover, the algorithm runs in time $O(|E| \log^3 d)$.*

4.2. Head-approximation algorithm

For our head-approximation algorithm, we also use the PCSF objective as a Lagrangian relaxation of the model-constraint problem (9). This time, we multiply the node prizes instead of the edge costs with a parameter λ . We perform a binary search similar to Alg. 1, but the final step of the algorithm requires an extra subroutine. At the end of the binary search, we are guaranteed to have a forest with good “density” $\frac{\pi(F)}{c(F)}$, but the good forest could correspond to either the lower bound λ_l or the upper bound λ_r . In the latter case, we have no bound on the cost of the corresponding forest F_r . However, it is always possible to extract a high-density sub-forest with bounded cost from F_r :

Algorithm 2 GRAPH-COSAMP

```

1: Input:  $y, X, G, s, g, B$ , number of iterations  $t$ .
2:  $\hat{\beta}_0 \leftarrow 0$ 
3: for  $i \leftarrow 1, \dots, t$  do
4:    $b \leftarrow X^T(y - X\hat{\beta}_{i-1})$ 
5:    $S' \leftarrow \text{supp}(\hat{\beta}_{i-1}) \cup \text{HEADAPPROX}'(b, G, s, g, B)$ 
6:    $z_{S'} \leftarrow X_{S'}^\dagger y, \quad z_{S'^c} \leftarrow 0$ 
7:    $S \leftarrow \text{TAILAPPROX}(z, G, s, g, B)$ 
8:    $\hat{\beta}_i \leftarrow z_S$ 
9: end for
10: return  $\hat{\beta} \leftarrow \hat{\beta}_i$ 

```

Lemma 10. *Let T be a tree and $C' \leq c(T)$. Then there is a subtree $T' \subseteq T$ such that $c(T') \leq C'$ and $\pi(T') \geq \frac{C'}{6} \cdot \frac{\pi(T)}{c(T)}$. Moreover, we can find such a subtree T' in linear time.*

The algorithm first converts the tree into a list of nodes corresponding to a tour through T . Then we can extract a good sub-tree by either returning a single, high-prize node or sliding a variable-size window across the list of nodes. See Section C.2 in the supplementary material for details. Combining these components, we get a head-approximation algorithm with the following properties.

Theorem 11. *Let \mathcal{M} be a (G, s, g, B) -WGM and let $b \in \mathbb{R}^d$. Then there is an algorithm that returns a support $S \subseteq [d]$ in the $(G, 2s + g, g, 2B)$ -WGM satisfying (9) with $c_H = \sqrt{1/14}$. The algorithm runs in time $O(|E| \log^3 d)$.*

5. Application in sparse recovery

We now instantiate the framework of (Hegde et al., 2014a) to give a sparse recovery algorithm using the WGM. The resulting algorithm (see Alg. 2) is a variant of CoSaMP (Needell & Tropp, 2009) and uses the head- and tail-approximation algorithms instead of the hard thresholding operators.⁴ In order to state the corresponding recovery guarantee in full generality, we briefly review the definition of the (model-) restricted isometry property (RIP) (Candès & Tao, 2005; Baraniuk et al., 2010). We say that a matrix X satisfies the (\mathcal{M}, δ) -model-RIP if for all $\beta \in \mathcal{M}$:

$$(1 - \delta) \cdot \|\beta\|^2 \leq \|X\beta\|^2 \leq (1 + \delta) \cdot \|\beta\|^2. \quad (13)$$

Theorem 12. *Let $\beta \in \mathbb{R}^d$ be in the (G, s, g, B) -WGM \mathcal{M} and let $X \in \mathbb{R}^{n \times d}$ be a matrix satisfying the model-RIP for a $(G, c_1 s, g, c_2 B)$ -WGM and a fixed constant δ , where c_1 and c_2 are also fixed constants. Moreover, let $e \in \mathbb{R}^n$ be an arbitrary noise vector and let $y \in \mathbb{R}^n$ be defined as in (2).*

⁴Strictly speaking, HEADAPPROX' is a "boosted" version of the head-approximation algorithm developed here. See (Hegde et al., 2014a) for details.

Then GRAPH-COSAMP returns a $\hat{\beta}$ in the $(G, 5s, g, 5B)$ -WGM such that $\|\beta - \hat{\beta}\| \leq c_3 \|e\|$, where c_3 is a fixed constant. Moreover, GRAPH-COSAMP runs in time

$$O\left((T_X + |E| \log^3 d) \log \frac{\|\beta\|}{\|e\|}\right),$$

where T_X is the time complexity of a matrix-vector multiplication with X .

In order to establish sample complexity bounds for concrete matrix ensembles (e.g., random Gaussian matrices as in Theorem 3), we use a result of (Baraniuk et al., 2010) that relates the sample complexity of sub-Gaussian matrix ensembles to the size of the model, i.e., the quantity $|\mathbb{M}|$. More precisely, $n = O(s + \log|\mathbb{M}|)$ rows / observations suffice for such matrices to satisfy the model-RIP for \mathcal{M} and a fixed constant δ . For the WGM, we use a counting argument to bound $|\mathbb{M}|$ (see Section B in the supplementary material). Together with Theorem 12, the following theorem establishes Theorem 3 from Section 2.2.

Theorem 13. *Let \mathbb{M} be the set of supports in the (G, s, g, B) -WGM. Then*

$$\log|\mathbb{M}| = O\left(s \left(\log \rho(G) + \log \frac{B}{s}\right) + g \log \frac{d}{g}\right).$$

Next, we turn our attention to the running time of GRAPH-COSAMP. Since our model-projection algorithms run in nearly-linear time, the matrix-vector products involving X can become the bottleneck in the overall time complexity:⁵ for a dense Gaussian matrix, we have $T_X = \Omega(sd)$, which would dominate the overall running time. If we can control the design matrix (as is often the case in compressive sensing), we can use the construction of (Hegde et al., 2014b) to get a sample-optimal matrix with nearly-linear T_X in the regime of $s \leq d^{1/2-\mu}$, $\mu > 0$. Such a matrix then gives an algorithm with nearly-linear running time. Note that the bound on s is only a theoretical restriction in this construction: as our experiments show, a partial Fourier matrix empirically performs well for significantly larger values of s .

6. Experiments

We focus on the performance of our algorithm Graph-CoSaMP for the task of recovering 2D data with clustered sparsity. Multiple methods have been proposed for this problem, and our theoretical analysis shows that our algorithm should improve upon the state of the art (see Table 1). We compare our results to StructOMP (Huang et al., 2011) and the heuristic Lattice Matching Pursuit (LaMP) (Cevher et al., 2009a). The implementations were supplied by the

⁵It is not necessary to compute a full pseudo-inverse X^\dagger . See (Needell & Tropp, 2009) for details.

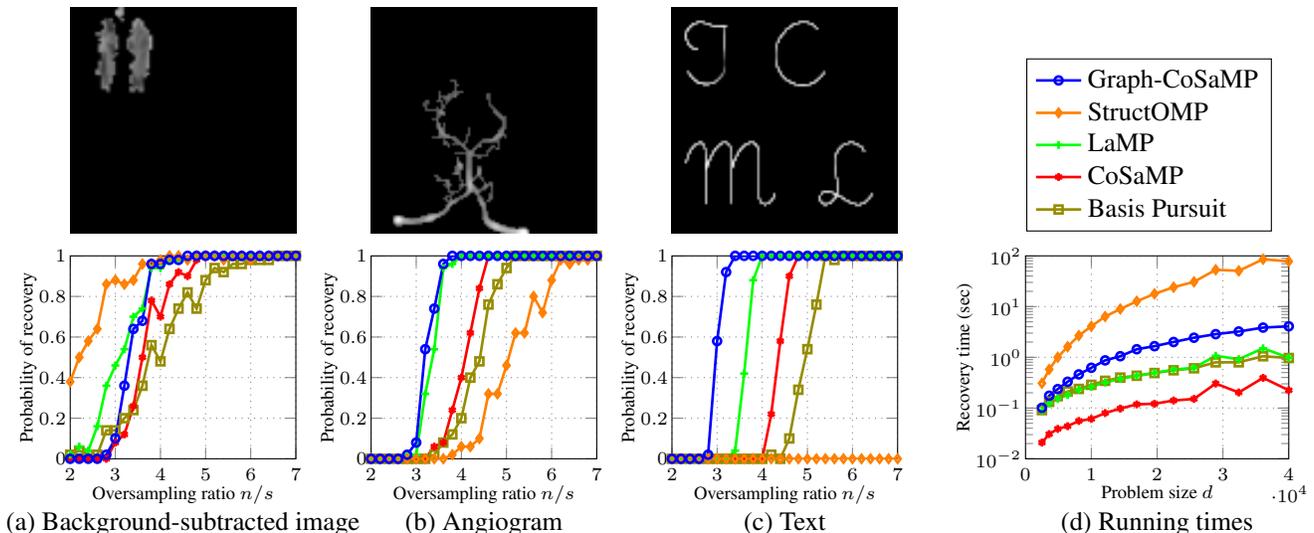


Figure 2. Sparse recovery experiments. The images in the top row are the original images β . In the regime where the algorithms recover with high probability, the estimates $\hat{\beta}$ are essentially identical to the original images. Our algorithm Graph-CoSaMP achieves consistently good recovery performance and offers the best sample complexity for images (b) and (c). Moreover, our algorithm is about 20 times faster than StructOMP, the other method with provable guarantees for the image cluster model.

authors and we used the default parameter settings. Moreover, we ran two common recovery algorithms for “standard” s -sparsity: Basis Pursuit (Candès et al., 2006) and CoSaMP (Needell & Tropp, 2009).

We follow a standard evaluation procedure for sparse recovery / compressive sensing: we record n observations $y = X\beta$ of the (vectorized) image $\beta \in \mathbb{R}^d$ using a subsampled Fourier matrix X . We assume that all algorithms possess prior knowledge of the sparsity s and the number of connected-components g in the true support of the image β . We declare a trial successful if the squared ℓ_2 -norm of the recovery error is at most 5% of the squared ℓ_2 -norm of the original vector β . The probability of successful recovery is then estimated by averaging over 50 trials. We perform several experiments with varying oversampling ratios n/s and three different images. See Section A in the supplementary material for a description of the dataset, experiments with noise, and a comparison with the graph Lasso.

Figure 2 demonstrates that Graph-CoSaMP yields consistently competitive phase transitions and exhibits the best sample complexity for images with “long” connected clusters, such as the angiogram image (b) and the text image (c). While StructOMP performs well on “blob”-like images such as the background-subtracted image (a), its performance is poor in our other test cases. For example, it can successfully recover the text image only for oversampling ratios $n/s > 15$. Note that the performance of Graph-CoSaMP is very consistent: in all three examples, the phase transition occurs between oversampling ratios 3 and 4. Other methods show significantly more variability.

We also investigate the computational efficiency of Graph-CoSaMP. We consider resized versions of the angiogram image and record $n = 6s$ observations for each image size d . Figure 2(d) displays the recovery times (averaged over 50 trials) as a function of d . We observe that the runtime of Graph-CoSaMP scales nearly linearly with d , comparable to the conventional sparse recovery methods. Moreover, Graph-CoSaMP is about $20\times$ faster than StructOMP.

7. Further applications

We expect our algorithms to be useful beyond sparse recovery and now briefly describe two promising applications.

Seismic feature extraction In (Schmidt et al., 2015), the authors use Steiner tree methods for a seismic feature extraction task. Our new algorithms for PCSF give a principled way of choosing tuning parameters for their proposed optimization problem. Moreover, our fast algorithms for PCSF can speed-up their method.

Event detection in social networks (Rozenshtein et al., 2014) introduce a method for event detection in social networks based on the PCST problem. Their method performs well but produces spurious results in the presence of multiple disconnected events because their PCST algorithm produces only a *single* tree instead of a forest. Our new algorithm for PCSF gives exact control over the number of trees in the solution and hence directly addresses this issue. Furthermore, the authors quote a running time of $O(|V|^2 \log|V|)$ for their GW scheme, so our nearly-linear time algorithm allows their method to scale to larger data.

Acknowledgements

We thank Jayadev Acharya, Stefanie Jegelka, Youssef Mroueh, Devavrat Shah, and the anonymous reviewers for many helpful comments on earlier versions of this paper. This work was supported by grants from the MITEL-Shell program, the MADALGO center, and a Simons Investigator Award.

References

- Bach, Francis, Jenatton, Rodolphe, Mairal, Julien, and Obozinski, Guillaume. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4(1):1–106, 2012a.
- Bach, Francis, Jenatton, Rodolphe, Mairal, Julien, and Obozinski, Guillaume. Structured sparsity through convex optimization. *Statistical Science*, 27(4):450–468, 11 2012b.
- Baraniuk, Richard G., Cevher, Volkan, Duarte, Marco F., and Hegde, Chinmay. Model-based compressive sensing. *IEEE Transactions on Information Theory*, 56(4):1982–2001, 2010.
- Bateni, MohammadHossein, Chekuri, Chandra, Ene, Alina, Hajiaghayi, Mohammad T., Korula, Nitish, and Marx, Daniel. Prize-collecting steiner problems on planar graphs. In *Proceedings of the Twenty-second Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 1028–1049, 2011.
- Bi, Wei and Kwok, James T. Multi-label classification on tree- and DAG-structured hierarchies. In *Proceedings of the 28th International Conference on Machine Learning (ICML)*, pp. 17–24, 2011.
- Candès, Emmanuel J. and Tao, Terence. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.
- Candès, Emmanuel J., Romberg, Justin, and Tao, Terence. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- Cevher, Volkan, Duarte, Marco F., Hegde, Chinmay, and Baraniuk, Richard. Sparse signal recovery using markov random fields. In *Advances in Neural Information Processing Systems 21 (NIPS)*, pp. 257–264. 2009a.
- Cevher, Volkan, Indyk, Piotr, Hegde, Chinmay, and Baraniuk, Richard G. Recovery of clustered sparse signals from compressive measurements. In *International Conference on Sampling Theory and Applications (SAMPTA)*, 2009b.
- Cole, Richard, Hariharan, Ramesh, Lewenstein, Moshe, and Porat, Ely. A faster implementation of the Goemans-Williamson clustering algorithm. In *Proceedings of the Twelfth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 17–25, 2001.
- Do Ba, Khanh, Indyk, Piotr, Price, Eric, and Woodruff, David P. Lower bounds for sparse recovery. In *Proceedings of the 21st Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 1190–1197, 2010.
- Duarte, Marco F. and Eldar, Yonina C. Structured compressed sensing: From theory to applications. *IEEE Transactions on Signal Processing*, 59(9):4053–4085, 2011.
- Eisenstat, David, Klein, Philip, and Mathieu, Claire. An efficient polynomial-time approximation scheme for steiner forest in planar graphs. In *Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 626–638, 2012.
- El Halabi, Marwa and Cevher, Volkan. A totally unimodular view of structured sparsity. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2015.
- Feofiloff, Paulo, Fernandes, Cristina G., Ferreira, Carlos E., and de Pina, José Coelho. A note on Johnson, Minkoff and Phillips’ algorithm for the prize-collecting Steiner tree problem. *Computing Research Repository (CoRR)*, abs/1004.1437, 2010.
- Goemans, Michel X. and Williamson, David P. A general approximation technique for constrained forest problems. *SIAM Journal on Computing*, 24(2):296–317, 1995.
- He, Lihan and Carin, Lawrence. Exploiting structure in wavelet-based bayesian compressive sensing. *IEEE Transactions on Signal Processing*, 57(9):3488–3497, 2009.
- Hegde, Chinmay, Indyk, Piotr, and Schmidt, Ludwig. Approximation algorithms for model-based compressive sensing. *Computing Research Repository (CoRR)*, abs/1406.1579, 2014a. Conference version appeared in the *Proceedings of the 25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*.
- Hegde, Chinmay, Indyk, Piotr, and Schmidt, Ludwig. Nearly linear-time model-based compressive sensing. In *Automata, Languages, and Programming (ICALP)*, volume 8572 of *Lecture Notes in Computer Science*, pp. 588–599. 2014b.
- Huang, Junzhou, Zhang, Tong, and Metaxas, Dimitris. Learning with structured sparsity. *The Journal of Machine Learning Research*, 12:3371–3412, 2011.

- Jacob, Laurent, Obozinski, Guillaume, and Vert, Jean-Philippe. Group Lasso with overlap and graph Lasso. In *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*, pp. 433–440, 2009.
- Johnson, David S., Minkoff, Maria, and Phillips, Steven. The prize collecting Steiner tree problem: Theory and practice. In *Proceedings of the Eleventh Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 760–769, 2000.
- Jünger, Michael and Pulleyblank, William R. New primal and dual matching heuristics. *Algorithmica*, 13(4):357–380, 1995.
- Karp, Richard M. Reducibility among combinatorial problems. In *Complexity of Computer Computations*, The IBM Research Symposia Series, pp. 85–103. 1972.
- Kim, Seyoung and Xing, Eric P. Tree-guided group Lasso for multi-task regression with structured sparsity. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pp. 543–550, 2010.
- Mosci, Sofia, Villa, Silvia, Verri, Alessandro, and Rosasco, Lorenzo. A primal-dual algorithm for group sparse regularization with overlapping groups. In *Advances in Neural Information Processing Systems 23 (NIPS)*, pp. 2604–2612. 2010.
- Needell, Deanna and Tropp, Joel A. CoSaMP: Iterative signal recovery from incomplete and inaccurate samples. *Applied and Computational Harmonic Analysis*, 26(3): 301–321, 2009.
- Negahban, Sahand N., Ravikumar, Pradeep, Wainwright, Martin J., and Yu, Bin. A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. *Statistical Science*, 27(4):538–557, 11 2012.
- Rao, Nikhil S., Recht, Ben, and Nowak, Robert D. Universal measurement bounds for structured sparse signal recovery. In *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 22 of *JMLR Proceedings*, pp. 942–950, 2012.
- Rozenshtein, Polina, Anagnostopoulos, Aris, Gionis, Aristides, and Tatti, Nikolaj. Event detection in activity networks. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 1176–1185, 2014.
- Schmidt, Ludwig, Hegde, Chinmay, Indyk, Piotr, Lu, Ligang, Chi, Xingang, and Hohl, Detlef. Seismic feature extraction using Steiner tree methods. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015.
- Simon, Noah, Friedman, Jerome, Hastie, Trevor, and Tibshirani, Robert. A sparse-group Lasso. *Journal of Computational and Graphical Statistics*, 22(2):231–245, 2013.
- Wainwright, Martin J. Structured regularizers for high-dimensional problems: Statistical and computational issues. *Annual Review of Statistics and Its Application*, 1 (1):233–253, 2014.
- Yuan, Ming and Lin, Yi. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.

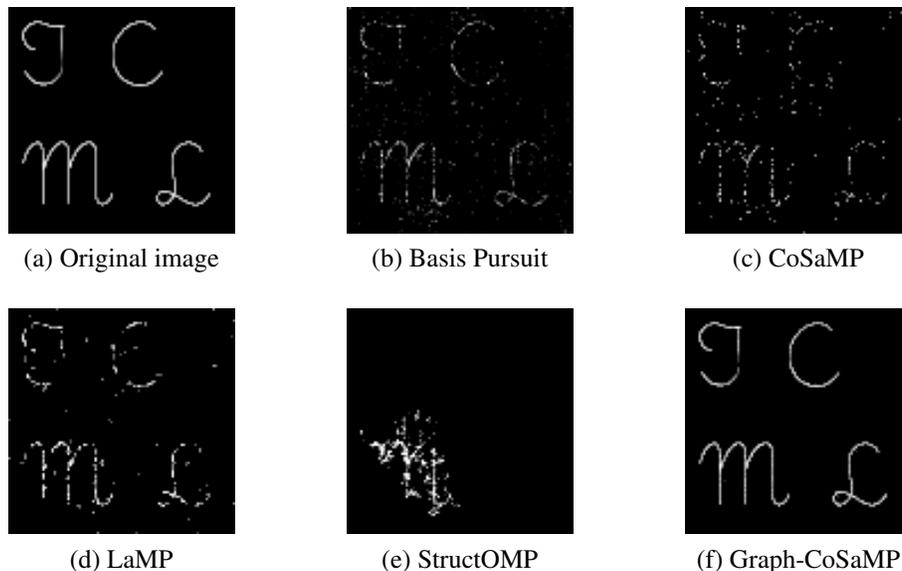


Figure 3. Recovery examples for the text image (see Figure 2) and $n = 3.3s$ noisy linear observations using different recovery algorithms. Only Graph-CoSaMP is able to recover the image correctly.

A. Further experimental results

We start with a more detailed description of our experimental setup. All three images used in Section 6 (Figure 2) are grayscale images of dimension 100×100 pixels with sparsity around 4% to 6%. The background-subtracted image was also used for the experimental evaluation in (Huang et al., 2011). The angiogram image is a slightly sparsified version of the image on the Wikipedia page about angiograms;⁶ it shows cerebral blood vessels. The text image was created by us.

We used SPGL1⁷ as implementation of Basis Pursuit. The implementation of CoSaMP was written by us, closely following (Needell & Tropp, 2009). Graph-CoSaMP and CoSaMP share the same code, only the projection methods differ (hard s -thresholding for CoSaMP and our model projections for Graph-CoSaMP). Empirically it is not necessary to “boost” the head-approximation algorithm as strongly as suggested by the analysis in (Hegde et al., 2014a), we use only a single approximate model projection in place of HEADAPPROX’ (see Alg. 2). The timing experiments in Figure 2(d) were conducted on a Windows machine with a 2.30 GHz Intel Core i7 CPU, 8 MB of cache, and 32 GB of RAM.

Recovered images In order to illustrate the outcomes of unsuccessful recovery trials, we show examples in the regime where Graph-CoSaMP recovers correctly but the other algorithms fail. This is the most relevant regime because it demonstrates that Graph-CoSaMP accurately recovers the image while other methods still show significant errors. See Figure 3 for the corresponding results.

Noise tolerance We also investigate the performance of the recovery algorithms in the noisy setting (the error term e in (2)). For this, we add Gaussian noise at a measurement-SNR level of roughly 15dB. Since we cannot hope for exact recovery in the noisy setting, we consider different tolerance levels for declaring a trial as successful (the ratio $\|\beta - \hat{\beta}\|^2 / \|\beta\|^2$). Figure 4 contains the phase transition plots for the text image from Figure 2(c). The results show that our algorithm also gives the best performance for noisy observations.

Graph Lasso Next, we compare our approach to the graph Lasso introduced in (Jacob et al., 2009). Since the implementation in the SPArse Modeling toolbox (SPAMS)⁸ focuses on dense design matrices, we limit our experiments to a smaller image than those in Figure 2. In particular, we use a 30×30 pixel synthetic image similar to the experiment in Section 9.3

⁶http://commons.wikimedia.org/wiki/File:Cerebral_angiography,_arteria_vertebralis_sinister_injection.JPG

⁷<https://www.math.ucdavis.edu/~mpf/spgl1/>

⁸<http://spams-devel.gforge.inria.fr/index.html>

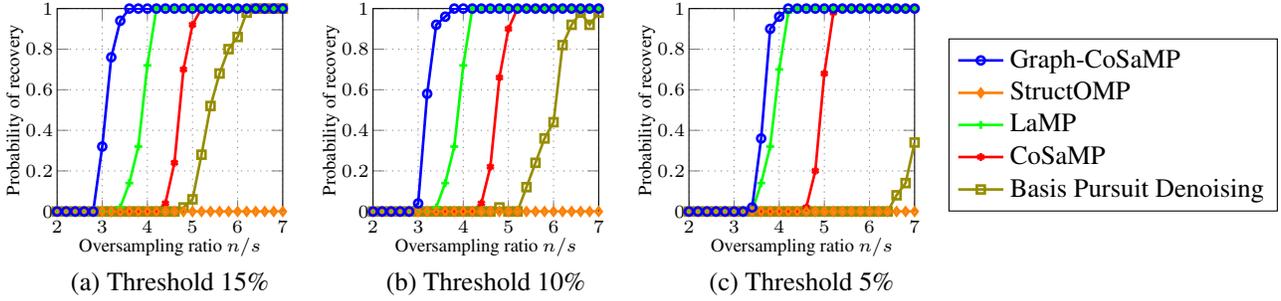


Figure 4. Phase transitions for successful recovery under noisy observations. The three plots are for the same image (the text image from Fig. 2 (c)) but use different thresholds for declaring a trial as successful (the ratio $\|\beta - \hat{\beta}\|^2 / \|\beta\|^2$). Our algorithm offers the best performance for all thresholds.

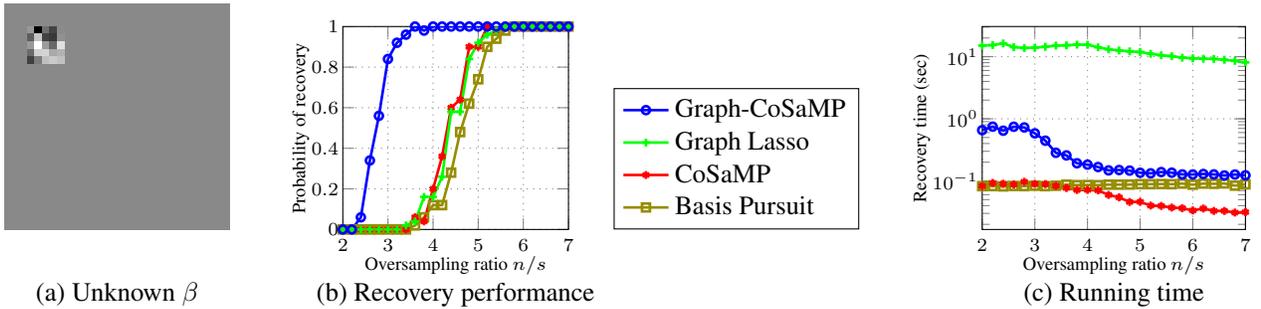


Figure 5. Comparison of our algorithm Graph-CoSaMP with the graph Lasso. Subfigure (a) shows the synthetic test image (30×30 pixels). Graph-CoSaMP recovers the vector β from significantly fewer measurements than the other approaches (phase transition plot (b)). Moreover, Graph-CoSaMP is significantly faster than the variable replication implementation of the graph Lasso and essentially matches the performance of Basis Pursuit in the regime where both algorithms succeed ($n/s \geq 5$ in subfigure (c)).

of (Jacob et al., 2009). The nonzeros form a 5×5 square and hence correspond to a single component in the underlying grid graph. As suggested in (Jacob et al., 2009), we encode the graph structure by using all 4-cycles as groups and use the variable replication approach to implement the overlapping group penalty.

We record n observations $y = X\beta$ with an i.i.d. Gaussian design matrix and follow the experimental procedure outlined in Section 6 (recovery threshold 5%, 50 trials per data point). See Figure 5 for our results. While the graph Lasso improves over Basis Pursuit, our algorithm Graph-CoSaMP recovers the unknown vector β from significantly fewer observations. Moreover, our algorithm is significantly faster than this implementation of the graph Lasso via variable replication.⁹ While there are faster algorithms for the overlapping group Lasso such as (Mosci et al., 2010), the recovery performance of the graph Lasso only matches Graph-CoSaMP for $n/s \geq 5$. In this regime, Graph-CoSaMP is already almost as fast as an efficient implementation of Basis Pursuit (SPGL1).

B. Sparse recovery with the WGM

We now give proofs for theorems in Section 5. First, we establish our general sample complexity bound.

Theorem 13. *Let \mathbb{M} be the set of supports in the (G, s, g, B) -WGM. Then*

$$\log |\mathbb{M}| = O\left(s \left(\log \rho(G) + \log \frac{B}{s}\right) + g \log \frac{d}{g}\right).$$

Proof. Note that every support in the WGM corresponds to a g -forests, which contains exactly $s - g$ edges. We prove the theorem by counting the possible locations of g tree roots in the graph G , and then the local arrangements of the $s - g$

⁹As suggested by the documentation of the SPAMS toolbox, we ran this set of experiments under Linux. The corresponding machine has an Intel Core 2 Duo CPU with 2.93 GHz, 3 MB of cache, and 8 GB of RAM.

edges in the g trees.

Consider the following process:

1. Choose g root nodes out of the entire graph. There are $\binom{d}{g}$ possible choices.
2. Consider the $s - g$ edges as an ordered list and distribute the total weight budget B to the edges. There are $\binom{B+s-g-1}{s-g}$ possible allocations.
3. Assign a “target index” $t_e \in [\rho(G)]$ to each edge. There are $\rho(G)^{s-g}$ possible assignments. Note that the combination of edge weight and target index uniquely determines a neighbor of a fixed node v because there are at most $\rho(G)$ neighbors of v connected with edges of the same weight.
4. We now split the list of edges (together with their weight and target index assignments) into s sets. There are $\binom{2s-g}{s-1}$ possible partitions of the edge list.

We now have a list L consisting of s edge sets together with weight assignments and target indices. Moreover, we have a list of root nodes. We convert this structure to a g -forest (and hence a support in the WGM) according to the following rules, which essentially form a breadth-first search:

While there is a remaining root node, repeat the following:

1. Add the root node to a queue Q .
2. Initialize a new empty tree T_i .
3. While Q is non-empty, repeat the following
 - (a) Let v be the first node in Q and remove v from Q .
 - (b) Add v to T_i .
 - (c) Let A be the first edge set in L and remove A from L .
 - (d) For each pair of target index and weight in A , add the corresponding neighbor to Q .

Note that this process does not always succeed: for some weight allocations, there might be no neighbor connected by an edge with the corresponding weight. Nevertheless, it is easy to see that every possible support in the (G, s, g, B) -WGM can be constructed from at least one allocation via the process described above. Hence we have a surjection from the set of allocations to supports in the (G, s, g, B) -WGM \mathbb{M} , which gives the following bound:

$$|\mathbb{M}| \leq \binom{B+s-g-1}{s-g} \cdot \rho^{s-g} \cdot \binom{2s+g}{s-1} \cdot \binom{d}{g}.$$

Taking a logarithm on both sides and simplifying yields the bound in the theorem. □

The proof of the recovery result in Theorem 12 directly follows by combining the guarantees established for our tail- and head-approximation algorithms (Theorems 9 and 11) with the framework of (Hegde et al., 2014a).

Theorem 12. *Let $\beta \in \mathbb{R}^d$ be in the (G, s, g, B) -WGM \mathcal{M} and let $X \in \mathbb{R}^{n \times d}$ be a matrix satisfying the model-RIP for a $(G, c_1 s, g, c_2 B)$ -WGM and a fixed constant δ , where c_1 and c_2 are also fixed constants. Moreover, let $e \in \mathbb{R}^n$ be an arbitrary noise vector and let $y \in \mathbb{R}^n$ be defined as in (2). Then GRAPH-COSAMP returns a $\hat{\beta}$ in the $(G, 5s, g, 5B)$ -WGM such that $\|\beta - \hat{\beta}\| \leq c_3 \|e\|$, where c_3 is a fixed constant. Moreover, GRAPH-COSAMP runs in time*

$$O\left((T_X + |E| \log^3 d) \log \frac{\|\beta\|}{\|e\|}\right),$$

where T_X is the time complexity of a matrix-vector multiplication with X .

Proof. Note that both our head- and tail-approximation algorithms project into an output model with parameters bounded by constant multiples of s and B (we always maintain that the support corresponds to a g -forest), see Theorems 9 and 11. This allows us to use the CoSaMP version of Corollary 19 in (Hegde et al., 2014a) to establish the recovery result in our theorem. The claim about the running time follows from the near-linear running time of our model-projection algorithms and the running time analysis of CoSaMP in (Needell & Tropp, 2009). The $\log \frac{\|\beta\|}{\|\epsilon\|}$ term in the running time comes from the geometric convergence of Graph-CoSaMP. \square

C. Approximate model-projection algorithms for the WGM

We now formally prove the head- and tail-approximation guarantees for our model-projection algorithms. We assume that we have access to an algorithm PCSF-GW for the PCSF problem with the approximation guarantee from Theorem 5, which we restate for completeness:

Theorem 5. *There is an algorithm for the PCSF problem that returns a g -forest F such that*

$$c(F) + 2\pi(\overline{F}) \leq \min_{F' \subseteq G, \gamma(F') \leq g} 2c(F') + 2\pi(\overline{F'}) . \quad (11)$$

We denote the running time of PCSF-GW with T_{PCSF} . See Section D for an algorithm that achieves guarantee (11) in nearly-linear time.

C.1. Tail-approximation

We first address the special case that there is a g -forest F^* with $c(F^*) \leq C$ and $\pi(\overline{F^*}) = 0$. In this case, we have to find a g -forest F with $\pi(\overline{F}) = 0$ in order to satisfy (12).

Lemma 14. *Let $\pi_{\min} = \min_{\pi(v) > 0} \pi(v)$ and $\lambda_0 = \frac{\pi_{\min}}{2C}$. If there is a g -forest F^* with $c(F^*) \leq C$ and $\pi(\overline{F^*}) = 0$, then PCSF-GW(G, c_{λ_0}, π, g) returns a g -forest F with $c(F) \leq 2C$ and $\pi(\overline{F}) = 0$.*

Proof. Applying the GW guarantee (11) gives

$$\begin{aligned} \lambda_0 \cdot c(F) + 2\pi(\overline{F}) &\leq 2\lambda_0 \cdot c(F^*) + 2\pi(\overline{F^*}) \\ \pi(\overline{F}) &\leq \lambda_0 C = \frac{\pi_{\min}}{2} . \end{aligned}$$

Since $\pi_{\min} > 0$, we must have $\pi(\overline{F}) < \pi_{\min}$ and hence $\pi(\overline{F}) = 0$.

Applying (11) again then gives $c_{\lambda_0}(F) \leq 2c_{\lambda_0}(F^*)$, which shows that $c(F) \leq 2c(F^*) \leq 2C$ as desired. \square

We can now proceed to prove an approximation guarantee for PCSF-TAIL.

Theorem 8. *Let $\nu > 2$ and $\delta > 0$. Then PCSF-TAIL returns a g -forest $F \subseteq G$ such that $c(F) \leq \nu \cdot C$ and*

$$\pi(\overline{F}) \leq \left(1 + \frac{2}{\nu - 2} + \delta\right) \min_{\gamma(F') = g, c(F') \leq C} \pi(\overline{F'}) . \quad (12)$$

Proof. We consider the three different cases in which PCSF-TAIL returns a forest. Note that the resulting forest is always the output of PCSF-GW with parameter g , so the resulting forest is always a g -forest. To simplify notation, in the following we use

$$OPT = \min_{\gamma(F') = g, c(F') \leq C} \pi(\overline{F'}) .$$

First, if PCSF-TAIL returns in Line 5, the forest F directly satisfies (12). Otherwise, there is no g forest F^* with $c(F^*) \leq C$ and $\pi(\overline{F^*}) = 0$ (contrapositive of Lemma 14). Hence in the following we can assume that $OPT \geq \pi_{\min}$.

If the algorithm returns in Line 10, we clearly have $c(F) \leq \nu \cdot C$. Moreover, the GW guarantee gives

$$\lambda_m \cdot c(F) + 2\pi(\overline{F}) \leq 2\lambda_m C + 2 \cdot OPT .$$

Since $c(F) \geq 2C$, we have $\pi(\overline{F}) \leq OPT$, satisfying (12).

Finally, consider the case that PCSF-TAIL returns in Line 13. Let F_l and F_r be the forests corresponding to λ_l and λ_r , respectively. We show that the final output F_l satisfies the desired approximation guarantee if $\lambda_r - \lambda_l$ is small. Note that during the binary search, we always maintain the invariant $c(F_l) \leq 2C$ and $c(F_r) \geq \nu \cdot C$.

Using the GW guarantee and $\pi(\overline{F_r}) \geq 0$ gives $\lambda_r c(F_r) \leq 2\lambda_r C + 2 \cdot OPT$. Therefore,

$$\lambda_r \leq \frac{2 \cdot OPT}{c(F_r) - 2C} \leq \frac{2 \cdot OPT}{C(\nu - 2)}. \quad (14)$$

At the end of the binary search, we have $\lambda_l \leq \lambda_r + \varepsilon$. Combining this with (14) above and the GW guarantee (11) gives

$$\pi(\overline{F}) \leq \lambda_l C + OPT \leq OPT + (\lambda_r + \varepsilon)C \leq OPT + \frac{2 \cdot OPT}{\nu - 2} + \varepsilon C \leq \left(1 + \frac{2}{\nu - 2} + \delta\right) OPT.$$

In the last inequality, we used $OPT \geq \pi_{\min}$ and $\varepsilon = \frac{\pi_{\min} \delta}{C}$. This concludes the proof. \square

Finally, we consider the running time of PCSF-TAIL.

Theorem 15. PCSF-TAIL runs in time $O(T_{\text{PCSF}} \cdot \log \frac{C \cdot \pi(G)}{\delta \cdot \pi_{\min}})$.

Proof. The time complexity is dominated by the number of calls to PCSF-GW. Hence we bound the number of binary search iterations in order to establish the overall time complexity. Let $\lambda^{(0)}$ be the initial value of λ_l in PCSF-TAIL. Then the maximum number of iterations of the binary search is

$$\left\lceil \log \frac{\lambda_l^{(0)}}{\varepsilon} \right\rceil = \left\lceil \log \frac{3\pi(G) \cdot C}{\delta \cdot \pi_{\min}} \right\rceil = O\left(\log \frac{C \cdot \pi(G)}{\delta \cdot \pi_{\min}}\right).$$

Since each iteration of the binary search takes $O(T_{\text{PCSF}})$ time, the time complexity stated in the theorem follows. \square

If the node prizes π and edge costs c are polynomially bounded in $|V|$, the running time of PCSF-TAIL simplifies to $O(T_{\text{PCSF}} \cdot \log |V|)$ for constant δ .

We now have all results to complete our tail-approximation algorithm for the WGM.

Theorem 9. Let \mathcal{M} be a (G, s, g, B) -WGM, let $b \in \mathbb{R}^d$, and let $\nu > 2$. Then there is an algorithm that returns a support $S \subseteq [d]$ in the $(G, 2\nu \cdot s + g, g, 2\nu \cdot B)$ -WGM satisfying (8) with $c_T = \sqrt{1 + 3/(\nu - 2)}$. Moreover, the algorithm runs in time $O(|E| \log^3 d)$.

Proof. We run the algorithm PCSF-TAIL on the graph G with node prizes $\pi(i) = b_i^2$, edge costs $c(e) = w(e) + \frac{B}{s}$, a cost budget $C = 2B$, and the parameter $\delta = \min(\frac{1}{2}, \frac{1}{\nu})$. Let F be the resulting forest and S the corresponding support. The running time bound follows from combining Theorems 15 and 28.

First, we show that S is in the $(G, 2\nu \cdot s + g, g, 2\nu \cdot B)$ -WGM. From Theorem 8 we know that F is a g -forest and that $c(F) \leq 2\nu \cdot B$. This directly implies that $w(F) \leq 2\nu \cdot B$. Moreover, the g -forest F has $|V_F| - g$ edges, all with cost at least $\frac{B}{s}$ because $w(e) \geq 0$ for all $e \in E$. Since $|V_F| = |S|$, this allows us to bound the sparsity of S as

$$(|S| - g) \cdot \frac{B}{s} \leq 2\nu \cdot B,$$

which gives $|S| \leq 2s + g$ as desired.

Now, let S^* be an optimal support in the (G, s, g, B) -WGM \mathbb{M} and let F^* be a corresponding g -forest, i.e.,

$$\pi(\overline{F^*}) = \|b - b_{S^*}\|^2 = \min_{S' \in \mathbb{M}} \|b - b_{S'}\|^2.$$

Then we have

$$\pi(\overline{F^*}) \geq \min_{\gamma(F')=g, c(F') \leq 2B} \pi(\overline{F'})$$

because by construction, every support in \mathbb{M} corresponds to a g -forest with cost at most $2B$. Since $\pi(\bar{F}) = \|b - b_S\|^2$, applying guarantee (12) gives

$$\|b - b_S\|^2 \leq \left(1 + \frac{2}{\nu - 2} + \delta\right) \min_{S' \in \mathbb{M}} \|b - b_{S'}\|^2.$$

Simplifying this inequality with our choice of δ then completes the proof. \square

C.2. Head-approximation

We first state our head-approximation algorithm (see Alg. 3 and Alg. 4). In addition to a binary search over the Lagrangian parameter λ , the algorithm also uses the subroutines PRUNETREE and PRUNEFORREST in order to extract sub-forests with good “density” $\frac{\pi(F)}{c(F)}$.

Algorithm 3 Head approximation for the WGM: main algorithm PCSF-HEAD

```

1: function PCSF-HEAD( $G, c, \pi, g, C, \delta$ )
2:   We write  $\pi_\lambda(i) = \lambda \cdot \pi(i)$ .
3:    $\pi_{\min} \leftarrow \min_{\pi(i) > 0} \pi(i)$ 
4:    $\lambda_r \leftarrow \frac{2C}{\pi_{\min}}$ 
5:    $F \leftarrow \text{PCSF-GW}(G, c, \pi_{\lambda_r}, g)$ 
6:   if  $c(F) \leq 2C$  then                                      $\triangleright$  Ensure that we have the invariant  $c(F_r) > 2C$  (see Theorem 17)
7:     return  $F$ 
8:   end if
9:    $\varepsilon \leftarrow \frac{\delta \cdot C}{2\pi(G)}$ 
10:   $\lambda_l \leftarrow \frac{1}{4\pi(G)}$ 
11:  while  $\lambda_r - \lambda_l > \varepsilon$  do                                      $\triangleright$  Binary search over the Lagrange parameter  $\lambda$ 
12:     $\lambda_m \leftarrow (\lambda_l + \lambda_r)/2$ 
13:     $F \leftarrow \text{PCSF-GW}(G, c, \pi_{\lambda_m}, g)$ 
14:    if  $c(F) > 2C$  then
15:       $\lambda_r \leftarrow \lambda_m$ 
16:    else
17:       $\lambda_l \leftarrow \lambda_m$ 
18:    end if
19:  end while
20:   $F_l \leftarrow \text{PCSF-GW}(G, c, \pi_{\lambda_l}, g)$ 
21:   $F_r \leftarrow \text{PCSF-GW}(G, c, \pi_{\lambda_r}, g)$ 
22:   $F'_r \leftarrow \text{PRUNEFORREST}(F, c, \pi, C)$                                       $\triangleright$  Prune the potentially large solution  $F_r$  (See Alg. 4)
23:  if  $\pi(F_l) \geq \pi(F'_r)$  then
24:    return  $F_l$ 
25:  else
26:    return  $F'_r$ 
27:  end if
28: end function

```

We start our analysis by showing that PRUNETREE extracts sub-trees of good density $\frac{\pi(T')}{c(T')}$.

Lemma 10. *Let T be a tree and $C' \leq c(T)$. Then there is a subtree $T' \subseteq T$ such that $c(T') \leq C'$ and $\pi(T') \geq \frac{C'}{6} \cdot \frac{\pi(T)}{c(T)}$. Moreover, we can find such a subtree T' in linear time.*

Proof. We show that PRUNETREE satisfies the guarantees in the theorem. We use the definitions of L , π' , c' , and ϕ given in PRUNETREE (see Alg. 4). Moreover, let T' be the tree returned by PRUNETREE. First, note that PRUNETREE clearly runs in linear time by definition. Hence it remains to establish the approximation guarantee

$$\pi(T') \geq \frac{C'}{6} \cdot \frac{\pi(T)}{c(T)} = \frac{C'}{6} \cdot \phi.$$

Algorithm 4 Head approximation for the WGM: subroutine PRUNEFORREST

```

1: function PRUNEFORREST( $F, c, \pi, C$ )
2:   Let  $\{T_1, \dots, T_{|F|}\}$  be the trees in  $F$  sorted by  $\frac{\pi(T_i)}{c(T_i)}$  descendingly.
3:    $C_r \leftarrow C$ 
4:   for  $i \leftarrow 1, \dots, |F|$  do
5:     if  $C_r \geq c(T_i)$  then
6:        $T'_i \leftarrow T_i$ 
7:        $C_r \leftarrow C_r - c(T_i)$  ▷ Cost budget  $C^{(i)} = c(T_i)$ 
8:     else if  $C_r > 0$  then
9:        $T'_i \leftarrow \text{PRUNETREE}(T_i, c, \pi, C_r)$ 
10:       $C_r \leftarrow 0$  ▷ Cost budget  $C^{(i)} = C_r$ 
11:     else
12:        $T'_i \leftarrow \{\arg \max_{j \in T_i} \pi(j)\}$  ▷ Cost budget  $C^{(i)} = 0$ 
13:     end if
14:   end for
15:   return  $\{T'_1, \dots, T'_{|F|}\}$ 
16: end function

17: function PRUNETREE( $T, c, \pi, C'$ )
18:   Let  $L = (v_1, \dots, v_{2|V_T|-1})$  be a tour through the nodes of  $T$ . ▷  $T = (V_T, E_T)$ 
19:   Let  $\pi'(j) = \begin{cases} \pi(v_j) & \text{if position } j \text{ is the first appearance of } v_j \text{ in } L \\ 0 & \text{otherwise} \end{cases}$ 
20:   Let  $c'(P) = \sum_{i=1}^{|P|-1} c(P_i, P_{i+1})$ 
21:   Let  $\phi = \frac{\pi(T)}{c(T)}$ 
22:   if there is a  $v \in V_T$  with  $\pi(v) \geq \frac{C' \cdot \phi}{6}$  then ▷ Check if there is a single good node (cost is automatically 0)
23:     return the tree  $\{v\}$ 
24:   end if
25:    $l \leftarrow 1$ 
26:    $P^1 = ()$  ▷ Empty list
27:   for  $i \leftarrow 1, \dots, 2|V_T| - 1$  do ▷ Search for good sublists of  $L$ 
28:     Append  $i$  to  $P^l$ 
29:     if  $c'(P^l) > C'$  then ▷ Start a new sublist if the cost reaches  $C'$ 
30:        $l \leftarrow l + 1$ 
31:        $P^l \leftarrow ()$ 
32:     else if  $\pi'(P^l) \geq \frac{C' \cdot \phi}{6}$  then ▷ Return if we have found a good sublist
33:       return the subtree of  $T$  on the nodes in  $P^l$ 
34:     end if
35:   end for
36:   Merge  $P^l$  and  $P^{l-1}$  ▷ The algorithm will never reach this point (see Lemma 10).
37: end function

```

Consider the first case in which PRUNETREE returns in line 23. Then T' is a tree consisting of a single node, so $c(T') = 0 \leq C'$. Moreover, we have $\pi(T') = \pi(v) \geq \frac{C' \cdot \phi}{6}$, which satisfies the guarantee in the theorem.

Next, we consider the case in which PRUNETREE returns in line 33. By definition of the algorithm, we have $c'(P^l) \leq C'$ and hence $c(T') \leq C'$ because the spanning tree T' of the nodes in P^l contains only edges that are also included at least once in $c'(P^l)$. Moreover, we have $\pi(T') \geq \pi'(P^l) \geq \frac{C' \cdot \phi}{6}$, so T' satisfies the guarantee in the theorem.

It remains to show that PRUNETREE always returns in one of the two cases above, i.e., never reaches line 36. We prove this statement by contradiction: assume that PRUNETREE reaches line 36. We first consider the partition of V_T induced by the lists P^i just before line 36. Note that there are no nodes $v \in V_T$ with $\pi(v) \geq \frac{C' \cdot \phi}{6}$ because otherwise PRUNETREE would have returned in line 23. Hence for every list P^i we have $\pi'(P^i) \leq \frac{C' \cdot \phi}{3}$ because the last element that was added to P^i can have increased $\pi'(P^i)$ by at most $\frac{C' \cdot \phi}{6}$, and we had $\pi'(P^i) < \frac{C' \cdot \phi}{6}$ before the last element was added to P^i because otherwise PRUNETREE would have returned in line 33. Moreover, every list P^i except P^l satisfies $c'(P^i) > C'$ by construction. Hence after merging the last two lists P^l and P^{l-1} , we have $c'(P^i) > C'$ for all P^i and also $\pi'(P^i) < \frac{C' \cdot \phi}{2}$.

We now derive the contradiction: note that all lists P^i have a low density $\frac{\pi'(P^i)}{c'(P^i)}$ but form a partition of the nodes in V_T . We can use this fact to show that the original tree had a density lower than $\frac{\pi(T)}{c(T)}$, which is a contradiction. More formally, we have

$$\pi(T) = \sum_{i=1}^{l-1} \pi'(P^i) < (l-1) \frac{C' \cdot \phi}{2}$$

and

$$2c(T) \geq \sum_{i=1}^{l-1} c'(P^i) > (l-1)C'.$$

Combining these two inequalities gives

$$\frac{\phi}{2} = \frac{\pi(T)}{2c(T)} < \frac{(l-1) \frac{C' \cdot \phi}{2}}{(l-1)C'} = \frac{\phi}{2},$$

which is a contradiction. Hence PRUNETREE always returns in line 23 or 33 and satisfies the guarantee of the theorem. \square

Extending the guarantee of PRUNETREE to forests is now straightforward: we can prune each tree in a forest F individually by assigning the correct cost budget to each tree. More formally, we get the following lemma.

Lemma 16. *Let F be a g -forest. Then $\text{PRUNEFORREST}(F, c, \pi, C)$ returns a g -forest F' with $c(F') \leq C$ and*

$$\pi(F') \geq \frac{C}{6 \cdot c(F)} \pi(F).$$

Proof. By construction, F' is a g -forest with $c(F') \leq C$. Let $C^{(i)}$ be the cost budget assigned to tree T'_i (see the comments in PRUNEFORREST). Using Lemma 10, we get

$$\pi(F') = \sum_{i=1}^g \pi(T'_i) \geq \sum_{i=1}^g \frac{C^{(i)}}{6} \frac{\pi(T_i)}{c(T_i)}.$$

Note that the $C^{(i)}$ are the optimal allocation of budgets to the ratios $\frac{\pi(T_i)}{c(T_i)}$ with $0 \leq C^{(i)} \leq c(T_i)$ and $\sum_{i=1}^g C^{(i)} = C$. In particular, we have

$$\sum_{i=1}^g \frac{C^{(i)}}{6} \frac{\pi(T_i)}{c(T_i)} \geq \sum_{i=1}^g \frac{C \cdot \frac{c(T_i)}{c(F)}}{6} \frac{\pi(T_i)}{c(T_i)} = \frac{C}{6 \cdot c(F)} \pi(F),$$

which completes the proof. \square

We can now prove our main theorem about PCSF-HEAD.

Theorem 17. Let $0 < \delta < \frac{1}{13}$. Then PCSF-HEAD returns a g -forest F such that $c(F) \leq 2C$ and

$$\pi(F) \geq \left(1 - \frac{12}{13(1-\delta)}\right) \max_{\gamma(F')=g, c(F') \leq C} \pi(F'). \quad (15)$$

Proof. Let F^* be an optimal g -forest with $c(F^*) \leq C$ and $\pi(F^*) = OPT$, where

$$OPT = \max_{\gamma(F')=g, c(F') \leq C} \pi(F').$$

In this proof, the following rearranged version of the GW guarantee 11 will be useful:

$$\begin{aligned} c(F) + 2(\pi(G) - \pi(F)) &\leq 2C + 2(\pi(G) - \pi(F^*)) \\ \pi(F) &\geq \pi(F^*) + \frac{c(F) - 2C}{2}. \end{aligned} \quad (16)$$

As in the definition of PCSF-HEAD, we write π_λ for the node prize function $\pi_\lambda(i) = \lambda \cdot \pi(i)$. Using such modified node prizes, (16) becomes

$$\pi(F) \geq OPT + \frac{c(F) - 2C}{2\lambda}. \quad (17)$$

We now analyze two cases: either PCSF-HEAD returns in line 7 or in one of the lines 24 and 26. Note that in all cases, the returned forest F is a g -forest because it is produced by PCSF-GW (and PRUNEFORREST maintains this property).

First, we consider the case that the algorithm returns in line 7. Then by definition we have $c(F) \leq 2C$. Moreover, the modified GW guarantee (17) gives

$$\pi(F) \geq OPT + \frac{c(F) - 2C}{2\lambda_r} \geq OPT - \frac{C}{\lambda_r} \geq OPT - \frac{\pi_{\min}}{2} \geq \frac{1}{2} OPT,$$

because clearly $OPT \geq \pi_{\min}$. Hence the guarantee in the theorem is satisfied.

Now, consider the case that the algorithm enters the binary search. Let F_l and F_r be the g -forests corresponding to λ_l and λ_r , respectively. During the binary search, we maintain the invariant that $c(F_l) \leq 2C$ and $c(F_r) > 2C$. Note that our initial choices for λ_l and λ_r satisfy this condition (provided the algorithm reaches the binary search).

When the algorithm terminates in line 24 or 26, we have $\lambda_r > \lambda_l > \lambda_r - \varepsilon$. Rearranging (17) gives

$$\begin{aligned} 2\lambda_r(\pi(F_r) - OPT) &\geq c(F_r) - 2C \\ \lambda_r &\geq \frac{c(F_r) - 2C}{2(\pi(F_r) - OPT)}. \end{aligned}$$

We now introduce a variable $\alpha \geq 2$ and distinguish two cases:

Case 1: Assume $\frac{c(F_r)}{\pi(F_r)} > \alpha \frac{C}{OPT}$. Then Equation (17) gives

$$\begin{aligned} \lambda_r &\geq \frac{\frac{1}{2} \frac{c(F_r)}{OPT} - \frac{C}{OPT}}{\frac{\pi(F_r)}{OPT} - 1} \\ &= \frac{\frac{1}{2} \frac{\pi(F_r)}{OPT} \frac{c(F_r)}{\pi(F_r)} - \frac{C}{OPT}}{\frac{\pi(F_r)}{OPT} - 1} \\ &\geq \frac{\frac{1}{2} \frac{\pi(F_r)}{OPT} \alpha - 1}{\frac{\pi(F_r)}{OPT} - 1} \frac{C}{OPT} \\ &\geq \frac{\frac{\alpha}{2} \frac{\pi(F_r)}{OPT} - \frac{\alpha}{2}}{\frac{\pi(F_r)}{OPT} - 1} \frac{C}{OPT} \\ &= \frac{\alpha}{2} \frac{C}{OPT}. \end{aligned}$$

So we get $\lambda_l \geq \lambda_r - \varepsilon \geq \frac{\alpha}{2} \frac{C}{OPT} - \frac{\delta \cdot C}{2\pi(G)} \geq (1 - \delta) \frac{\alpha C}{2OPT}$. We can now use this together with (17) to get:

$$\begin{aligned} \pi(F_l) &\geq OPT - \frac{C}{\lambda_l} \\ &\geq OPT - \frac{2OPT}{(1 - \delta)\alpha} \\ &= \left(1 - \frac{2}{(1 - \delta)\alpha}\right) OPT. \end{aligned} \tag{18}$$

Case 2: Assume $\frac{c(F_r)}{\pi(F_r)} \leq \alpha \frac{C}{OPT}$, which is equivalent to $\frac{\pi(F_r)}{c(F_r)} \geq \frac{1}{\alpha} \frac{OPT}{C}$. Since $c(F_r) > 2C$, we can invoke PRUNEFOR-EST on F_r with cost budget C . Let F'_r be the resulting g -forest. From Lemma 16 we have $c(F'_r) \leq C$. Moreover,

$$\pi(F'_r) \geq \frac{C}{6 \cdot c(F_r)} \pi(F_r) = \frac{C}{6} \frac{\pi(F_r)}{c(F_r)} \geq \frac{1}{6\alpha} OPT. \tag{19}$$

Either case 1 or case 2 must hold. Since HEADAPPROX chooses the better forest among F_l and F'_r , we can combine Equations (18) and (19) to get the following guarantee on the final result F :

$$\pi(F) \geq \min\left(1 - \frac{2}{(1 - \delta)\alpha}, \frac{1}{6\alpha}\right) OPT.$$

Choosing $\alpha = \frac{13}{6}$ to balance the two expressions (assuming δ is close to 0) then gives the approximation guarantee stated in the theorem. \square

Next, we consider the running time of PCSF-HEAD.

Theorem 18. PCSF-HEAD runs in time $O\left(T_{\text{PCSF}} \cdot \log \frac{\pi(G)}{\delta \cdot \pi_{\min}}\right)$.

Proof. As in Theorem 15 it suffices to bound the number of iterations of the binary search. Let $\lambda_r^{(0)}$ be the initial value of λ_r in PCSF-HEAD. Then the maximum number of iterations is

$$\left\lceil \log \frac{\lambda_r^{(0)}}{\varepsilon} \right\rceil = \left\lceil \log \frac{4 \cdot C \cdot \pi(G)}{\delta \cdot C \cdot \pi_{\min}} \right\rceil = O\left(\frac{\pi(G)}{\delta \cdot \pi_{\min}}\right).$$

\square

As before, the running time simplifies to $O(T_{\text{PCSF}} \cdot \log|V|)$ for constant δ if the node prizes and edge costs are polynomially bounded in the size of the graph.

We can now conclude with our head-approximation algorithm for the WGM.

Theorem 11. Let \mathcal{M} be a (G, s, g, B) -WGM and let $b \in \mathbb{R}^d$. Then there is an algorithm that returns a support $S \subseteq [d]$ in the $(G, 2s + g, g, 2B)$ -WGM satisfying (9) with $c_H = \sqrt{1/14}$. The algorithm runs in time $O(|E| \log^3 d)$.

Proof. We embed the WGM into a PCSF instance similar to Theorem 9: we run PCSF-HEAD on the graph G with node prizes $\pi(i) = b_i^2$, edge costs $c(e) = w(e) + \frac{B}{s}$, a cost budget $C = 2B$, and the parameter $\delta = \frac{1}{169}$. Let F be the resulting forest and S be the corresponding support. The running time bound follows from combining Theorems 18 and 28.

From Theorem 17 we directly have that F is a g -forest with $w(F) \leq 2B$. Following a similar argument as in Theorem 9, we also get $|S| \leq 2s + g$. So S is in the $(G, 2s + g, g, 2B)$ -WGM.

Now, let S^* be an optimal support in the (G, s, g, B) -WGM \mathbb{M} and let F^* be a corresponding g -forest, i.e.,

$$\pi(F^*) = \|b_{S^*}\|^2 = \max_{S' \in \mathbb{M}} \|b_{S'}\|^2.$$

By construction, every support in \mathbb{M} corresponds to a g -forest with cost at most $2B$. Hence we have

$$\pi(F^*) \leq \max_{\gamma(F')=g, c(F') \leq 2B} \pi(F')$$

Since $\pi(F) = \|b_S\|^2$, applying Theorem 17 gives

$$\|b_S\|^2 \geq \left(1 - \frac{12}{13(1-\delta)}\right) \max_{S' \in \mathbb{M}} \|b_{S'}\|^2.$$

Substituting $\delta = \frac{1}{169}$ completes the proof. \square

D. The prize-collecting Steiner forest problem (PCSF)

For completeness, we first review the relevant notation and the definition of the PCSF problem. Let $G = (V, E)$ be an undirected, weighted graph with edge costs $c : E \rightarrow \mathbb{R}_0^+$ and node prizes $\pi : V \rightarrow \mathbb{R}_0^+$. For a subset of edges $E' \subseteq E$, we write $c(E') = \sum_{e \in E'} c(e)$ and adopt the same convention for node subsets. Moreover, for a node subset $V' \subseteq V$, let $\overline{V'}$ be the complement $\overline{V'} = V \setminus V'$. We denote the number of connected components in the (sub-)graph F with $\gamma(F)$.

Definition 4 (The prize-collecting Steiner forest problem). *Let $g \in \mathbb{N}$ be the target number of connected components. Then the goal of the prize-collecting Steiner forest (PCSF) problem is to find a subgraph $F = (V', E')$ with $\gamma(F) = g$ that minimizes $c(E') + \pi(\overline{V'})$.*

We divide our analysis in two parts: we first modify the Goemans-Williamson (GW) scheme to get an efficient algorithm with provable approximation guarantee for the PCSF problem (Subsection D.1). Then we show how to simulate the GW scheme in nearly-linear time (Subsection D.2).

D.1. The Goemans-Williamson (GW) scheme for PCSF

Before we introduce our variant of the GW scheme and prove the desired approximation guarantee, we introduce additional notation. For a set of nodes $U \subseteq V$ and a set of edges $D \subseteq E$, we write $\delta_D U$ to denote the set of edges contained in D with exactly one endpoint in U . If $D = E$, we write δU . The degree of a node v in an edge set D is $\deg_D(v) = |\delta_D\{v\}|$. We say that a (sub-)graph F is a g -forest if F is a forest with $\gamma(F) = g$.

At its core, the GW algorithm produces three results: a *laminar family* of clusters, a *dual value* for each cluster, and a forest connecting the nodes within each cluster.

Definition 19 (Laminar family). *A family \mathcal{L} of non-empty subsets of V is a laminar family if one of the following three cases holds for all $L_1, L_2 \in \mathcal{L}$: either $L_1 \cap L_2 = \{\}$, or $L_1 \subseteq L_2$, or $L_2 \subseteq L_1$.*

Let U be a subset of V . Then we define the following two subsets of \mathcal{L} :

- $\mathcal{L}_{\downarrow U} := \{L \in \mathcal{L} \mid L \subseteq U\}$ (going “down” in the laminar hierarchy).
- $\mathcal{L}_{\uparrow U} := \{L \in \mathcal{L} \mid U \subseteq L\}$ (going “up” in the laminar hierarchy).

Let \mathcal{L}^* be the family of maximal sets in \mathcal{L} , i.e., $L \in \mathcal{L}^*$ iff there is no $L' \in \mathcal{L}$ with $L \subsetneq L'$. If $\bigcup_{L \in \mathcal{L}} L = V$, then \mathcal{L}^* is a partition of V .

Let $e \in E$, then we write $\mathcal{L}(e) := \{L \in \mathcal{L} \mid e \in \delta L\}$ for the sub-family of sets that contain exactly one endpoint of e .

Definition 20 (Dual values). *Let \mathcal{L} be a laminar family. Then the dual values are a function $y : \mathcal{L} \rightarrow \mathbb{R}_0^+$ with the following two properties (as before, we write $y(\mathcal{L}') := \sum_{L \in \mathcal{L}'} y(L)$ for a sub-family $\mathcal{L}' \subseteq \mathcal{L}$).*

- $y(\mathcal{L}(e)) \leq c(e)$ for each $e \in E$.
- $y(\mathcal{L}_{\downarrow L}) \leq \pi(L)$ for each $L \in \mathcal{L}$.

We also define several properties of g -forests related to the new concepts introduced above.

Let \mathcal{L} be a laminar family. We say a tree T is \mathcal{L} -connected iff for every $L \in \mathcal{L}$, the subgraph on $V_T \cap L$ is connected (we consider an empty graph to be connected). A g -forest F is \mathcal{L} -connected iff every $T \in F$ is \mathcal{L} -connected.

Let $L \in \mathcal{L}^*$ and let $L(F)$ be the trees in F with at least one node in L , i.e., $L(F) = \{T \in F \mid V_T \cap L \neq \{\}\}$. A g -tree F is \mathcal{L}^* -disjoint iff $|L(F)| \leq 1$ for every $L \in \mathcal{L}^*$.

Let \mathcal{D} be a family of subsets of V . A tree T has a leaf component in \mathcal{D} iff there is a $D \in \mathcal{D}$ with $|\delta_T D| = 1$. A g -forest F has a leaf component in \mathcal{D} iff there is a tree $T \in F$ that has a leaf component in \mathcal{D} .

A tree T is contained in \mathcal{D} iff there is a $D \in \mathcal{D}$ such that $V_T \subseteq D$. A g -forest F is contained in \mathcal{D} iff there is a tree $T \in F$ that is contained in \mathcal{D} .

D.1.1.1. ALGORITHM

Our algorithm is a modification of the unrooted GW PCST algorithm in (Johnson et al., 2000). In contrast to their unrooted prize-collecting Steiner tree algorithm, our algorithm stops the growth phase when exactly g active clusters are left. We use these active clusters as starting point in the pruning phase to identify a g -forest as the final result.

Since the focus of this section is the approximation guarantee rather than the time complexity, the pseudo code in Algorithm 5 is intentionally stated at a high level.

Algorithm 5 Prize-collecting Steiner forest

```

1: function PCSF-GW( $G, c, \pi, g$ )
2:    $\mathcal{L} \leftarrow \{\{v\} \mid v \in V\}$                                 ▷ Laminar family of clusters.
3:    $y(C) \leftarrow 0$  for all  $C \in \mathcal{L}$ .                                ▷ Initial dual values.
4:    $V_F \leftarrow V, E_F \leftarrow \{\}$                                 ▷ Initial forest.
5:    $\mathcal{D} \leftarrow \{\}$                                                 ▷ Family of inactive clusters.
6:    $\mathcal{A} \leftarrow \mathcal{L}^* \setminus \mathcal{D}$                                        ▷ Family of active clusters.
7:   while  $|\mathcal{A}| > g$  do                                           ▷ Growth phase.
8:      $\varepsilon_d \leftarrow \min_{C \in \mathcal{A}} \pi(C) - y(\mathcal{L}_{\downarrow C})$            ▷ Next cluster deactivation time
9:      $\varepsilon_m \leftarrow \min_{\substack{e \in \delta C \\ C \in \mathcal{A}}} c(e) - y(\mathcal{L}(e))$        ▷ Next cluster merge time
10:     $\varepsilon \leftarrow \min(\varepsilon_d, \varepsilon_m)$ 
11:    for  $C \in \mathcal{A}$  do
12:       $y(C) \leftarrow y(C) + \varepsilon$                                 ▷ Increase dual variables for active clusters.
13:    end for
14:    if  $\varepsilon_c \leq \varepsilon_m$  then                                       ▷ Cluster deactivation next.
15:      Let  $C \in \mathcal{A}$  be such that  $\pi(C) - y(\mathcal{L}_{\downarrow C}) = 0$ .
16:       $\mathcal{D} \leftarrow \mathcal{D} \cup \{C\}$                                        ▷ Mark cluster as inactive.
17:    else
18:      Let  $e$  be such that  $c(e) - y(\mathcal{L}(e)) = 0$  and  $e \in \delta C$  for some  $C \in \mathcal{A}$ .
19:      Let  $C_1$  and  $C_2$  be the endpoints of  $e$  in  $\mathcal{L}^*$ .
20:       $\mathcal{L} \leftarrow \mathcal{L} \cup \{C_1 \cup C_2\}$                                ▷ Merge the two clusters.
21:       $y(C_1 \cup C_2) \leftarrow 0$ 
22:       $E_F \leftarrow E_F \cup \{e\}$                                        ▷ Add  $e$  to the forest.
23:    end if
24:     $\mathcal{A} \leftarrow \mathcal{L}^* \setminus \mathcal{D}$                                        ▷ Update active clusters.
25:  end while
26:  Restrict  $F$  to the  $g$  trees contained in  $\mathcal{A}$ .                                ▷ Discard trees spanning inactive clusters.
27:  while there is a  $D \in \mathcal{D}$  such that  $|\delta_F D| = 1$  do
28:     $V_F \leftarrow V_F \setminus D$                                        ▷ Pruning phase.
29:    Remove all edges from  $E_F$  with at least one endpoint in  $D$ .
30:  end while
31:  return  $F$ 
32: end function

```

D.1.2. ANALYSIS

We now show that the forest returned by PCSF-GW has the desired properties: it is a g -forest and satisfies the guarantee in Equation (11). Our analysis follows the overall approach of (Feofiloff et al., 2010).

Lemma 21. *Let $H = (V_H, E_H)$ be a graph and let $A, B \subseteq V_H$ be a partition of V_H . Moreover, let $F = \{T_1, \dots, T_g\}$ be a g -forest such that each T_i has no leaves in B and is not contained in B . Then*

$$\sum_{v \in A} \deg_F(v) + 2|A \setminus V_F| \leq 2|A| - 2g .$$

Proof. Since each T_i has no leaf in B and is not contained in B , every $v \in V_F \cap B$ satisfies $\deg_F(v) \geq 2$. Therefore,

$$\sum_{v \in V_F \cap B} \deg_F(v) \geq 2|V_F \cap B| .$$

Note that $\sum_{v \in V_F} \deg_F(v) = 2(|V_F| - g)$ because F divides V_F into g connected components. Hence

$$\begin{aligned} \sum_{v \in V_F \cap A} \deg_F(v) &= \sum_{v \in V_F} \deg_F(v) - \sum_{v \in V_F \cap B} \deg_F(v) \\ &\leq 2(|V_F| - g) - 2|V_F \cap B| \\ &= 2|V_F \cap A| - 2g . \end{aligned}$$

Moreover, $|A| = |A \cup V_F| + |A \setminus V_F|$. Combining this with the inequality above gives

$$\begin{aligned} \sum_{v \in V_F \cap A} \deg_F(v) + 2|A \setminus V_F| &\leq 2|V_F \cap A| - 2g + 2|A \setminus V_F| \\ &\leq 2|A| - 2g . \end{aligned}$$

Since $\sum_{v \in A} \deg_F(v) = \sum_{v \in V_F \cap A} \deg_F(v)$, the statement of the lemma follows. \square

Lemma 22. *Let \mathcal{L} be a laminar family, let $\mathcal{D} \subseteq \mathcal{L}$ be a sub-family, and let $\mathcal{A} = \mathcal{L}^* \setminus \mathcal{D}$. Let F be a g -forest which is \mathcal{L} -connected, \mathcal{L}^* -disjoint, has no leaf component in \mathcal{D} , and is not contained in \mathcal{D} . Then*

$$\sum_{C \in \mathcal{A}} |\delta_F C| + 2 \left| \{C \in \mathcal{A} \mid C \in \mathcal{L}_{\downarrow \overline{F}}\} \right| \leq 2|\mathcal{A}| - 2g .$$

Proof. Contract each set $C \in \mathcal{L}^*$ into a single node, keeping only edges with endpoints in distinct sets in \mathcal{L}^* . Call the resulting graph H and let A and B be the sets of vertices corresponding to \mathcal{A} and $\mathcal{L}^* \cap \mathcal{D}$, respectively.

Note that F is still a g -forest in H . Since F is \mathcal{L}^* -disjoint, no trees in T are connected by the contraction process. Moreover, no cycles are created because F is \mathcal{L} -connected. Let F' be the resulting g -forest in H . Since F has no leaf component in \mathcal{D} , F' has no leaves in B . Furthermore, no tree in F is contained in \mathcal{D} and thus no tree in F' is contained in B . Therefore, F' satisfies the conditions of Lemma 21.

Since F is \mathcal{L} -connected, there is a bijection between edges in F' and edges in F with endpoints in distinct elements of \mathcal{L}^* . Thus we have

$$\sum_{C \in \mathcal{A}} |\delta_F C| = \sum_{v \in A} \deg_{F'}(v) .$$

Furthermore, the contraction process gives

$$\left| \{C \in \mathcal{A} \mid C \in \mathcal{L}_{\downarrow \overline{F}}\} \right| = |A \setminus V_{F'}|$$

and $|\mathcal{A}| = |A|$. Now the statement of the lemma follows directly from applying Lemma 21. \square

Lemma 23. *At the beginning of every iteration of the growth phase in PCSF-GW (lines 7 to 24), the following invariant (I) holds:*

Let F be a g -forest which is \mathcal{L} -connected, \mathcal{L}^ -disjoint, has no leaf component in \mathcal{D} , and is not contained in \mathcal{D} . Moreover, let $A = \{v_1, \dots, v_g\}$ be an arbitrary set of g nodes in G and let $\mathcal{B} = \bigcup_{v \in A} \mathcal{L}_{\uparrow\{v\}}$. Then*

$$\sum_{e \in E_F} y(\mathcal{L}(e)) + 2y(\mathcal{L}_{\downarrow\overline{F}}) \leq y(\mathcal{L} \setminus \mathcal{B}). \quad (20)$$

Proof. Clearly, (I) holds at the beginning of the first iteration because the dual values y are 0 for every element in \mathcal{L} . We now assume that (I) holds at the beginning of an arbitrary iteration and show that (I) then also holds at the beginning of the next iteration. By induction, this then completes the proof.

Let \mathcal{L}' , \mathcal{D}' , \mathcal{A}' , and y' be the values of \mathcal{L} , \mathcal{D} , \mathcal{A} , and y at the beginning of the iteration. We analyze two separate cases based on the current event in this iteration of the loop: either a cluster is deactivated (lines 15 to 16) or two clusters are merged (lines 18 to 22).

First, we consider the cluster deactivation case. Let F be a g -forest satisfying the conditions of invariant (I). Since $\mathcal{L}' = \mathcal{L}$ and $\mathcal{D}' \subseteq \mathcal{D}$, F is also \mathcal{L}' -connected, \mathcal{L}'^* -disjoint, has no leaf component in \mathcal{D}' , and is not contained in \mathcal{D}' . Hence we can invoke Equation (20):

$$\sum_{e \in E_F} y'(\mathcal{L}'(e)) + 2y'(\mathcal{L}'_{\downarrow\overline{F}}) \leq y'(\mathcal{L}' \setminus \mathcal{B}). \quad (21)$$

Note that y and y' differ only on sets in $\mathcal{A}' = \mathcal{L}'^* \setminus \mathcal{D}'$. Therefore, we have the following three equations quantifying the differences between the three terms in Equations (20) and (21):

$$\bullet \sum_{e \in E_F} y(\mathcal{L}'(e)) - \sum_{e \in E_F} y'(\mathcal{L}'(e)) = \sum_{e \in E_F} \sum_{C \in \mathcal{A}'} \varepsilon \cdot \mathbf{1}[e \in \delta_F C] = \varepsilon \sum_{C \in \mathcal{A}'} |\delta_F C| \quad (22)$$

$$\bullet y(\mathcal{L}'_{\downarrow\overline{F}}) - y'(\mathcal{L}'_{\downarrow\overline{F}}) = \sum_{C \in \mathcal{A}'} \varepsilon \cdot \mathbf{1}[C \in \mathcal{L}'_{\downarrow\overline{F}}] = \varepsilon \left| \{C \in \mathcal{A}' \mid C \in \mathcal{L}'_{\downarrow\overline{F}}\} \right| \quad (23)$$

$$\bullet y(\mathcal{L}' \setminus \mathcal{B}) - y'(\mathcal{L}' \setminus \mathcal{B}) = \sum_{C \in \mathcal{A}'} \varepsilon \cdot \mathbf{1}[C \notin \mathcal{B}] = \varepsilon |\mathcal{A}'| - \varepsilon |\mathcal{A}' \cap \mathcal{B}| \geq \varepsilon |\mathcal{A}'| - \varepsilon g \quad (24)$$

In the last inequality, we used the fact that $|A| = g$ and hence \mathcal{B} can contain at most g maximal sets in the laminar family \mathcal{L}' . Combining the three equations above with Equation (21) and Lemma 22 then gives:

$$\sum_{e \in E_F} y(\mathcal{L}'(e)) + 2y(\mathcal{L}'_{\downarrow\overline{F}}) \leq y(\mathcal{L}' \setminus \mathcal{B}). \quad (25)$$

Since $\mathcal{L}' = \mathcal{L}$, this is equivalent to Equation (20), completing the proof for this case.

Now we consider the cluster merge case. As before, let F be a g -forest satisfying the conditions of invariant (I). Since $\mathcal{L} = \mathcal{L}' \cup \{C_1 \cup C_2\}$ and $\mathcal{D} = \mathcal{D}'$, F is also \mathcal{L}' -connected, \mathcal{L}'^* -disjoint, has no leaf component in \mathcal{D}' , and is not contained in \mathcal{D}' . Therefore, we can invoke Equation (20) again. Moreover, Equations (22), (23), and (24) also hold in this case. Combining these equations with (21) and Lemma 22 then again results in Equation (25). Furthermore, we have $y(C_1 \cup C_2) = 0$ and thus $y(\mathcal{L}(e)) = y(\mathcal{L}'(e))$, $y(\mathcal{L}_{\downarrow\overline{F}}) = y(\mathcal{L}'_{\downarrow\overline{F}})$, and $y(\mathcal{L} \setminus \mathcal{B}) = y(\mathcal{L}' \setminus \mathcal{B})$. Applying these equalities to Equation (25) completes the proof. \square

The following lemma is essential for proving a lower bound on the value of the optimal solution.

Lemma 24. *Let \mathcal{L} be a laminar family with dual values y . Let F be a g -forest and let $\mathcal{B} = \bigcup_{T \in F} \mathcal{L}_{\uparrow T}$. Then*

$$c(E_F) + \pi(\overline{V}_F) \geq y(\mathcal{L} \setminus \mathcal{B}).$$

Proof. Let $\mathcal{M} = \{C \in \mathcal{L} \mid \delta_F C \neq \{\}\}$ and $\mathcal{N} = \mathcal{L}_{\downarrow \overline{F}}$. Then $\mathcal{L} = \mathcal{M} \cup \mathcal{N} \cup \mathcal{B}$.

Since the y are dual values, we have $c(e) \geq y(\mathcal{L}(e))$ for every $e \in E_F$. Therefore,

$$\begin{aligned} c(E_F) &= \sum_{e \in E_F} c(e) \geq \sum_{e \in E_F} y(\mathcal{L}(e)) = \sum_{e \in E_F} \sum_{C \in \mathcal{L}(e)} y(C) \\ &= \sum_{C \in \mathcal{L}} \sum_{e \in \delta_F C} y(C) \geq \sum_{C \in \mathcal{M}} y(C) = y(\mathcal{M}). \end{aligned}$$

Moreover, we have $\pi(C) \geq y(\mathcal{L}_{\downarrow C})$ for every $C \in \mathcal{L}$. Thus,

$$\pi(\overline{V_F}) \geq \sum_{\substack{C \in \mathcal{L}^* \\ C \subseteq \overline{V_F}}} \pi(C) \geq \sum_{\substack{C \in \mathcal{L}^* \\ C \subseteq \overline{V_F}}} y(\mathcal{L}_{\downarrow C}) = y(\mathcal{L}_{\downarrow \overline{V_F}}) = y(\mathcal{N}).$$

Finally, we get

$$c(E_F) + \pi(\overline{V_F}) \geq y(\mathcal{M}) + y(\mathcal{N}) \geq y(\mathcal{M} \cup \mathcal{N}) = y(\mathcal{L} \setminus (\mathcal{L} \setminus (\mathcal{M} \cup \mathcal{N}))) \geq y(\mathcal{L} \setminus \mathcal{B}),$$

where we used $\mathcal{L} = \mathcal{M} \cup \mathcal{N} \cup \mathcal{B}$ in the final step. \square

We can now prove the main theorem establishing an approximation guarantee for PCSF-GW, which also proves Theorem 5 from the main text of the paper.

Theorem 25. *Let F be the result of PCSF-GW(G, c, π, g). Then F is a g -forest and*

$$c(F) + 2\pi(\overline{F}) \leq 2c(F_{OPT}) + 2\pi(\overline{F_{OPT}}),$$

where F_{OPT} is a g -forest minimizing $c(F_{OPT}) + \pi(\overline{F_{OPT}})$.

Proof. By construction in the growth phase of PCSF-GW (lines 7 to 24), F is a \mathcal{L} -connected forest at the end of the growth phase. Since at most one element is added to \mathcal{D} in each iteration of the growth phase, we have $|\mathcal{A}| = g$ at the end of the growth phase. Hence restricting F to \mathcal{A} in line 26 leads to F being a g -forest which is still \mathcal{L} -connected. Furthermore, F is \mathcal{L}^* -disjoint and no tree in F is contained in \mathcal{D} .

The pruning phase (lines 27 to 29) maintains that F is a g -forest, \mathcal{L} -connected, \mathcal{L}^* -disjoint, and not contained in \mathcal{D} . Moreover, the pruning phase removes all leaf components of F in \mathcal{D} . Hence at the end of the pruning phase, F satisfies the conditions of Lemma 23 (\mathcal{L} , \mathcal{D} , and \mathcal{A} did not change in the pruning phase).

Now let $F_{OPT} = (T_1^{OPT}, \dots, T_g^{OPT})$ be a g -forest minimizing $c(F_{OPT}) + \pi(\overline{F_{OPT}})$ and let $A = \{v_1, \dots, v_g\}$ with $v_i \in T_i^{OPT}$. Moreover, let $\mathcal{B}_1 = \bigcup_{v \in A} \mathcal{L}_{\uparrow \{v\}}$ as in Lemma 23. Invoking the Lemma then gives

$$\sum_{e \in E_F} y(\mathcal{L}(e)) + 2y(\mathcal{L}_{\downarrow \overline{F}}) \leq 2y(\mathcal{L} \setminus \mathcal{B}_1). \quad (26)$$

Now, note that every $e \in E_F$ was added to F when we had $c(e) = y(\mathcal{L}(e))$. Hence

$$\sum_{e \in E_F} y(\mathcal{L}(e)) = \sum_{e \in E_F} c(e) = c(F). \quad (27)$$

Moreover, $\overline{V_F}$ can be decomposed into elements in \mathcal{D} : after restricting F to $\mathcal{A} = \mathcal{L}^* \setminus \mathcal{D}$ in line 26 this clearly holds. During the pruning phase, all subtrees that are removed from trees in F are elements of \mathcal{D} . Therefore, there is a family of pairwise disjoint sets $\mathcal{Z} \subseteq \mathcal{D}$ such that $\bigcup_{C \in \mathcal{Z}} C = \overline{V_F}$. Note that for every $C \in \mathcal{D}$ we have $\pi(C) = y(\mathcal{L}_{\downarrow C})$ because C was deactivated at some point in the growth phase. Therefore,

$$\pi(\overline{F}) = \sum_{C \in \mathcal{Z}} \pi(C) = \sum_{C \in \mathcal{Z}} y(\mathcal{L}_{\downarrow C}) \leq y(\mathcal{L}_{\downarrow \overline{F}}). \quad (28)$$

Combining Equations (26), (27), and (28) then gives

$$c(F) + 2\pi(\overline{F}) \leq 2y(\mathcal{L} \setminus \mathcal{B}_1). \quad (29)$$

We now relate this upper bound to the optimal solution F_{OPT} . Let $\mathcal{B}_2 = \bigcup_{T \in F_{OPT}} \mathcal{L}_{\uparrow T}$ as in Lemma 24. The y are valid dual values due to their construction in PCSF-GW. Thus Lemma 24 gives

$$y(\mathcal{L} \setminus \mathcal{B}_2) \leq c(F_{OPT}) + \pi(\overline{F_{OPT}}). \quad (30)$$

Note that $\mathcal{B}_2 \subseteq \mathcal{B}_1$ and therefore $y(\mathcal{L} \setminus \mathcal{B}_1) \leq y(\mathcal{L} \setminus \mathcal{B}_2)$. The guarantee in the theorem now follows directly from Equations (29) and (30). \square

D.2. A fast algorithm for Goemans-Williamson

We now introduce our fast variant of the GW scheme. To the best of our knowledge, our algorithm is the first practical implementation of a GW-like algorithm that runs in nearly linear time.

On a high level, our algorithm uses a more aggressive and *adaptive* dynamic edge splitting scheme than (Cole et al., 2001): our algorithm moves previously inserted splitting points in order to reach a tight edge constraint quicker than before. By analyzing the precision needed to represent merge and deactivation events in the GW algorithm, we prove that our algorithm runs in $O(\alpha \cdot |E| \log |V|)$ time, where α is the number of bits used to specify each value in the input. For constant bit precision α (as is often the case in practical applications) our algorithm hence has a running time of $O(|E| \log |V|)$. Furthermore, our algorithm achieves the approximation guarantee (11) exactly without the additional $\frac{2}{n^k}$ term present in the work of (Cole et al., 2001). From an empirical point of view, our more aggressive edge splitting scheme produces only very few additional edge pieces: we observed that the number of processed edge events is usually close to $2|E|$, the number of edge events initially created. We demonstrate this empirical benefit in our experiments (see Section D.3).

Since the pruning stage of the GW scheme can be implemented relatively easily in linear time (Johnson et al., 2000), we focus on the moat growing stage here. We also remark that there are algorithms for the PCST problem that achieve a nearly-linear time for *planar* graphs (Bateni et al., 2011; Eisenstat et al., 2012).

D.2.1. ALGORITHM

Similar to (Cole et al., 2001), our algorithm divides each edge $e = (u, v)$ into two *edge parts* e_u and e_v corresponding to the endpoints u and v . We say an edge part p is *active* if its endpoint is in an active cluster, otherwise the edge part p is *inactive*. The key advantage of this approach over considering entire edges is that *all active edge parts always grow at the same rate*. For each edge part p , we also maintain an *event value* $\mu(p)$. This event value is the total amount that the moats on edge part p are allowed to grow until the next event for this edge occurs. In order to ensure that the moats growing on the two corresponding edge parts e_u and e_v never overlap, we always set the event values so that $\mu(e_u) + \mu(e_v) = c(e)$. As for edges, we define the remaining slack of edge part e_u as $\mu(e_u) - \sum_{C \in \mathcal{C}} y_C$, where \mathcal{C} is the set of clusters containing node u .

We say that an *edge event* occurs when an edge part has zero slack remaining. However, this does not necessarily mean that the corresponding edge constraint has become tight as the edge event might be “stale” since the other edge parts has become inactive and stopped growing since the last time the edge event was updated. Nevertheless, we will be able to show that the total number of edge events to be processed over the course of the algorithm is small. Note that we can find the next edge event by looking at the edge events with smallest remaining slack values in their clusters. This is an important property because it allows us to organize the edge parts in an efficient manner. In particular, we maintain a priority queue Q_C for each cluster C that contains the edge parts with endpoint in C , sorted by the time at which the next event on each edge part occurs. Furthermore, we arrange the cluster priority queues in an overall priority queue resulting in a “heap of heaps” data structure. This data structure allows us to quickly locate the next edge event and perform the necessary updates after cluster deactivation or merge events.

In addition to the edge events, we also maintain a priority queue of *cluster events*. This priority queue contains each active cluster with the time at which the corresponding cluster constraint becomes tight. Using these definitions, we can now state the high-level structure of our algorithm in pseudo code (see Algorithm 6) and then describe the two subroutines MERGECLUSTERS and GENERATENEWEDGEEVENTS in more detail.

Algorithm 6 Fast variant of the GW algorithm for PCSF.

```

1: function PCSF-FAST( $G, c, \pi, g$ )
2:   INITPCST( $G, c, \pi$ )
3:    $t \leftarrow 0$   $\triangleright$  Current time
4:    $g' \leftarrow |V|$   $\triangleright$  Number of active clusters
5:   while  $g' > g$  do
6:      $\triangleright$  Returns event time and corresponding edge part
7:      $(t_e, p_u) \leftarrow \text{GETNEXTEDGEEVENT}()$ 
8:      $\triangleright$  Returns event time and corresponding cluster
9:      $(t_c, C) \leftarrow \text{GETNEXTCLUSTEREVENT}()$ 
10:    if  $t_e < t_c$  then
11:       $t \leftarrow t_e$ 
12:      REMOVENEXTEDGEEVENT()
13:       $p_v \leftarrow \text{GETOTHEREDGEPART}(p_u)$ 
14:       $\triangleright$  GETSUMONEDGEPART returns the current moat sum on the edge part
15:       $\triangleright p_u$  and the maximal cluster containing  $u$ 
16:       $(s, C_u) \leftarrow \text{GETSUMONEDGEPART}(p_u)$ 
17:       $(s', C_v) \leftarrow \text{GETSUMONEDGEPART}(p_v)$ 
18:       $r \leftarrow \text{GETEDGE}(\text{COST}(p_u) - s - s')$   $\triangleright$  Remaining amount on the edge
19:      if  $C_u = C_v$  then  $\triangleright$  The two endpoints are already in the same cluster
20:        continue  $\triangleright$  Skip to beginning of while-loop
21:      end if
22:      if  $r = 0$  then
23:        MERGECLUSTERS( $C_u, C_v$ )
24:      else
25:        GENERATENEWEDGEEVENTS( $p_u, p_v$ )
26:      end if
27:    else
28:       $t \leftarrow t_c$ 
29:      REMOVENEXTCLUSTEREVENT()
30:      DEACTIVATECLUSTER( $C$ )
31:       $g' \leftarrow g' - 1$ 
32:    end if
33:  end while
34:  PRUNING()
35: end function

```

MERGECLUSTERS (C_u, C_v) : As a first step, we mark C_u and C_v as inactive and remove them from the priority queue keeping track of cluster deactivation events. Furthermore, we remove the priority queues Q_{C_u} and Q_{C_v} from the heap of heaps for edge events. Before we merge the heaps of C_u and C_v , we have to ensure that both heaps contain edge events on the “global” time frame. If C_u (or C_v) is inactive since time t' when the merge occurs, the edge event times in Q_{C_u} will have become “stale” because the moat on edge parts incident to C_u did not grow since t' . In order to correct for this offset and bring the keys in Q_{C_u} back to the global time frame, we first increase all keys in Q_{C_u} by $t - t'$. Then, we merge Q_{C_u} and Q_{C_v} , which results in the heap for the new merged cluster. Finally, we insert the new heap into the heap of heaps and add a new entry to the cluster deactivation heap.

GENERATENEWEDGEEVENTS (p_u, p_v) : This function is invoked when an edge event occurs, but the corresponding edge constraint is not yet tight. Since the edge part p_u has no slack left, this means that there is slack remaining on p_v . Let \mathcal{C}_u and \mathcal{C}_v be the set of clusters containing u and v , respectively. Then $r = c(e) - \sum_{C \in \mathcal{C}_u \cup \mathcal{C}_v} y_C$ is the length of the part of edge e not covered by moats yet. We distinguish two cases:

1. The cluster containing the endpoint v is active.

Since both endpoints are active, we expect both edge parts to grow at the same rate until they meet and the edge constraint becomes tight. Therefore, we set the new event values to $\mu(p_u) = \sum_{C \in \mathcal{C}_u} + \frac{r}{2}$ and $\mu(p_v) = \sum_{C \in \mathcal{C}_v} + \frac{r}{2}$. Note that this maintains the invariant $\mu(p_u) + \mu(p_v) = c(e)$. Using the new event values for p_u and p_v , we update the priority queues Q_{C_u} and Q_{C_v} accordingly and then also update the heap of heaps.

2. The cluster containing the endpoint v is inactive.

In this case, we assume that v stays inactive until the moat growing on edge part p_u makes the edge constraint for e tight. Hence, we set the new event values to $\mu(p_u) = \sum_{C \in \mathcal{C}_u} + r$ and $\mu(p_v) = \sum_{C \in \mathcal{C}_v}$. As in the previous case, this maintains the invariant $\mu(p_u) + \mu(p_v) = c(e)$ and we update the relevant heaps accordingly. It is worth noting our setting of $\mu(p_v)$ reduces the slack for p_v to zero. This ensures that as soon as the cluster C_v becomes active again, the edge event for p_v will be processed next.

Crucially, in **GENERATENEWEDGEEVENTS**, we set the new event values for p_u and p_v so that the next edge event on e would merge the clusters C_u and C_v , *assuming both clusters maintain their current activity status*. If one of the two clusters changes its activity status, this will not hold:

1. If both clusters were active and cluster C_u has become inactive since then, the next event on edge e will be part p_v reaching the common midpoint. However, due to the deactivation of C_u , the edge part p_u will not have reached the common midpoint yet.
2. If C_v was inactive and becomes active before the edge event for p_u occurs, the edge event for p_v will also immediately occur after the activation for C_v . At this time, the moat on p_u has not reached the new, size-0 moat of C_v , and thus the edge constraint is not tight.

However, in the next section we show that if all input values are specified with d bits of precision then at most $O(d)$ edge events can occur per edge. Moreover, even in the general case our experiments in Section 6 show that the pathological cases described above occur very rarely in practice. In most instances, only two edge events are processed per edge on average.

D.2.2. ANALYSIS

We now study the theoretical properties of our algorithm PCSF-FAST. Note that by construction, the result of our algorithm exactly matches the output of PCSF-GW and hence also satisfies guarantee (11).

First, we establish the following structural result for the growth stage of the GW algorithm (the “exact” algorithm PCSF-GW, not yet PCSF-FAST). Informally, we show that a single additional bit of precision suffices to exactly represent all important events in the moat growth process. The following result is equivalent to Theorem 6.

Theorem 26. *Let all node prizes $\pi(v)$ and edge costs $c(e)$ be even integers. Then all cluster merge and deactivation events occur at integer times.*

Proof. We prove the theorem by induction over the cluster merge and deactivation events occurring in the GW scheme, sorted by the time at which the events happen. We will show that the updates caused by every event maintain the following invariant:

Induction hypothesis Based on the current state of the algorithm, let t_e be the time at which the edge constraint for edge e becomes tight and t_C be the time at which the cluster constraint for cluster C becomes tight. Then t_e and t_C are integers. Moreover, if the merge event at t_e is a merge event between an active cluster and an inactive cluster C , then $t_e - t_{\text{inactive}(C)}$ is even, where $t_{\text{inactive}(C)}$ is the time at which cluster C became inactive.

Clearly, the induction hypothesis holds at the beginning of the algorithm: all edge costs are even, so $t_e = \frac{c(e)}{2}$ is an integer. Since the node prizes are integers, so are the t_C . The assumption on merge events with inactive clusters trivially holds because there are no inactive clusters at the beginning of the algorithm. Next, we perform the induction step by a case analysis over the possible events:

- **Active-active:** a merge event between two active clusters. Since this event modifies no edge events, we only have to consider the new deactivation event for the new cluster C . By the induction hypothesis, all events so far have occurred at integer times, so all moats have integer size. Since the sum of prizes in C is also an integer, the new cluster constraint becomes tight at an integer time.
- **Active-inactive:** a merge event between an active cluster and an inactive cluster. Let e be the current edge, t_e be the current time, and C be the inactive cluster. The deactivation time for the new cluster is the same as that of the current active cluster, so it is also integer. Since every edge e' incident to C now has a new growing moat, we have to consider the change in the event time for e' . We denote the previous event time of e' with $t'_{e'}$. We distinguish two cases:
 - If the other endpoint of e' is in an active cluster, the part of e' remaining has size $t'_{e'} - t_e$ and e' becomes tight at time $t_e + \frac{t'_{e'} - t_e}{2}$ because e' has two growing moats. We have

$$t'_{e'} - t_e = (t'_{e'} - t_{\text{inactive}(C)}) - (t_e - t_{\text{inactive}(C)}).$$

Note that both terms on the right hand side are even by the induction hypothesis, and therefore their difference is also even. Hence the new event time for edge e' is an integer.

- If the other endpoint of e' is an inactive cluster, say C' , we have to show that $t_{e'} - t_{\text{inactive}(C')}$ is even, where $t_{e'}$ is the new edge event time for e' . We consider whether C or C' became inactive last:
 - * C became inactive last: from the time at which C became inactive we know that $t'_{e'} - t_{\text{inactive}(C')}$ is even. Moreover, we have that $t_{e'} = t'_{e'} + (t_e - t_{\text{inactive}(C)})$. Since $t_e - t_{\text{inactive}(C)}$ is even by the induction hypothesis, so is $t_{e'} - t_{\text{inactive}(C')}$.
 - * C' became inactive last: from the time at which C' became inactive we know that $t'_{e'} - t_{\text{inactive}(C)}$ is even. The time of the new edge event can be written as $t_{e'} = t_e + t'_{e'} - t_{\text{inactive}(C')}$ (an integer by the induction hypothesis), which is equivalent to $t_{e'} - t'_{e'} = t_e - t_{\text{inactive}(C')}$. We now use this equality in the second line of the following derivation:

$$\begin{aligned} t_{e'} - t_{\text{inactive}(C')} &= t_{e'} - t'_{e'} + t'_{e'} - t_e + t_e - t_{\text{inactive}(C')} \\ &= 2(t_{e'} - t'_{e'}) + t'_{e'} - t_e \\ &= 2(t_{e'} - t'_{e'}) + (t'_{e'} - t_{\text{inactive}(C)}) - (t_e - t_{\text{inactive}(C)}). \end{aligned}$$

Since $t_e - t_{\text{inactive}(C)}$ is even by the induction hypothesis, all three terms on the right hand side are even.

- **Cluster deactivation:** Clearly, a deactivation of cluster C leads to no changes in other cluster deactivation times. Moreover, edges incident to C and another inactive cluster will never become tight based on the current state of the algorithm. The only quantities remaining are the edge event times for edges e with another cluster endpoint that is active. Note that up to time t_C , the edge e had two growing moats and t_e was an integer. Therefore, the part of e remaining has length $2(t_e - t_C)$, which is an even integer. The new value of t_e is $t_C + 2(t_e - t_C)$, and since $t_{\text{inactive}(C)} = t_C$ the induction hypothesis is restored.

Since the induction hypothesis is maintained throughout the algorithm and implies the statement of the theorem, the proof is complete. \square

We now use this result to show that the number of edge part events occurring in PCSF-FAST is small.

Corollary 27. *Let all node prizes $\pi(v)$ and edge costs $c(e)$ be specified with α bits of precision. Then the number of edge part events processed in PCSF-FAST is bounded by $O(\alpha \cdot |E|)$.*

Proof. We look at each edge e individually. For every edge part event A on e that does not merge two clusters, the following holds: either A reduces the remaining slack of e by at least a factor of two or the event directly preceding A reduced the remaining slack on e by at least a factor of two. In the second case, we charge A to the predecessor event of A .

So after $O(\alpha)$ edge parts events on e , the remaining slack on e is at most $\frac{c(e)}{2^\alpha}$. Theorem 26 implies that the minimum time between two cluster merge or deactivation events is $\frac{c(e)}{2^{\alpha+1}}$. So after a constant number of additional edge part events on e , the edge constraint of e must be the next constraint to become tight, which is the last edge part event on e to be processed. Therefore, the total number of edge part events on e is $O(\alpha)$. \square

We now show that all subroutines in PCSF-FAST can be implemented in $O(\log|V|)$ amortized time, which leads to our final bound on the running time.

Theorem 28. *Let all node prizes $\pi(v)$ and edge costs $c(e)$ be specified with α bits of precision. Then PCSF-FAST runs in $O(\alpha \cdot |E| \log|V|)$ time.*

Proof. The requirements for the priority queue maintaining edge parts events are the standard operations of a mergeable heap data structure, combined with an operation that adds a constant offset to all elements in a heap in $O(\log|V|)$ amortized time. We can build such a data structure by augmenting a pairing heap with an offset value at each node. Due to space constraints, we omit the details of this construction here. For the outer heap in the heap of heaps and the priority queue containing cluster deactivation events, a standard binomial heap suffices.

We represent the laminar family of clusters in a tree structure: each cluster C is a node, the child nodes are the two clusters that were merged to form C , and the parent is the cluster C was merged into. The initial clusters, i.e., the individual nodes, form the leaves of the tree. By also storing the moat values at each node, the GETSUMONEDGEPART operation for edge part p_u can be implemented by traversing the tree from leaf u to the root of its subtree. However, the depth of this tree can be up to $\Omega(|V|)$. In order to speed up the data structure, we use path compression in essentially the same way as standard union-find data structures. The resulting amortized running time for GETSUMONEDGEPART and merging clusters then becomes $O(\log|V|)$ via a standard analysis of union-find data structures (with path compression only).

This shows that all subroutines in PCSF-FAST (Algorithm 6) can be implemented to run in $O(\log|V|)$ amortized time. Since there are at most $O(\alpha|E|)$ events to be processed in total, the overall running time bound of $O(\alpha \cdot |E| \log|V|)$ follows. \square

D.3. Experimental results

We also investigate the performance of our algorithm PCSF-FAST outside sparse recovery. As test data, we use the public instances of the DIMACS challenge on Steiner tree problems¹⁰. We record both the total running times and the number of edge events processed by our algorithm. All experiments were conducted on a laptop computer from 2010 (Intel Core i7 with 2.66 GHz, 4 MB of cache, and 8 GB of RAM). All reported running times are averaged over 11 trials after removing the slowest run. Since the GW scheme has a provable approximation guarantee, we focus on the running time results here.

Running times Figure 6 shows the running times of our algorithm on the public DIMACS instances for the unrooted prize-collecting Steiner tree problem (PCSPG). For a single instance, the maximum running time of our algorithm is roughly 1.3 seconds and most instances are solved significantly faster. The scatter plots also demonstrates the nearly-linear scaling of our running time with respect to the input size.

Effectiveness of our edge splitting heuristic As pointed out in our running time analysis, the number of edge part events determines the overall running time of our algorithm. For input values specified with α bits of precision, our analysis shows that the algorithm encounters at most $O(\alpha)$ events per edge. In order to get a better understanding of our

¹⁰<http://dimacs11.cs.princeton.edu/>

empirical performance, we now look at the number of edge part events encountered by our algorithm (see Figure 7). The scatter plots show that the average number of events per edge is less than 3 for all instances. These results demonstrate the effectiveness of our more adaptive edge splitting heuristics. Moreover, the number of edge events encountered explains the small running times on the large i640 instances in Figure 6.

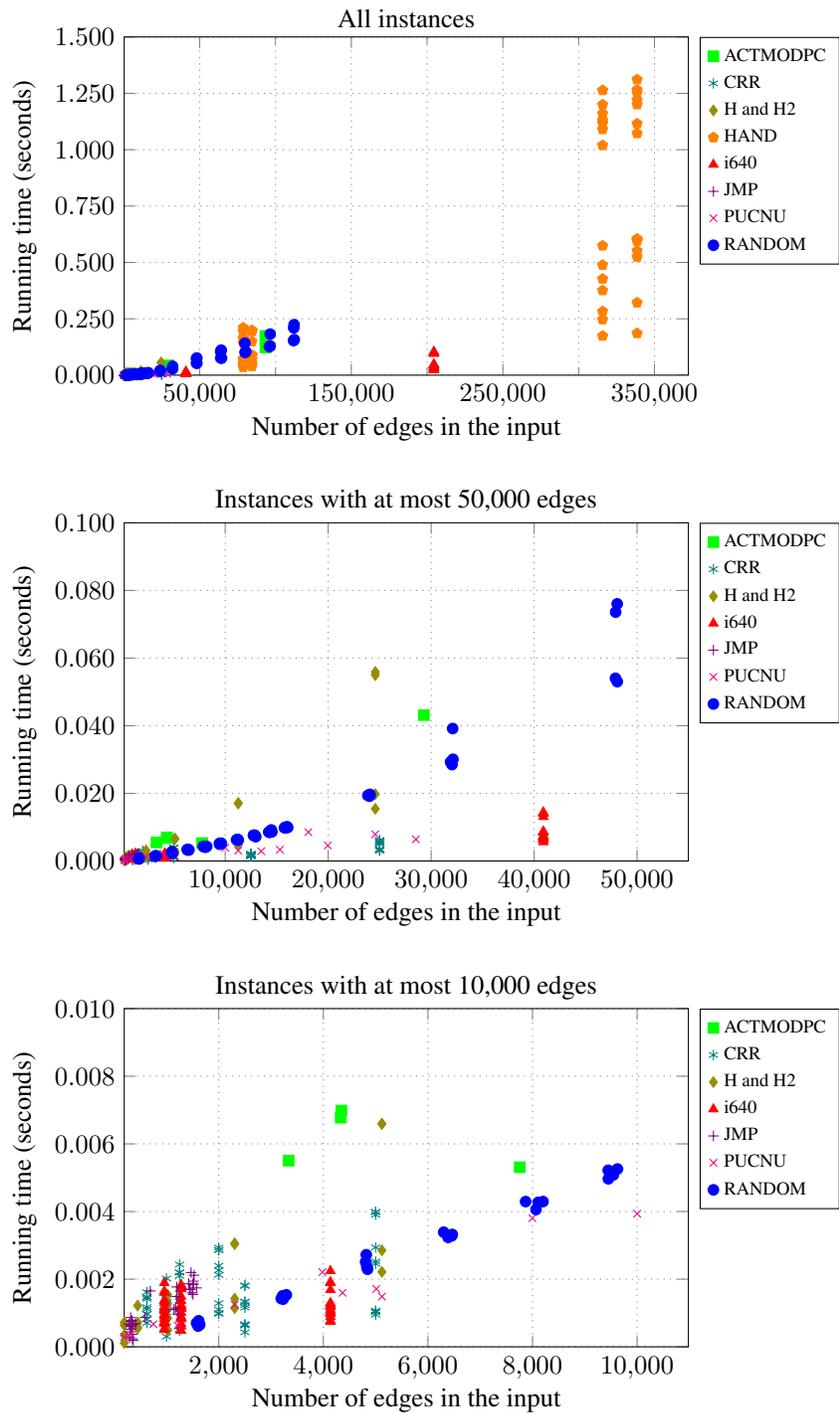


Figure 6. Running times for the PCSPG instances of the DIMACS challenge. Each color corresponds to one test case group. Our algorithm runs for at most 1.3s on any instance and clearly shows nearly-linear scaling with the input size.

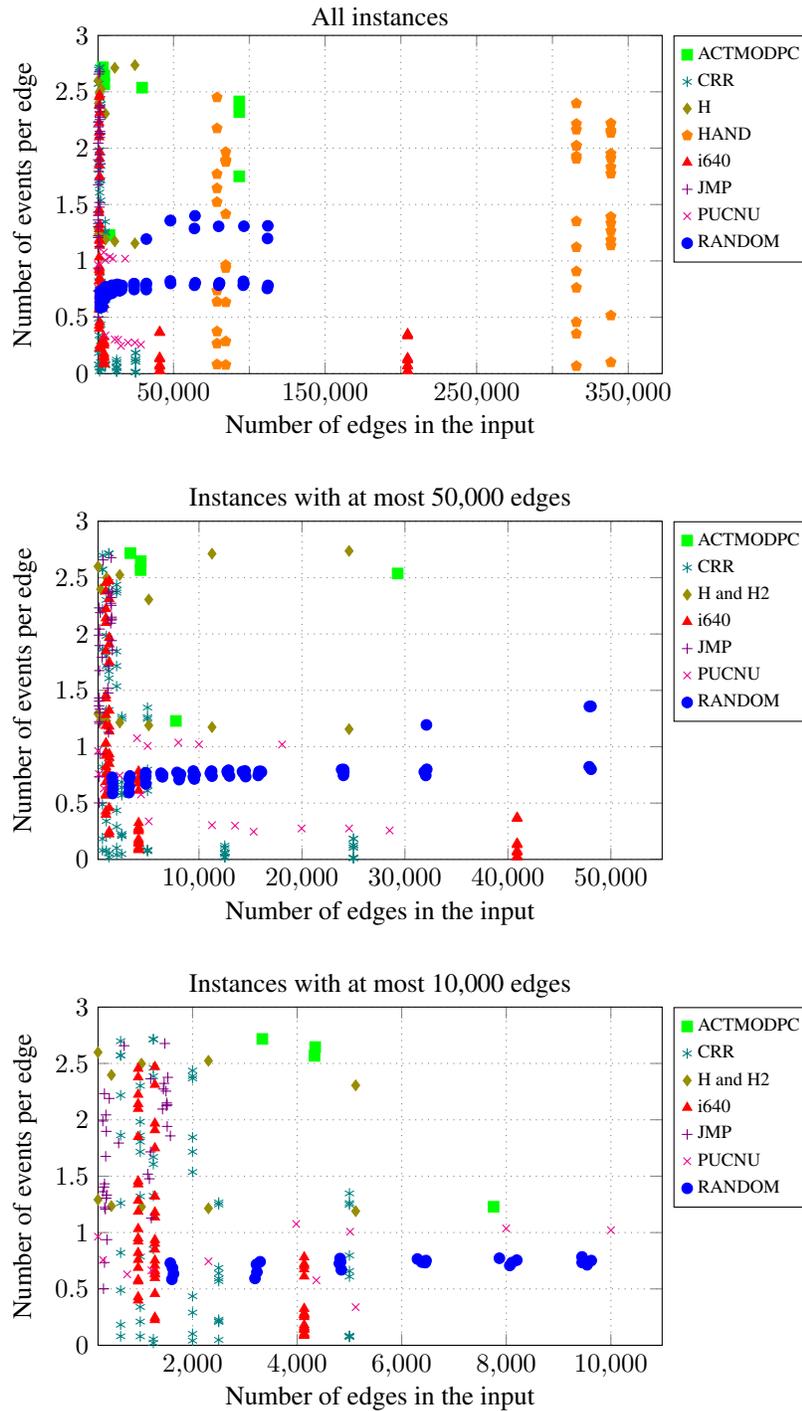


Figure 7. Average number of edge events processed per edge for the PCSPG instances of the DIMACS challenge. Each color corresponds to one test case group. The results demonstrate the effectiveness of our edge splitting approach and show that the average number of edge events is less than 3 for every instance.