# Improved Bounds for Sampling Solutions of Random CNF Formulas

Kun He[*]        Kewen Wu[†]        Kuan Yang[‡]

### Abstract

Let $\Phi$ be a random $k$-CNF formula on $n$ variables and $m$ clauses, where each clause is a disjunction of $k$ literals chosen independently and uniformly. Our goal is to sample an approximately uniform solution of $\Phi$ (or equivalently, approximate the partition function of $\Phi$).

Let $\alpha = m/n$ be the density. The previous best algorithm runs in time $n^{\mathsf{poly}(k,\alpha)}$ for any $\alpha \lesssim 2^{k/300}$ [Galanis, Goldberg, Guo, and Yang, SIAM J. Comput.'21]. Our result significantly improves both bounds by providing an almost-linear time sampler for any $\alpha \lesssim 2^{k/3}$.

The density $\alpha$ captures the *average degree* in the random formula. In the worst-case model with bounded *maximum degree*, current best efficient sampler works up to degree bound $2^{k/5}$ [He, Wang, and Yin, FOCS'22 and SODA'23], which is, for the first time, superseded by its average-case counterpart due to our $2^{k/3}$ bound. Our result is the first progress towards establishing the intuition that the solvability of the average-case model (random $k$-CNF formula with bounded average degree) is better than the worst-case model (standard $k$-CNF formula with bounded maximal degree) in terms of sampling solutions.

## 1   Introduction

A random $k$-CNF formula $\Phi = \Phi(k, n, m)$ is a formula on $n$ Boolean variables and $m$ clauses, where each clause is a disjunction of $k$ literals sampled from all $2n$ possible literals uniformly and independently. Let $\alpha = m/n$ be the *density* of the formula, which captures the *average degree* for variables in $\Phi$.

The random $k$-CNF model exhibits a fascinating phenomenon of a sharp phase transition in satisfiability. Based on numerical simulations and non-rigorous arguments in physics [MPZ02, MMZ05], it was conjectured that there exists a critical value $\alpha_\star = \alpha_\star(k)$ such that for all $\varepsilon > 0$, it holds that

$$\lim_{n \to \infty} \mathbf{Pr}\left[\Phi(k, n, m) \text{ is satisfiable}\right] = \begin{cases} 1 & \text{if } \alpha = \alpha_\star - \varepsilon, \\ 0 & \text{if } \alpha = \alpha_\star + \varepsilon. \end{cases}$$

It has been a well-known challenge to prove the conjecture and determine the critical value $\alpha_\star$. Following a line of work, [KKKS98, Fri99, AM02, AP03, COP16], this conjecture is proved by Ding, Sly, and Sun [DSS22] for sufficiently large $k$, where the exact value of $\alpha_\star$ is also established. Roughly speaking, we have $\alpha_\star = 2^k \ln(2) - (1 + \ln(2))/2 + o_k(1)$ as $k \to +\infty$.

However, the method for showing the sharp lower bound of $\alpha_\star$ is not constructive, and thus does not provide efficient algorithms to find solutions. The current best polynomial-time algorithm for searching solutions is the FIX algorithm given by Coja-Oghlan [Coj10], which succeeds with high

probability if $\alpha \lesssim 2^k \ln(k)/k$.[1] This is conjectured to be the search threshold, i.e., finding a solution is conjectured computationally hard if $\alpha$ goes beyond $2^k \ln(k)/k$. It is known [ACO08] that the solution space of random formulas has long-range correlations beyond density bound $2^k \ln(k)/k$, which suggests that local search algorithms are unlikely to succeed in polynomial time. Later, some particular algorithms have been ruled out (See e.g., [Het16, COHH17]). To date, the strongest negative result is given by Bresler and Huang [BH22], which proves that a class of *low-degree polynomial algorithms* (including FIX) cannot efficiently solve random $k$-CNF formulas beyond density $4.911 \cdot 2^k \ln(k)/k$. This gives a strong evidence that $2^k \ln(k)/k$ is the correct algorithmic phase transition.

Beyond decision and search, it is a natural next step to sample a satisfying assignment uniformly from the solution space. This is closely related to approximating the number of solutions of the formula $\Phi$, denoted by $Z(\Phi)$, and falls under the algorithmic study of partition functions in statistical physics. Montanari and Shah [MS07] presented the first efficient algorithm to approximately compute the partition function $\log(Z_\beta(\Phi))/n$ for a weighted model of random $k$-CNF, where the weight of an assignment $\sigma$ is $\mathbf{e}^{-\beta \cdot H(\sigma)}$ and $H(\sigma)$ is the number of unsatisfied clauses under $\sigma$. The number of satisfying assignments $Z(\Phi)$ then corresponds to $\lim_{\beta \to +\infty} Z_\beta(\Phi)$. However, their algorithm is based on the correlation decay method and only works within the *uniqueness regime* of the Gibbs distribution of the random $k$-CNF model. This uniqueness regime is $\alpha \lesssim 2 \ln(k)/k$, exponentially lower than the satisfiability and search thresholds. The first significant improvement was given by Galanis, Goldberg, Guo, and Yang [GGGY21], who designed a fully polynomial-time approximation scheme for $Z(\Phi)$ with runtime $n^{\mathsf{poly}(k,\alpha)}$ assuming $\alpha \lesssim 2^{k/300}$.

**Comparison with the Worst-Case Model.** Since the density $\alpha$ is defined to be the ratio between the number of clauses and variables, it is easy to see that $k \cdot \alpha$ equals the *average degree* of variables in the random $k$-CNF model. Here we compare this average-case model (i.e., random $k$-CNF formulas with average degree $k\alpha$) with its worst-case counterpart (i.e., standard $k$-CNF formulas with maximum degree $d$). Since randomness kills structures in the worst-case examples, intuitively the average-case model should have advantages over the worst-case model in terms of solvability under the same (average/maximum) degree assumption. This brings out the following intriguing question:

*Is it true that the average-case model is easier to solve than the worst-case model?*

This question has been answered affirmatively for satisfiability and search:

- The satisfiability threshold of the average-case model is $k\alpha \approx k2^k \ln(2)$ [DSS22], whereas it shrinks to $d \approx 2^{k+1}/(\mathbf{e}k)$ in the worst-case model by the lopsided Lovász local lemma [ES91, GST16].

- The search threshold for the average-case model is (at least) $k\alpha \approx 2^k \ln(k)$ [Coj10], which is still beyond the above $d \approx 2^{k+1}/(\mathbf{e}k)$ satisfiability threshold of the worst-case model.[2]

Given these, it is reasonable to speculate that the task of sampling solutions is also easier in the average-case model than the worst-case model, which, however, is less clear before our work.

Moitra [Moi19] designed the the first sampling algorithm for the worst-case model, which works whenever $d \lesssim 2^{k/60}$ and runs in time $n^{\mathsf{poly}(k,d)}$. Since then, both the degree bound and the runtime have been significantly improved. After [FGYZ21, FHY21, JPV21, HSW21], the state-of-the-art

---

[1]We use $\lesssim$ to informally and flexibly hide low-order terms to simplify expressions.

[2]In fact, the search threshold for the worst-case model here is indeed this bound [GST16].

bound is $d \lesssim 2^{k/5}$ and $n \cdot \mathsf{poly}(k, d, \log(n))$ runtime by He, Wang, Yin [HWY22, HWY23]. In terms of the computational hardness, Bezáková, Galanis, Goldberg, Guo, and Štefankovič [BGG+19] showed that the sampling task becomes intractable if $d$ can go beyond $2^{k/2}$ assuming $\mathsf{NP} \neq \mathsf{RP}$.

In contrast, for the average-case model, there is no improvement after [GGGY21]. The best bound is still $\alpha \lesssim 2^{k/300}$ and $n^{\mathsf{poly}(k,\alpha)}$ runtime, which falls short of the solvability intuition. Indeed, [GGGY21] builds upon the techniques of [Moi19], and thus has the similar runtime bound; on the other hand, the existence of high-degree variables in the random setting poses significant challenges in carrying over the previous analysis, which results in the even worse $2^{k/300} \ll 2^{k/60}$ degree bound. Moreover, the ideas leading to subsequent improvements [FGYZ21, FHY21, JPV21, HSW21] over [Moi19] do not seem to extend here. We will elaborate in more detail in Subsection 1.2.

Therefore, it remains an intriguing open problem whether the "average-case easier than worst-case" conjecture is also true for sampling thresholds. Our result is the first evidence towards this direction: Our algorithms works up to $\alpha \lesssim 2^{k/3}$ and runs in time $n^{1+o_k(1)} \cdot \mathsf{poly}(k, \alpha, \log(n))$. This not only drastically improves both degree and runtime bounds in [GGGY21], but outperforms the current best $2^{k/5}$ degree bound in the worst-case model [HWY22, HWY23] as predicted by the intuition above.[3]

**Independent Works.** Independent of our work, there are two recent works on sampling solutions of random $k$-CNF formulas [GGGH22, CMM23] improving [GGGY21]. The algorithm from [GGGH22] works when $\alpha \lesssim 2^{0.039k}$ and runs in almost-linear time; and the algorithm from [CMM23] requires $\alpha \lesssim 2^{0.0134k}$ and runs in $n^{\mathsf{poly}(k,\alpha)}$ time. In terms of results, our density bound $\alpha \lesssim 2^{k/3}$ and almost-linear runtime subsume both of them.

Both [GGGH22] and [CMM23] use Markov-chain-based algorithms in line with [FGYZ21, FHY21, JPV21, HSW21], while our algorithm follows the recursive sampling approach recently developed in [AJ22, HWY22, HWY23]. Therefore both the analysis and bounds of the papers are very different.

## 1.1 Our Results and Future Directions

Our main result is a Monte Carlo algorithm with almost-linear runtime for sampling solutions of a random CNF formula with large density.

Theorem 1.1 characterizes the extreme case where $\alpha$ is close to $2^{k/3}$ up to $\mathsf{poly}(k)$ factors.

**Theorem 1.1** (Informal). *Assume $\alpha \approx 2^{k/3}$ and $k, n$ sufficiently large. Then with high probability, we can sample an approximate uniform solution of $\Phi$ in time $n^{1+1/k} \cdot \mathsf{poly}(\alpha, \log(n))$.*

The runtime of our algorithm improves as the gap between the density and $2^{k/3}$ becomes larger. Theorem 1.2 obtains extremely efficient runtime with a slight exponential sacrifice on the density.

**Theorem 1.2** (Informal). *Assume $\alpha \approx 2^{0.33 \cdot k}$ and $k, n$ sufficiently large. Then with high probability, we can sample an approximate uniform solution of $\Phi$ in time $n^{1+2^{-0.001 \cdot k}} \cdot \mathsf{poly}(\alpha, \log(n))$.*

Both Theorem 1.1 and Theorem 1.2 are the informal and special cases of the following Theorem 1.3,[4] which achieves a smooth interpolation between the slack $\xi$ on the density and the efficiency on the runtime. It also makes the "approximate uniform" precise by an explicit total variation distance measure $\varepsilon$.

---

[3]We do *not* claim that our result validates the intuition. On the one hand, it is very possible that our bounds can be further improved. On the other hand, the bounds for the worst-case model may also be far from the truth considering the hardness results [BGG+19].

[4]In the statement of Theorem 1.3, we only hide absolute constants in $\Omega(\cdot), o(\cdot)$ and fixed polynomial in $\mathsf{poly}(\cdot)$. These do not depend on any parameter we introduce.

**Theorem 1.3.** *There exists a Monte Carlo algorithm $\mathcal{A} = \mathcal{A}(\varepsilon, k, \alpha, n, \Phi)$ for $\varepsilon \in (0, 1)$, $k \geq 2^{20}$, and $n \geq 2^{\Omega(k)}$ such that the following holds: If*

$$\alpha \leq \frac{2^{k/3}}{k^{50}} \cdot \xi \quad \text{where } 2^{-k/8} \leq \xi \leq 1,$$

*then $\mathcal{A}$ runs in time*

$$(n/\varepsilon)^{1+\xi/k}/\xi \cdot \mathsf{poly}(k, \alpha, \log(n/\varepsilon)).$$

*Moreover, let $\mu'$ be the output distribution of $\mathcal{A}$ and let $\mu$ be a uniform solution of $\Phi$. Then*

$$\Pr_{\Phi} \left[ \Phi \text{ is not satisfiable } \vee \ d_{\mathsf{TV}} \left( \mu, \mu' \right) \leq \varepsilon \right] \geq 1 - o(1/n),$$

*where $d_{\mathsf{TV}}(\cdot, \cdot)$ is the total variation distance and $\mu$ is a uniform random solution of $\Phi$.*

The $o(1/n)$ factor in Theorem 1.3 can be improved to any $n^{-\Omega(1)}$ by slightly changing constants in our analysis for the structural properties in Section 3. Similarly, the denominator $k^{50}$ in the density bound or the $1/k$ on the exponent of the runtime bound can be polynomially improved by more refined calculation.

Our sampling algorithm can be turned into an efficient approximate counting algorithm. This can be achieved by executing the algorithm multiple times to get approximations for marginal probabilities of variables in partial assignments, then applying well-known reductions between marginals and total number of solutions. We refer interested readers to [GGGY21, Section 9] for detail.

Curiously, our result holds in a stronger sense that we allow adversaries to change the signs of the literals in the clauses, i.e., an adversary can add or remove negations arbitrarily. Indeed, we identify the good formula purely based on the structural properties of the underlying hypergraphs on variables, regardless of the negations. This feature may be of independent interests.

**Future Directions.** We highlight some interesting future directions regarding sampling solutions of random formulas:

- **Better Density Bounds.** Our sampling algorithm is efficient for density up to $2^{k/3}$. In contrast, the satisfiability and search thresholds are roughly $2^k$. We believe that there exist better sampling algorithms that goes beyond $2^{k/3}$. A milestone will be to get around $2^{k/2}$, which, if true, would match the hardness in the worst case setting [BGG$^+$19]. In fact, it is speculative that the sampling threshold is also near $2^k$, since the random $k$-CNF formula is locally sparse and, in the bounded-degree model, solutions of $k$-CNF formulas on linear hypergraphs admits efficient sampling for variable degree up to $2^k$ [HSZ19, QWZ22].

- **Random Monotone Formulas.** As mentioned above, our algorithm works even when the signs of the literals are chosen adversarially. This is partially due to our use of Lovász local lemma which is oblivious to the signs. It is possible that better algorithms arises from better understanding on the patterns of negations. Towards this direction, we ask if better density bounds are obtainable for random monotone $k$-CNF formulas, which should be the easiest due to its trivial satisfiability. For its bounded-degree counterpart, it is indeed known that the sampling threshold is $2^{k/2}$ [BGG$^+$19, HSZ19, QWZ22], much larger than the $2^{k/3}$ bound obtained here.

- **Better Error Bounds.** The $o(1/n)$ error bound in Theorem 1.3 can be easily improved to any $n^{-\Omega(1)}$. It is even imaginable to obtain a bound scales with $k$, say, $n^{-\sqrt{k}}$. However, it is not clear how to go beyond $n^{-\Omega(k)}$. This is because our analysis crucially replies on Lovász

4

local lemma which in turn needs an $\Omega(k)$ lower bound on the clause width, i.e., the number of distinct literals in a clause. Whereas, once the error bound becomes smaller than $n^{-\Omega(k)}$, we may get many clauses of very small width.

- **Small Input Regimes.** Our result holds for large inputs that has both large clause width $k \geq 2^{20}$ and large amount of variables $n \geq 2^{\Omega(k)}$. It is an intriguing question whether we can weaken these assumptions. The former large-$k$ assumption appears commonly in the study of satisfiability and search thresholds (See e.g., [DSS22, Coj10]), and there are non-rigorous arguments and experimental evidence [AZ08] showing the difficulty and distinction for small $k$'s, which may carry over to the sampling task as well. The second large-$n$ assumption comes from our pursuit for *highly efficient* algorithms. Indeed, if we are satisfied with arbitrary overhead on $k$ in the runtime, say, $2^{2^{O(k)}} \cdot n^{1+\xi/k}$, it can be removed as we can trivially go over all possible $2^n = 2^{2^{O(k)}}$ assignments when $n \leq 2^{O(k)}$. But it is not clear how to do it if we want $n^{1+o(1)} \cdot \mathsf{poly}(k, \alpha, \log(n))$ or even $\mathsf{poly}(n, k, \alpha)$ runtime.

## 1.2 Proof Overview

Our algorithm is inspired by a recursive sampling scheme recently developed in [AJ22, HWY22, HWY23]. We first identify the technical difficulties in applying the techniques from [FGYZ21, FHY21, JPV21, HSW21], which have proved successful in the worst-case model. Then we show how [GGGY21] circumvent some of the issues using techniques from [Moi19] and what makes their bound much worse than [Moi19]. Finally we discuss our approach and technique novelties leading to near $2^{k/3}$ density and almost-linear runtime.

**Bottlenecks in Previous Algorithms.** The algorithms in [FGYZ21, FHY21, JPV21, HSW21] are based on Markov chains. Recall that in the worst-case model, the variables have a worst-case degree bound $d$. Their algorithm can be summarized as follows: (1) Classify the variables as marked and unmarked ones. (2) Construct a Markov chain on the marked variables where each time we update a (random) marked variable based on its marginal distribution conditioned on the partial assignment at that point. (3) When the Markov chain on the marked variables mixes, we sample the unmarked variables to obtain a solution.

The core of their analysis is the *local uniformity* for the marked variables: Once we guarantee that every clause has enough unmarked variables, by Lovász local lemma [EL73, HSS11], the marginal distribution of a marked variable is close to an unbiased coin, assuming the unmarked variables are untouched and regardless of the value of the other marked variables. We also need to guarantee that every clause has enough marked variables to ensure that the update in Step (2) and the sampling in Step (3) are efficient. In addition, this marking needs to be provided in advance of the Markov chain, which makes the mark-vs-unmark trade-off static and thus restrict their final degree bounds.

The most challenging part is to establish bounds for the mixing time for Step (3). To this end, [FGYZ21, FHY21] rely on path coupling arguments, i.e., showing large contraction for one-step update of neighboring Markov chain states; and [JPV21, HSW21] uses information percolation arguments, i.e., bounding the probability of long-range uncoupling in the time series. Both these arguments face severe obstacles in the average-case model due to the existence of high-degree variables which appear with high probability and do not have the local uniformity property. As a consequence, the contraction in path coupling arguments could be vanishing, and the long-range uncoupling could actually appear.

Aside from the proof strategies, there is some evidence that this kind of one-step-update Markov chain relying on the local uniformity property may be slow mixing. Consider a star graph of degree

$D$, i.e., a node $v$ connecting to nodes $u_1, u_2, \ldots, u_D$. The node $v$ models a high-degree variable or a component consisting of mostly high-degree variables, and nodes $u_1, u_2, \ldots, u_D$ are the surrounding low-degree neighbors. Then it is likely that this structure appears in the underlying hypergraph of a random $k$-CNF formula for $D = \omega_{k,\alpha}(1)$ or even $D = \mathsf{poly}(k, \alpha) \cdot \log(n)$. Let $\sigma$ and $\sigma'$ be two distinct assignments that do not touch $v$. Now the Markov chain will ignore $v$ and only update $u_i$'s due to the local uniformity constraint. For each $u_i$, even if its current value is the same in $\sigma$ and $\sigma'$, the one-step-update may make it differ. The probability of this uncoupling is a small constant (independent of $n$) provided by the local uniformity, which means the estimate of the mixing time is $O(1)^D \gtrsim \mathsf{poly}(n)$.

Note that the recent independent works [GGGH22, CMM23] bypass this issue by making the Markov chain update more (actually, constant fraction of) variables a time. To argue the mixing time, they leverage recently developed *spectral independence* techniques. We refer interested readers to their paper for detail. Unfortunately, their bounds still suffer from the static mark-vs-unmark trade-off and are thus much weaker than our result.

**How [GGGY21] Circumvents the Barrier.** The algorithm in [GGGY21] does not involve Markov chains, and it samples an assignment by fixing a variable once at a time according to its (approximate) marginal distribution conditioned on the previous assignment.

To obtain the marginal distribution of a variable $v$, they adapt the linear programming framework from [Moi19]. Intuitively, starting from $v$, they gradually expand the possible values of its neighboring variables in a tree fashion. Using this tree, they formulate a system of linear inequalities regarding the marginal probabilities provided by the local uniformity property,[5] where the marginals of the leaf nodes can be directly computed. Then it is shown that any feasible solution to the linear system is a good approximation of the actually marginals, and in addition, it suffices to expand the tree up to logarithmic depth. Therefore, a good approximate of the marginal of $v$ can be obtained by solving the linear programming.

This approach can be carried out in the average-case model. In particular, the above star graph example is no longer an issue if the formulated linear system includes all the $2^D$ partial assignments on $u_1, \ldots, u_D$. Since $D$ is also upper bounded by $\mathsf{poly}(k, \alpha) \cdot \log(n)$ with high probability, we just expand the tree to this depth. This explains their runtime being $n^{\mathsf{poly}(k,\alpha)}$: They need to solve a linear system of size $2^{\mathsf{poly}(k,\alpha) \cdot \log(n)}$.

For the density bound, the analysis in [Moi19] already suffers from the loss in the static marking scheme required for local uniformity and to control the error of the linear system. In the average-case model, [GGGY21] needs to first separate the high-degree variables, and then impose a stronger local uniformity assumption on the rest to make sure the error analysis goes through. As a consequence, the bound in [GGGY21] is even worse than the one in [Moi19].

**How We Improve [GGGY21].** Our sampling algorithm follows the outline in [Moi19, GGGY21] by gradually fixing the variables towards a full assignment. However, we replace the linear programming framework with the recursive sampling framework recently developed in [AJ22, HWY22, HWY23], which can be seen as a dynamic marking scheme as opposed to the static one above. The benefit is two-fold: The runtime is significantly improved since we no longer need to solve giant linear systems, and the density bound is much better since the recursive sampling approach allows us to weaken the local uniformity assumption.

---

[5]In fact, the linear inequalities are about the ratio of the marginal probabilities. But since this is not important for us, we do not expand here.

The first step of our algorithm is to start with high-degree variables and include all the bad variables which are influenced by them and do not possess local uniformity properties. This part is similar to [GGGY21, CF14] but we tighten their analysis in the study of structure properties of the random formula. In particular, each remaining clause, after removing these bad variables, still has width $(1 - o(1)) \cdot k$.

To give a quantitative sense on the local uniformity property, we introduce $\theta \in (0, 1)$ as the parameter for maximum possible "marked" variables in a clause. Note that our algorithm does not compute a static marking, and thus $\theta k$ is only used to upper bound the number of fixed variables in a clause at any point (or equivalently, $(1 - \theta)k$ lower bounds the number of untouched variables in a clause). Then by Lovász local lemma [EL73, HSS11], the local uniformity parameter is

$$\delta \approx \alpha \cdot 2^{-(1-\theta)k}, \tag{1}$$

which means the correct marginal distribution $\mu_v$ conditioned on the previous assignment for any remaining good variable $v$ is $\delta$-close to an unbiased coin.

Now we sample $\mu_v$ sequentially for good variables $v$ as [GGGY21]. By local uniformity, we can already fix its value $\sigma(v)$ to 0/1 with probability $(1 - \delta)/2$ each, and set $\sigma(v) = \star$ for the remaining $\delta$ uncertainty. We denote this distribution as $\tau$. Since ultimately we need to complete the $\star$ to 0/1 to obtain a sample from $\mu_v$, we will need to sample from $\tau_v \propto \mu_v - \tau$. This part is similar to [HWY22]. With the *Bernoulli factory* technique [NP05, Hub16, DHKN21], samples from $\tau_v$ can be obtained efficiently provided samples from $\mu_v$. This alone is merely a self-referencing: Sampling from $\mu_v$ circles back to samples from $\mu_v$. But the trick here is to postpone sampling from $\tau_v$ and perform more sampling from $\tau$.

Let $v_1$ be a different variable with local uniformity property conditioned on $v = \star$. We can tentatively sample its value $\sigma(v_1) \sim \tau$, and, if $\sigma(v_1) = \star$, update it by $\sigma(v_1) \sim \tau_{v_1|v=\star}$. Then we turn to the next variable $v_2$, sample $\sigma(v_2) \sim \tau$, and update $\sigma(v_2) \sim \tau_{v_2|v=\star,v_1=\sigma(v_1)}$ if necessary. Iteratively doing so gives us $\sigma(v_1), \sigma(v_2), \ldots, \sigma(v_t)$. Now if we update $\sigma(v) \sim \tau_{v|v_1=\sigma(v_1),v_2=\sigma(v_2),\ldots,v_t=\sigma(v_t)}$, it follows the correct distribution $\tau_v$ in general by the law of conditional probability. The hope here is that, after fixing $\sigma(v_1), \ldots, \sigma(v_t)$, the CNF formula decomposes into components and $v$ belongs to a small one, which allows us to efficiently obtain samples from $\mu_{v|v_1=\sigma(v_1),v_2=\sigma(v_2),\ldots,v_t=\sigma(v_t)}$ using rejection sampling for Bernoulli factory. We remark that this algorithm incurs many recursions as, for example, sampling $\sigma(v_1) \sim \tau_{v_1|v=\star}$ will also be postponed and implemented by the same recursive sampling idea.

The correctness of the above marginal sampling algorithm is evident from the description and can be proved rigorously by induction. The difficulty lies in the efficiency analysis. Indeed, we face two issues regarding the runtime: (1) The recursion may dive too deep such that branches into too many possibilities, and (2) the final Bernoulli factory may still require exponential time. To address them, we keep track of the component $\mathcal{C}_{con}^\sigma$ containing variables and clauses we visited during the recursion and relate its size $|\mathcal{C}_{con}^\sigma|$ to the depth of the recursion and the efficiency of the final Bernoulli factory. By a similar analysis as [HWY22], we show that a deep recursion produces a large component $\mathcal{C}_{con}^\sigma$. Therefore, to address both (1) and (2), it suffices to truncate the program once $|\mathcal{C}_{con}^\sigma|$ exceeds certain size.

Then the issue comes back to the correctness: Is the output of the algorithm close to a uniform solution? Observe that the difference between the new algorithm and the original one only lies in the place where truncation happens. Therefore, it suffices to bound the probability that a large component $\mathcal{C}_{con}^\sigma$ appears in the original algorithm. To this end, we will construct a succinct witness $\mathcal{W}^\sigma$ that enjoys the following properties: (a) Each large $\mathcal{C}_{con}^\sigma$ gives rise to a witness $\mathcal{W}^\sigma$, (b) each fixed $\mathcal{W}$ appears as a witness of some $\sigma$ with small probability during the algorithm, and (c) there

are not many possible $\mathcal{W}$. The construction of $\mathcal{W}^\sigma$ is the place where we significantly deviate from (and simplify) the previous analysis and leverage the structural properties of random formulas.

Our witness $\mathcal{W}$ consists of two sets of clauses $\mathcal{C}_{\text{int}}$ and $\mathcal{C}_{\star\text{-int}}$.

(i) $\mathcal{C}_{\text{int}}$ contains some unsatisfied clauses.

This is helpful for Property (b). When we execute the algorithm and are about to fix a variable appearing in some clause $C \in \mathcal{C}_{\text{int}}$, the variable cannot be fixed to the bit that satisfies $C$. Thus intuitively, the probability that the algorithm proceeds in the direction consistent with $\mathcal{W}$ halves in this step.

(ii) $\mathcal{C}_{\star\text{-int}}$ contains some clauses containing $\star$'s.

Then similar to the $\mathcal{C}_{\text{int}}$ case, this intuitively requires the algorithm to go into the direction that assigns $\star$ from $\tau$ whenever we encounter a variable indicated as $\star$ in $\mathcal{C}_{\star\text{-int}}$. The proper transition probability in this step is governed by the local uniformity $\tau(\star) = \delta$.

(iii) $\mathcal{C}_{\star\text{-int}}$ connects $\mathcal{C}_{\text{int}}$ in the underlying hypergraph.

This is helpful for Property (c). Using structural properties of the random formula, it can be shown that

$$\# \text{ possible } \mathcal{W}\text{'s} \lesssim \mathsf{poly}(n) \cdot \alpha^{|\mathcal{W}|}. \tag{2}$$

Assume $\mathcal{W}$ is the witness for $\mathcal{C}_{\text{con}}^\sigma$, i.e., $\mathcal{W} = \mathcal{W}^\sigma$ from Property (a). Item (i) tells us to include more visited variables in $\sigma$, since each one of them represents a probability decay for Property (b). Recall that $\theta \in (0,1)$ controls the fraction of variables we can visit for each clause during the algorithm. Then we have a trivial bound: The number of visited variables in $\mathcal{C}_{\text{int}}$ is at most $\theta k \cdot |\mathcal{C}_{\text{int}}|$. Perhaps surprisingly, by the locally sparse properties of the random formula, we can almost achieve this bound! More precisely, we show that one can carefully select a subset of $\mathcal{C}_{\text{con}}^\sigma$ to form $\mathcal{C}_{\text{int}}$ such that the number of visited variables in $\mathcal{C}_{\text{int}}$ is at least $(\theta - o(1))k \cdot |\mathcal{C}_{\text{int}}|$, which means the accumulated probability drop from Item (i) is roughly

$$2^{-\theta k \cdot |\mathcal{C}_{\text{int}}|}. \tag{3}$$

Item (ii) also requires us to include more $\star$'s in $\sigma$ for Property (b). For this, we investigate the connectivity $\mathcal{C}_{\text{int}}$ inside $\mathcal{C}_{\text{con}}^\sigma$ and show that one can make it connected by only inserting clauses that contain $\star$'s. Thus, by including the minimum amount of such clauses in a spanning tree fashion, we can additionally guarantee that the number of visited $\star$'s in $\mathcal{C}_{\star\text{-int}}$ is at least $|\mathcal{C}_{\star\text{-int}}|$. This means the accumulated probability drop from Item (ii) is roughly

$$\delta^{|\mathcal{C}_{\star\text{-int}}|}. \tag{4}$$

Combining Equations (3) and (4), we derive the Property (b) of $\mathcal{W}$ as

$$\mathbf{Pr}[\text{encounter this } \mathcal{W}] \lesssim 2^{-\theta k |\mathcal{C}_{\text{int}}|} \cdot \delta^{|\mathcal{C}_{\star\text{-int}}|} \leq \max\left\{2^{-\theta k |\mathcal{W}|}, \delta^{|\mathcal{W}|}\right\}.$$

Now to offset the number of possible $\mathcal{W}$'s in the union bound, by (2) and (1), we need to ensure

$$\max\left\{2^{-\theta k}, \delta\right\} \approx \max\left\{2^{-\theta k}, \alpha \cdot 2^{-(1-\theta)k}\right\} \lesssim 1/\alpha \quad \text{and} \quad |\mathcal{W}| \gtrsim \log(n).$$

The former gives $\theta \leq 1/3$ and $\alpha \lesssim 2^{k/3}$ as foreshadowed. The latter, through some additional arguments, implies that the truncation threshold should be set to $\mathsf{poly}(k, \alpha) \cdot \log(n)$, which also explains why the above star graph example is not an obstacle here.

8

There are some technical difficulties that we choose to omit here for simplicity. For example, the locally sparse property only holds up to certain size, and we need additional pruning ideas to make sure our witness enjoys the property. After pruning, our witness is doomed to have an upper bound on its size. This means, through final union bound in the witness analysis, the distance between the algorithm's output and a uniform solution has an inevitable lower bound. Therefore, to handle the case where we want an extremely small output difference, we need another algorithm. Similarly, some of the structural properties we use require a lower bound on the density $\alpha$. Thus we also need a different algorithm for small densities. We fix these issues by analyzing the naive rejection sampling algorithm and carefully balancing parameters for different algorithms.

**Organization.** We give formal definitions in Section 2. Useful structural properties of random CNF formulas are provided in Section 3 and their proofs are deferred to Appendix A. In Section 4, we present the pre-processing algorithm to construct variable and clause separators. In Section 5, we analyze the naive rejection sampling on random formulas which gives the algorithms for the atypical setting. In Section 6, we introduce our main algorithms for the typical setting, the most technical part of which is the truncation analysis and is carried out in Section 7. Finally we put everything together and prove Theorem 1.3 in Section 8.

## 2 Preliminaries

We use $\mathbf{e} \approx 2.71828$ to denote the *natural* base, and we will frequently use the inequality $\binom{a}{b} \leq (\mathbf{e}a/b)^b$ for all $a, b \geq 0$ where $0^0$ is defined as 1. We use $\log(\cdot)$ and $\ln(\cdot)$ to denote the logarithm with base 2 and $\mathbf{e}$ respectively. For positive integer $n$, we use $[n]$ to denote the set $\{1, 2, \ldots, n\}$.

For a finite set $\mathcal{X}$ and a distribution $\mathcal{D}$ over $\mathcal{X}$, we use $x \sim \mathcal{D}$ to denote that $x$ is a random variable sampled from $\mathcal{X}$ according to distribution $\mathcal{D}$. We also use $x \sim \mathcal{X}$ when $\mathcal{D}$ is the uniform distribution.

**Asymptotics.** We only use $O(\cdot), \Omega(\cdot), o(\cdot), \omega(\cdot)$ to hide absolute constants that does not depend on any parameters we introduce. In addition, $\widetilde{O}(\cdot)$ is only used to bound algorithms' runtime which hides polynomial factors in $k, \alpha, \log(n/\varepsilon)$, i.e., $\widetilde{O}(f) = \mathsf{poly}(k, \alpha, \log(n/\varepsilon)) \cdot f$ for some fixed $\mathsf{poly}$.

**(Random) CNF Formula.** A CNF formula is a disjunction of clauses. Each clause is a conjunction of literals, and a literal is either a Boolean variable or the negation of a Boolean variable. Given a CNF formula $\Phi = (\mathcal{V}, \mathcal{C})$ with variable set $\mathcal{V}$ and clause set $\mathcal{C}$, we define the following measure for $\Phi$:

- The *width* is $k(\Phi) = \max_{C \in \mathcal{C}} |\mathsf{vbl}(C)|$, where $\mathsf{vbl}(C)$ denotes the variables that $C$ depends on.
- The *variable degree* is $d(\Phi) = \max_{v \in \mathcal{V}} |\{C \in \mathcal{C} \mid v \in \mathsf{vbl}(C)\}|$.
- The *constraint degree* is $\Delta(\Phi) = \max_{C \in \mathcal{C}} |\{C' \in \mathcal{C} \mid \mathsf{vbl}(C) \cap \mathsf{vbl}(C') \neq \emptyset\}|$.[6]
- The *maximum violation probability* is $p(\Phi) = \max_{C \in \mathcal{C}} \mathbf{Pr}[C(\sigma) = \mathsf{False}] = \max_{C \in \mathcal{C}} 2^{-|\mathsf{vbl}(C)|}$.

In addition, we use $\mu(\Phi)$ to denote the uniform distribution over the solutions of $\Phi$. Note that $\mu$ is well defined whenever $\Phi$ is satisfiable. In the rest of the paper, we will simply use $k, d, \Delta, p, \mu$ when $\Phi$ is clear from the context.

We use $\Phi(k, n, m)$ to denote a random $k$-CNF formula on $n$ variables and $m$ clauses, where $\mathcal{V} = \{v_1, v_2, \ldots, v_n\}$ is the variable set, $\mathcal{C} = \{C_1, C_2, \ldots, C_m\}$, and each clause is an independent

---

[6]Note that in our definition, $\Delta$ is one plus the maximum degree of the dependency graph of $\Phi$.

disjunction of $k$ literals chosen independently and uniformly from $\{v_1, v_2, \ldots, v_n, \neg v_1, \neg v_2, \ldots, \neg v_n\}$. We will simply use $\Phi$ to denote $\Phi(k, n, m)$ when context is clear.

**Partial Assignments and Restrictions.** Our algorithm will sample an assignment by gradually fixing coordinates. To this end, we will work with partial assignments and restrictions of the formula on partial assignments. We use $\star$ for *unaccessed variables* and use $\bigstar$ for *accessed but unassigned variables* and a partial assignment $\sigma$ lies in the space $\{0, 1, \bigstar, \star\}^{\mathcal{V}}$. We define

$$\Lambda(\sigma) = \{v \in \mathcal{V} \mid \sigma(v) \in \{\bigstar, \star\}\}$$

to be the set of unassigned variables. We then abuse the notation to say $C(\sigma) = \mathsf{True}$ if fixing $v$ to $\sigma(v)$ for all $v \notin \Lambda(\sigma)$ already satisfies $C$.

For a partial assignment $\sigma$, let $\Phi^\sigma = (\mathcal{V}^\sigma, \mathcal{C}^\sigma)$ be the CNF formula after we fix $v$ to be $\sigma(v)$ for each $v \notin \Lambda(\sigma)$. Note that $\mathcal{V}^\sigma = \Lambda(\sigma)$ and each clause in $\mathcal{C}^\sigma$ depends only on variables in $\Lambda(\sigma)$.

We use $\mu^\sigma = \mu(\Phi^\sigma)$ to denote the uniform distribution over solutions of $\Phi^\sigma$. For each $v \in \Lambda(\sigma)$, we write $\mu_v^\sigma$ as the marginal distribution of $v$ under $\mu^\sigma$. Then $\mu_v^\sigma(b)$ denotes the probability that $v$ is fixed to $b \in \{0, 1\}$ under $\mu_v^\sigma$. For multiple variables $S \subseteq \Lambda(\sigma)$, we use $\mu_S^\sigma$ to denote the marginal distribution of $S$ under $\mu^\sigma$.

**Incidence Graphs.** Given a formula $\Phi = (\mathcal{V}, \mathcal{C})$, we define two incidence graphs $G_\Phi$ and $H_\Phi$:

- The vertex set of $G_\Phi$ is $\mathcal{C}$, and two clauses $C_1, C_2 \in \mathcal{C}$ are adjacent iff $\mathsf{vbl}(C_1) \cap \mathsf{vbl}(C_2) \neq \emptyset$. We say a set $S \subseteq \mathcal{C}$ of clauses is connected if the induced sub-graph $G_\Phi[S]$ is connected.

- The vertex set of $H_\Phi$ is $\mathcal{V}$, and two variables $v, v' \in \mathcal{V}$ are adjacent iff there exists some $C \in \mathcal{C}$ with $v, v' \in \mathsf{vbl}(C)$. We say a set $T \subseteq \mathcal{V}$ of variables is connected if the induced sub-graph $H_\Phi[T]$ is connected.

**Lovász Local Lemma.** The celebrated Lovász local lemma [EL73] provides a sufficient condition for the existence of a solution of a constraint satisfaction problem. Here we use a more general version for CNF formulas due to [HSS11]:

**Theorem 2.1** ([HSS11, Theorem 2.1]). *Let $\Phi = (\mathcal{V}, \mathcal{C})$ be a CNF formula. If $\mathbf{e}p\Delta \leq 1$, then $\Phi$ is satisfiable. Moreover, for any event $B$ (not necessarily from $\mathcal{C}$) we have*

$$\Pr_{\sigma \sim \mu}[B(\sigma) = \mathsf{True}] \leq (1 - \mathbf{e}p)^{-|\Gamma(B)|} \Pr_{\sigma \sim \{0,1\}^{\mathcal{V}}}[B(\sigma) = \mathsf{True}],$$

*where $\Gamma(B) = \{C \in \mathcal{C} \mid \mathsf{vbl}(C) \cap \mathsf{vbl}(B) \neq \emptyset\}$.*

## 3  Properties of Random CNF Formulas

For the rest of the paper, we will use $\Phi$ to denote a random $k$-CNF formula on $n$ variables $\mathcal{V} = \{v_1, \ldots, v_n\}$ and $m$ clauses $\mathcal{C} = \{C_1, \ldots, C_m\}$. We reserve $\alpha = \alpha(\Phi) = m/n$ as the *density* of $\Phi$.

For convenience and later reference, we list desirable properties of $\Phi$ here. In the next sections, we will assume $\Phi$ satisfies these structural properties, which happens with high probability, and prove the correctness and efficiency of our algorithm.

We first cite the following celebrated satisfiability result.

**Theorem 3.1** ([DSS22, Theorem 1]). *For $k \geq \Omega(1)$, $\Phi$ has a sharp satisfiability threshold $\alpha_\star(k)$ such that for all $\varepsilon > 0$, it holds that*

$$\lim_{n \to +\infty} \mathbf{Pr}\left[\Phi(k, n, m) \text{ is satisfiable}\right] = \begin{cases} 1 & \text{if } \alpha \leq \alpha_\star(k) - \varepsilon, \\ 0 & \text{if } \alpha \geq \alpha_\star(k) + \varepsilon. \end{cases}$$

*Roughly, $\alpha_\star(k) = 2^k \ln(2) - (1 + \ln(2))/2 + o_k(1) < 2^k$ as $k \to +\infty$.*[7]

By Theorem 3.1, it is reasonable to focus our attention to the case where $\alpha \leq 2^{O(k)}$. In particular, this justifies our assumption $\alpha \leq 2^k$ used below. We remark that the proofs for the following properties are similar to the ones in [GGGY21, CF14]. Therefore we defer them to Appendix A.

The first property states that every clause in $\Phi$ has at most two duplicate variables.

**Proposition 3.2.** *Assume $\alpha \leq 2^k$ and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, $|\mathsf{vbl}(C)| \geq k - 2$ holds for every $C \in \mathcal{C}$.*

Intuitively, Proposition 3.3 and Proposition 3.4 show that typically the clauses in $\Phi$ are spread out in that they do not share many common variables.

**Proposition 3.3.** *Let $\eta = \eta(k) > 0$ be a parameter. Assume $\alpha \leq 2^k$, $\frac{k}{\log(k)} \geq 14 \cdot \left(1 + \frac{1}{\eta}\right)$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, the following holds:*

1. *For every $\mathcal{V}' \subset \mathcal{V}$ with $1 \leq |\mathcal{V}'| \leq n/2^{k/\log(k)}$, we have $|\{C \in \mathcal{C} \mid \mathsf{vbl}(C) \subseteq \mathcal{V}'\}| \leq (1+\eta)|\mathcal{V}'|/k$.*
2. *For every $\mathcal{C}' \subset \mathcal{C}$ with $1 \leq |\mathcal{C}'| \leq n/2^{2k/\log(k)}$, we have $\left|\bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)\right| \geq k|\mathcal{C}'|/(1+\eta)$.*

**Proposition 3.4.** *Let $\eta = \eta(k) \in (0, 1)$ be a parameter. Assume $\alpha \leq 2^k$, $\frac{k}{\log(k)} \geq 14 \cdot \left(1 + \frac{1}{\eta}\right)$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, the following holds: For any $b \geq \eta$ and every $\mathcal{V}' \subset \mathcal{V}$ with $1 \leq |\mathcal{V}'| \leq n/2^{3k/\log(k)}$, we have*

$$|\mathcal{V}'| \geq (b - \eta)k \cdot \left|\left\{C \in \mathcal{C} \mid |\mathsf{vbl}(C) \cap \mathcal{V}'| \geq bk\right\}\right|.$$

Recall our definition of incidence graph $G_\Phi$ from Section 2, we can bound the number of induced connected sub-graphs in $G_\Phi$.

**Proposition 3.5.** *With probability $1 - o(1/n)$ over the random $\Phi$, the following holds: For every $C \in \mathcal{C}$ and $\ell \geq 1$, there are at most $\alpha^2 n^4 (\mathbf{e}k^2\alpha)^\ell$ many connected sets of clauses in $G_\Phi$ with size $\ell$ containing $C$.*

In terms of incidence graph $H_\Phi$, we can bound the expansion of any connected set.

**Proposition 3.6.** *Assume $k \geq 30$ and $\alpha \geq 1/k^3$. Then with probability $1 - o(1/n)$ over the random $\Phi$, the following holds: For any $\mathcal{V}' \subseteq \mathcal{V}$ connected in $H_\Phi$, we have*

$$\left|\{v \in \mathcal{V} \mid v \in \mathcal{V}' \text{ or } v \text{ is adjacent to } \mathcal{V}'\}\right| \leq 3k^4\alpha \cdot \max\left\{|\mathcal{V}'|, \lfloor k\log(n)\rfloor\right\}.$$

Given a set of clauses $\mathcal{C}' \in \mathcal{C}$ and a variable $v \in \mathcal{V}$, we define the degree of $v$ in $\mathcal{C}'$ as $\deg_{\mathcal{C}'}(v) = |\{C \in \mathcal{C}' \mid v \in \mathsf{vbl}(C)\}|$. Then $d(\Phi) = \max_{v \in \mathcal{V}} \deg_{\mathcal{C}}(v)$. We first note a classical bound (See e.g., [RS98, Theorem 1]) on $d(\Phi)$.

---

[7]The explicit value of $\alpha_\star(k)$ is characterized by a complicated proposition in [DSS22]. We omit it here to simplify the statement. This asymptotic estimation is given by [KKKS98] as an upper bound and by [COP16] as a lower bound.

**Proposition 3.7.** *With probability $1 - o(1/n)$ over the random $\Phi$, we have $d(\Phi) \leq 4k\alpha + 6\log(n)$.*

We also need the following control over the number of high-degree variables.

**Proposition 3.8.** *Let $D = D(k, \alpha)$ be a parameter satisfying $D \geq 8k(\alpha + 1)$. Assume $k \geq 2$, $\alpha \leq 2^k$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, we have*

$$|\{v \in \mathcal{V} \,|\, \deg_{\mathcal{C}}(v) \geq D\}| \leq n/2^{4k}.$$

We can also bound the fraction of high-degree variables in any connected set.

**Proposition 3.9.** *Let $D = D(k, \alpha)$ be a parameter satisfying $6k^7(\alpha + 1) \leq D \leq 2^{2k}$. Assume $k \geq 2^{10}$, $1/k^3 \leq \alpha \leq 2^k$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, the following holds: Let $\mathcal{V}' \subseteq \mathcal{V}$ be connected in $H_\Phi$ and $|\mathcal{V}'| \geq \log(n)$. Then*

$$\left|\left\{v \in \mathcal{V}' \,\middle|\, \deg_{\mathcal{C}}(v) \geq D\right\}\right| \leq |\mathcal{V}'|/k^2.$$

Finally, the following proposition characterizes peeling procedures: It shows that the process of introducing new variables by including more clauses should stop soon.

**Proposition 3.10.** *Assume $k \geq 12$, $\alpha \leq 2^k$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, the following holds: Fix an arbitrary $\mathcal{C}' \subset \mathcal{C}$ with $|\mathcal{C}'| \leq n/2^{4k}$. Let $C_{i_1}, \ldots, C_{i_\ell} \in \mathcal{C} \backslash \mathcal{C}'$ be clauses with distinct indices. For each $s \in [\ell]$, define $\mathcal{V}_s = \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C) \cup \bigcup_{j=1}^{s-1} \mathsf{vbl}(C_{i_j})$. If $|\mathsf{vbl}(C_{i_s}) \cap \mathcal{V}_s| \geq 6$ holds for all $s \in [\ell]$, then $\ell \leq |\mathcal{C}'|$.*

## 3.1 Good and Nice Instances

At this point, we can assume $\Phi$ satisfies certain structural properties which exist with high probability over the random $\Phi$.

To be specific, we define the following Definition 3.11 and Definition 3.12: The former provides structural properties for $\Phi$ when $\alpha$ has an upper bound, and the latter guarantees more structural properties by further assuming $\alpha \geq 1/k^3$. For convenience, we include $\eta$ and $D$ to be consistent with Section 3.

**Definition 3.11** (Good Instances). We say $(\Phi, k, \alpha, n, \xi, \eta, D)$ is *good* if:

- $k \geq 2^{20}$, $n \geq 2^{\Omega(k)}$, $2^{-k/8} \leq \xi \leq 1$, and $\alpha \leq \xi \cdot 2^{k/3}/k^{50}$.
- $\eta = 15\log(k)/k$ and $D = k^8(\alpha + 1)/\xi$.
- $\Phi$ is satisfiable and has the properties in Proposition 3.2, Proposition 3.3, Proposition 3.4, Proposition 3.5, Proposition 3.7, Proposition 3.8, and Proposition 3.10.

**Definition 3.12** (Nice Instances). We say $(\Phi, k, \alpha, n, \xi, \eta, D)$ is *nice* if:

- $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good.
- $\alpha \geq 1/k^3$ and $\Phi$ has properties in Proposition 3.6 and Proposition 3.9 additionally.

**Remark 3.13.** By the choice of $\eta$ and $k \geq 2^{20}$, $\eta$ satisfies $\frac{k}{\log(k)} \geq 14 \cdot \left(1 + \frac{1}{\eta}\right)$. Since $2^{-k/8} \leq \xi \leq 1$, $k \geq 2^{20}$, and $\alpha \leq 2^{k/3} \cdot \xi$, we always have $6k^7(\alpha + 1) \leq D \leq 2^{2k}$. This means that $\eta$ and $D$ are consistent with the structural statements in Section 3.

By the bounds in Section 3 and Remark 3.13, we can indeed focus on good/nice instances. Though checking whether it is indeed a good/nice instance may actually need exponential time, we will not do it in our algorithm. Instead, we will assume the input enjoys the property, then run algorithm anyways and terminate it upon the prescribed maximum runtime. The correctness of our algorithm is only guaranteed when the input is actually good/nice.

**Corollary 3.14.** *Assume $k, \alpha, n, \xi, \eta, D$ satisfy the relations in Definition 3.11 (resp., Definition 3.12). Then with probability $1 - o(1/n)$ over the random $\Phi$, either $\Phi$ is not satisfiable or $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good (resp., nice).*

We remark that though Theorem 3.1 asserts that the satisfiability probability of $\Phi$ approaches 1 as $n$ goes to infinity, it only holds for $k$ sufficiently large (potentially much larger than $2^{20}$ in our setting). In addition, it does not control the convergence rate. Therefore we cannot simply say we have good/nice instances with probability $1 - o(1/n)$.

# 4    Separating High-Degree Variables

Define $\mathsf{HD}(\mathcal{V}') = \{v \in \mathcal{V}' \mid \deg_{\mathcal{C}}(v) \geq D\}$ to be the set of high-degree variables in $\mathcal{V}'$. Similar to [GGGY21], our algorithm will start with high-degree variables and propagates them to form a separator. We use $\mathcal{V}_{\mathsf{sep}}$ and $\mathcal{C}_{\mathsf{sep}}$ to denote the variable separators and clause separators obtained from $\texttt{ConstructSep}(\mathcal{V})$ respectively.

---

**Algorithm 1:** The $\texttt{ConstructSep}$ Algorithm

**Input:** Variables $\mathcal{V}' \subseteq \mathcal{V}$
**Output:** Variable separators $\mathcal{V}'_{\mathsf{sep}} \subseteq \mathcal{V}$ and clause separators $\mathcal{C}'_{\mathsf{sep}} \subseteq \mathcal{C}$

1  Initialize $\mathcal{V}'_{\mathsf{sep}} \leftarrow \mathsf{HD}(\mathcal{V}')$ and $\mathcal{C}'_{\mathsf{sep}} \leftarrow \emptyset$
2  **while** $\exists C \in \mathcal{C} \setminus \mathcal{C}'_{\mathsf{sep}}$ *such that* $\left|\mathsf{vbl}(C) \cap \mathcal{V}'_{\mathsf{sep}}\right| \geq 2\eta k$ **do**
3  $\quad$ Update $\mathcal{V}'_{\mathsf{sep}} \leftarrow \mathcal{V}'_{\mathsf{sep}} \cup \mathsf{vbl}(C)$ and $\mathcal{C}'_{\mathsf{sep}} \leftarrow \mathcal{C}'_{\mathsf{sep}} \cup \{C\}$
   **end**
4  **return** $\mathcal{V}'_{\mathsf{sep}}$ *and* $\mathcal{C}'_{\mathsf{sep}}$

---

By dynamically monitoring and updating $\left|\mathsf{vbl}(C) \cap \mathcal{V}'_{\mathsf{sep}}\right|$ for each $C \in \mathcal{C}$, Algorithm 1 can be done efficiently.

**Fact 4.1.** *The runtime of $\texttt{ConstructSep}(\mathcal{V}')$ is $\widetilde{O}(n)$ for any $\mathcal{V}' \subseteq \mathcal{V}$.*

Here we list some useful properties regarding the variable separators and clause separators for future referencing.

**Fact 4.2.** $\mathcal{V}'_{\mathsf{sep}} \subseteq \mathcal{V}_{\mathsf{sep}}$ *and* $\mathcal{C}'_{\mathsf{sep}} \subseteq \mathcal{C}_{\mathsf{sep}}$ *hold for any* $\mathcal{V}' \subseteq \mathcal{V}$.

We first bound the number of variable separators in terms of the number of high-degree variables.

**Lemma 4.3.** *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good. Then $|\mathcal{V}'_{\mathsf{sep}}| \leq 2|\mathcal{V}'|/\eta$ holds for any $\mathcal{V}' \subseteq \mathsf{HD}(\mathcal{V})$.*

*Proof.* By Proposition 3.8, $|\mathcal{V}'| \leq |\mathsf{HD}(\mathcal{V})| \leq n/2^{4k} < n/2^{3k/\log(k)}$. Let

$$\mathcal{C}' = \left\{ C \in \mathcal{C} \mid |\mathsf{vbl}(C) \cap \mathcal{V}'| \geq 2\eta k \right\}.$$

Thus by Proposition 3.4 with $b = 2\eta$, we have $|\mathcal{C}'| \leq |\mathcal{V}'|/(\eta k)$. Moreover, $|\mathcal{C}'| \leq |\mathcal{V}'| \leq n/2^{4k}$.

Observe that, starting from $\mathcal{C}'$, each clause newly added to $\mathcal{C}'_{\mathsf{sep}}$ intersects at least $2\eta k \geq 6$ variables with existing clauses. Then by Proposition 3.10 with $C_{i_1}, \ldots, C_{i_\ell}$ being $\mathcal{C}'_{\mathsf{sep}} \setminus \mathcal{C}'$, we have $|\mathcal{C}'_{\mathsf{sep}} \setminus \mathcal{C}'| \leq |\mathcal{C}'|$, which implies $|\mathcal{C}'_{\mathsf{sep}}| \leq 2|\mathcal{C}'| \leq 2|\mathcal{V}'|/(\eta k)$. Thus $|\mathcal{V}'_{\mathsf{sep}}| \leq k|\mathcal{C}'_{\mathsf{sep}}| \leq 2|\mathcal{V}'|/\eta$. $\qquad\square$

**Lemma 4.4** ([GGGY21, Lemma 8.9])**.** *Let $\mathcal{V}' \subseteq \mathcal{V}_{\mathsf{sep}}$ be an arbitrary maximal connected component in $H_\Phi[\mathcal{V}_{\mathsf{sep}}]$. Then $\mathcal{V}'_{\mathsf{sep}} = \mathcal{V}'$.*

**Lemma 4.5.** *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good. Let $\mathcal{V}' \subseteq \mathcal{V}_{\mathsf{sep}}$ consist of maximal connected components in $H_\Phi[\mathcal{V}_{\mathsf{sep}}]$. Then $|\mathcal{V}'| \leq 2|\mathsf{HD}(\mathcal{V}')|/\eta$.*

*Proof.* Note that it suffices to prove the bound for every maximal connected component in $H_\Phi[\mathcal{V}_{\mathsf{sep}}]$ and then add them up. Therefore we assume without loss of generality $\mathcal{V}'$ is connected in $H_\Phi[\mathcal{V}_{\mathsf{sep}}]$.

Observe that $\mathcal{V}'_{\mathsf{sep}}$ and $\mathcal{C}'_{\mathsf{sep}}$ equal the output of $\mathtt{ConstructSep}(\mathsf{HD}(\mathcal{V}'))$. By Lemma 4.4, we have $\mathcal{V}' = \mathcal{V}'_{\mathsf{sep}}$ and the desired bound follows immediately from Lemma 4.3. $\square$

Now we bound the fraction of $\mathcal{V}_{\mathsf{sep}}$ in any large connected component in $H_\Phi$.

**Lemma 4.6.** *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is nice. Let $\mathcal{V}' \subseteq \mathcal{V}$ be connected in $H_\Phi$ of size $|\mathcal{V}'| \geq \log(n)$. Then $|\mathcal{V}' \cap \mathcal{V}_{\mathsf{sep}}| \leq |\mathcal{V}'|/k$.*

*Proof.* Let $\mathcal{V}_1, \ldots, \mathcal{V}_\ell \subseteq \mathcal{V}_{\mathsf{sep}}$ be distinct maximal connected components in $H_\Phi[\mathcal{V}_{\mathsf{sep}}]$ and they intersect $\mathcal{V}'$. Let $\widetilde{\mathcal{V}} = \mathcal{V}' \cup \mathcal{V}_1 \cup \cdots \cup \mathcal{V}_\ell$. Then $\widetilde{\mathcal{V}}$ is connected in $H_\Phi$ and $\mathsf{HD}(\widetilde{\mathcal{V}}) = \mathsf{HD}(\mathcal{V}') \cup \mathsf{HD}(\mathcal{V}_1) \cup \cdots \cup \mathsf{HD}(\mathcal{V}_\ell)$.

Now by Lemma 4.5, we have $|\mathcal{V}_i| \leq 2|\mathsf{HD}(\mathcal{V}_i)|/\eta$. By Proposition 3.9, we also have $|\mathsf{HD}(\widetilde{\mathcal{V}})| \leq |\widetilde{\mathcal{V}}|/k^2$. Note that $\eta k \geq 2$, we have

$$|\widetilde{\mathcal{V}} \cap \mathcal{V}_{\mathsf{sep}}| = \sum_{i=1}^\ell |\mathcal{V}_i| \leq \frac{2}{\eta} \sum_{i=1}^\ell |\mathsf{HD}(\mathcal{V}_i)| \leq \frac{2|\mathsf{HD}(\widetilde{\mathcal{V}})|}{\eta} \leq \frac{2|\widetilde{\mathcal{V}}|}{\eta k^2} \leq \frac{|\widetilde{\mathcal{V}}|}{k}.$$

Since $\widetilde{\mathcal{V}} \setminus \mathcal{V}' \subset \mathcal{V}_{\mathsf{sep}}$, we have

$$\frac{|\mathcal{V}' \cap \mathcal{V}_{\mathsf{sep}}|}{|\mathcal{V}'|} \leq \frac{|\mathcal{V}' \cap \mathcal{V}_{\mathsf{sep}}| + |\widetilde{\mathcal{V}} \setminus \mathcal{V}'|}{|\mathcal{V}'| + |\widetilde{\mathcal{V}} \setminus \mathcal{V}'|} = \frac{|\widetilde{\mathcal{V}} \cap \mathcal{V}_{\mathsf{sep}}|}{|\widetilde{\mathcal{V}}|} \leq \frac{1}{k}. \qquad \square$$

As a corollary, we obtain the following bound on the fraction of $\mathcal{C}_{\mathsf{sep}}$ in any large connected component in $G_\Phi$.

**Corollary 4.7.** *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is nice. Let $\mathcal{C}' \subseteq \mathcal{C}$ be connected in $G_\Phi$ of size $|\mathcal{C}'| \geq \log(n)$. Then $|\mathcal{C}' \cap \mathcal{C}_{\mathsf{sep}}| \leq (1 + \eta)|\mathcal{C}'|/k$.*

*Proof.* Let $\mathcal{V}' = \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)$. Then $\mathcal{V}'$ is connected in $H_\Phi$.

First we prove for the case $|\mathcal{C}'| \leq n/2^{2k/\log(k)}$. By Item 2 of Proposition 3.3, we have $|\mathcal{V}'| \geq k|\mathcal{C}'|/(1+\eta) \geq \log(n)$. Then by Lemma 4.6, we have $|\mathcal{V}' \cap \mathcal{V}_{\mathsf{sep}}| \leq |\mathcal{V}'|/k$. Since $\mathcal{C}' \cap \mathcal{C}_{\mathsf{sep}}$ supports on $\mathcal{V}' \cap \mathcal{V}_{\mathsf{sep}}$, applying Item 1 of Proposition 3.3, we have

$$|\mathcal{C}' \cap \mathcal{C}_{\mathsf{sep}}| \leq \frac{1+\eta}{k} \cdot |\mathcal{V}' \cap \mathcal{V}_{\mathsf{sep}}| \leq \frac{1+\eta}{k^2} \cdot |\mathcal{V}'| \leq \frac{1+\eta}{k} \cdot |\mathcal{C}'|,$$

where $|\mathcal{V}'| \leq k|\mathcal{C}'| \leq n/2^{k/\log(k)}$ as required.

Now we turn to the case $|\mathcal{C}'| \geq n/2^{2k/\log(k)}$. By Lemma 4.5 and Proposition 3.8, we have

$$|\mathcal{V}_{\mathsf{sep}}| \leq \frac{2|\mathsf{HD}(\mathcal{V}_{\mathsf{sep}})|}{\eta} = \frac{2|\mathsf{HD}(\mathcal{V})|}{\eta} \leq \frac{2 \cdot n}{\eta \cdot 2^{4k}}.$$

Since $\frac{k}{\log(k)} \geq 14 \cdot \left(1 + \frac{1}{\eta}\right)$, we have $|\mathcal{V}_{\mathsf{sep}}| \leq n/2^{3k/\log(k)}$. Meanwhile, $\mathcal{C}_{\mathsf{sep}}$ supports on $\mathcal{V}_{\mathsf{sep}}$. Thus by Proposition 3.4 with $b = 1$, we have

$$|\mathcal{C}_{\mathsf{sep}}| \leq \frac{|\mathcal{V}_{\mathsf{sep}}|}{(1-\eta)k} \leq \frac{2 \cdot n}{\eta(1-\eta)k \cdot 2^{4k}} \leq \frac{n}{k \cdot 2^{2k/\log(k)}} \leq \frac{|\mathcal{C}'|}{k}.$$

Thus $|\mathcal{C}' \cap \mathcal{C}_{\mathsf{sep}}| \leq |\mathcal{C}_{\mathsf{sep}}| \leq |\mathcal{C}'|/k$. $\square$

# 5 The Naive Rejection Sampling Algorithm

The naive way to sample a solution is the rejection sampling algorithm, where we simply sample a uniform assignment and check if it happens to be a solution.

Starting with a (possibly empty) partial assignment $\sigma$, we can factorize $\Phi^\sigma$ into maximal connected components $\Phi_1, \Phi_2, \ldots$, where each $\Phi_i$ supports on disjoint subsets of the unassigned variables $\Lambda(\sigma)$. Then $\mu^\sigma = \mu_1 \times \mu_2 \times \cdots$ is a product distribution where $\mu_i$ is the uniform distribution over solutions of $\Phi_i$.

Now assume we want to get a sample from $\mu_S^\sigma$, i.e., the marginal distribution of variables in $S \subseteq \Lambda(\sigma)$ in a uniform solution of $\Phi^\sigma$. Assume $S = S_1 \cup S_2 \cup \cdots$ and each $S_i$ is contained in the support of $\Phi_i$. Then it suffices to get a sample from the marginal distribution of $S_i$ under $\mu_i$ for each $i$ independently and glue them together. This is formalized in Algorithm 2.

Recall that $\Lambda(\sigma)$ is the set of unassigned (i.e., $\bigstar$ or $\star$) variables in $\sigma$. Our rejection sampling algorithm does not distinguish $\bigstar$ and $\star$.

---

**Algorithm 2:** The `RejectionSampling` Algorithm

    **Input:** $\sigma \in \{0, 1, \bigstar, \star\}^{\mathcal{V}}$ and $S \subseteq \Lambda(\sigma)$
    **Output:** A random assignment distributed as $\mu_S^\sigma$
**1** Let $\Phi_i = (\mathcal{V}_i, \mathcal{C}_i), i = 1, 2, \ldots$ be the maximal connected components in $\Phi^\sigma$ intersecting $S$
**2** **foreach** $\Phi_i$ **do**
**3**     **repeat** Sample $\pi(\mathcal{V}_i) \sim \{0, 1\}^{\mathcal{V}_i}$ **until** $\pi(\mathcal{V}_i)$ *is a solution of* $\Phi_i$
    **end**
**4** **return** $\pi(S)$

---

We first note the simple correctness guarantee of Algorithm 2.

**Fact 5.1.** *If $\Phi^\sigma$ is satisfiable, then* `RejectionSampling`$(\sigma, S)$ *terminates almost surely and has output distribution exactly $\mu_S^\sigma$.*

To analyze the efficiency, we will make the following assumption on the partial assignment and it will be preserved throughout our algorithm. The intuition here is that, the partial assignment will not touch $\mathcal{C}_{\mathsf{sep}}$ which involves high-degree variables, and for the other clauses it leaves enough number of variables alive that guarantees satisfiability and efficient sampling using Theorem 2.1.

**Assumption 5.2.** $\mathcal{V}_{\mathsf{sep}} \subseteq \Lambda(\sigma)$ *and for every clause $C \in \mathcal{C} \setminus \mathcal{C}_{\mathsf{sep}}$, either $C(\sigma) = \mathsf{True}$ or $|\mathsf{vbl}(C) \cap \Lambda(\sigma) \setminus \mathcal{V}_{\mathsf{sep}}| \geq k'$ for some $k' \leq (1 - 2\eta)k - 2$.*

We remark that the condition $k' \leq (1-2\eta)k-2$ is for analysis convenience and is also reasonable considering Proposition 3.2 and Lines 2 and 3 of `ConstructSep`$(\mathcal{V})$. Later we will use it with $k' = (1 - 2\eta)k - 2$ and $k' = (2/3 - 2\eta)k$ respectively in different scenarios.

**Lemma 5.3.** *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good and $\sigma$ satisfies Assumption 5.2. If $\mathbf{e}2^{-k'} \cdot kD \leq 1$, then $\Phi^\sigma$ is satisfiable. Moreover, for each $b \in \{0, 1\}$ and $v \in \Lambda(\sigma) \setminus \mathcal{V}_{\mathsf{sep}}$, we have*

$$\frac{1 - \mathbf{e}2^{-k'}D}{2} \leq \mu_v^\sigma(b) \leq \frac{1 + \mathbf{e}2^{-k'}D}{2}.$$

*Proof.* Note that clauses in $\mathcal{C}_{\mathsf{sep}}$ only depend on $\mathcal{V}_{\mathsf{sep}}$. Since $\Phi$ is satisfiable, there exists a partial assignment $\pi$ extending $\sigma$ by fixing values of $\mathcal{V}_{\mathsf{sep}}$ to 0/1 and satisfying all clauses in $\mathcal{C}_{\mathsf{sep}}$.

Observe that $\pi$ only additionally fixes variables in $\mathcal{V}_{\text{sep}}$. Thus $\Lambda(\pi) = \Lambda(\sigma) \setminus \mathcal{V}_{\text{sep}}$. Now it suffices to show for any such $\pi$, $\Phi^\pi$ is satisfiable and we have

$$\frac{1-\delta}{2} \leq \mu_v^\sigma(b \,|\, \pi) = \mu_v^\pi(b) \leq \frac{1+\delta}{2},$$

where $\mu_v^\sigma(\cdot \,|\, \pi)$ is $\mu_v^\sigma$ conditioned on $\pi(\Lambda(\sigma) \cap \mathcal{V}_{\text{sep}})$.

Since each clause $C \in \mathcal{C}^\pi$ satisfies $|\mathsf{vbl}(C)| \geq k'$ where $\mathsf{vbl}(C) \subseteq \Lambda(\pi)$ is the set of remaining variables. Thus

$$\Pr_{\pi' \sim \{0,1\}^{\Lambda(\pi)}} \left[ C(\pi') = \mathsf{False} \right] \leq 2^{-k'}.$$

Since $\Lambda(\pi) \cap \mathcal{V}_{\text{sep}} = \emptyset$, every variable in $\Lambda(\pi)$ has variable degree at most $D$ in $\Phi^\pi$, and the constraint degree of $\Phi^\pi$ is at most $kD$.

Assuming $\mathbf{e}2^{-k'} \cdot kD \leq 1$ and by Theorem 2.1, $\Phi^\pi$ is satisfiable. Moreover, with $B$ being event "$v$ is assigned to $b$" which correlates with at most $D$ clauses in $\Phi^\pi$, we have

$$\mu_v^\pi(b) \leq \frac{\left(1 - \mathbf{e}2^{-k'}\right)^{-D}}{2} \leq \frac{1 + \mathbf{e}2^{-k'}D}{2}$$

and the other direction follows from $\mu_v^\pi(b) = 1 - \mu_v^\pi(1-b)$ and the upper bound of $\mu_v^\pi(1-b)$. $\qquad \square$

Now we show Line 3 of $\mathtt{RejectionSampling}(\sigma, S)$ is efficient if $\Phi_i$ is small and $\sigma$ satisfies Assumption 5.2.

**Lemma 5.4.** *Assume* $(\Phi, k, \alpha, n, \xi, \eta, D)$ *is good and* $\sigma$ *satisfies* Assumption 5.2. *If* $\mathbf{e}2^{-k'} \cdot kD \leq 1$, *then for* $\Phi_i = (\mathcal{V}_i, \mathcal{C}_i)$ *from* Line 1 *of* $\mathtt{RejectionSampling}(\sigma, S)$ *we have*

$$\Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} \left[ \pi(\mathcal{V}_i) \text{ is a solution of } \Phi_i \right] \geq \exp\left\{ -\min\left\{ \frac{\xi|\mathcal{C}_i|}{k^6(\alpha+1)}, \frac{kn}{2^{4k}} + \frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\} \right\}.$$

*Proof.* Let $\overline{\mathcal{V}} = \left\{ v \in \mathcal{V}_i \,\middle|\, \deg_{\mathcal{C}_i}(v) \geq D \right\} \subseteq \mathsf{HD}(\mathcal{V}_i) \subseteq \mathsf{HD}(\mathcal{V})$. Then $|\overline{\mathcal{V}}| \leq k|\mathcal{C}_i|/D$ and by Proposition 3.8, $|\overline{\mathcal{V}}| \leq |\mathsf{HD}(\mathcal{V})| \leq n/2^{4k}$. Recall that $\overline{\mathcal{V}}_{\text{sep}}$ and $\overline{\mathcal{C}}_{\text{sep}}$ are the outputs of $\mathtt{ConstructSep}(\overline{\mathcal{V}})$. Then by Lemma 4.3 and $2/\eta \leq k$,

$$|\overline{\mathcal{V}}_{\text{sep}}| \leq \frac{2|\overline{\mathcal{V}}|}{\eta} \leq \frac{2}{\eta} \cdot \min\left\{ \frac{k|\mathcal{C}_i|}{D}, \frac{n}{2^{4k}} \right\} \leq \min\left\{ \frac{2k|\mathcal{C}_i|}{\eta D}, \frac{kn}{2^{4k}} \right\}. \tag{5}$$

Let $\mathcal{V}' = \mathcal{V}_i \cap \overline{\mathcal{V}}_{\text{sep}}$ and $\mathcal{C}' = \mathcal{C}_i \cap \overline{\mathcal{C}}_{\text{sep}}$. By Lemma 5.3, $\Phi^\sigma$ is satisfiable, and thus $\Phi_i$ is also satisfiable. Therefore there exists a partial assignment $\widetilde{\pi}$ extending $\sigma$ by fixing values of $\mathcal{V}'$ to $0/1$ and satisfying all clauses in $\mathcal{C}'$. Then

$$\Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} [\Phi_i(\pi(\mathcal{V}_i)) = \mathsf{True}] \geq 2^{-|\mathcal{V}'|} \Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} \left[ \Phi_i(\pi(\mathcal{V}_i)) = \mathsf{True} \,\middle|\, \pi(\mathcal{V}') = \widetilde{\pi}(\mathcal{V}') \right]$$

$$\geq 2^{-|\overline{\mathcal{V}}_{\text{sep}}|} \Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} \left[ \Phi_i(\pi(\mathcal{V}_i)) = \mathsf{True} \,\middle|\, \pi(\mathcal{V}') = \widetilde{\pi}(\mathcal{V}') \right]$$

$$\geq \mathbf{e}^{-|\overline{\mathcal{V}}_{\text{sep}}|} \Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} \left[ \Phi_i(\pi(\mathcal{V}_i)) = \mathsf{True} \,\middle|\, \pi(\mathcal{V}') = \widetilde{\pi}(\mathcal{V}') \right]. \tag{6}$$

By our choice of $\widetilde{\pi}$, clauses in $\mathcal{C}'$ are already satisfied. On the other hand, since $\mathcal{V}' \subseteq \overline{\mathcal{V}}_{\text{sep}} \subseteq \mathcal{V}_{\text{sep}}$ by Fact 4.2, every clause $C \in \mathcal{C}_i \setminus \mathcal{C}'$ that is not satisfied by $\widetilde{\pi}$ falls into one of the following cases:

16

- If $C$ was not originally in $\mathcal{C}_{\mathsf{sep}}$, then it contains at least $k'$ unassigned variables in $\widetilde{\pi}$ by Assumption 5.2 since $\mathcal{V}' \subseteq \mathcal{V}_{\mathsf{sep}}$ and $\mathcal{V}_{\mathsf{sep}} \subseteq \Lambda(\sigma)$.

- Otherwise, $C$ was originally in $\mathcal{C}_{\mathsf{sep}}$. Then in $\sigma$, it contains at least $k - 2$ unassigned variables by Proposition 3.2 and Assumption 5.2. Now in $\widetilde{\pi}$, at most $2\eta k$ variables are in $\overline{\mathcal{V}}_{\mathsf{sep}}$ and thus fixed, which means at least $k - 2 - 2\eta k \geq k'$ variables remain.

In addition, all the remaining variables $\mathcal{V}_i \setminus \mathcal{V}'$ have degree at most $D$.

Let $\Phi'' = (\mathcal{V}'', \mathcal{C}'')$ where $\mathcal{V}'' = \mathcal{V}_i \setminus \mathcal{V}'$ and $\mathcal{C}'' \subseteq \mathcal{C}_i \setminus \mathcal{C}'$. Then $p(\Phi'') \leq 2^{-k'}$ and $\Delta(\Phi'') \leq kD$. Since $\mathbf{e}2^{-k'} \cdot kD \leq 1$, by Theorem 2.1 with $B$ being the event "$\Phi''$ is satisfied" which correlates with all $|\mathcal{C}''| \leq |\mathcal{C}_i|$ clauses, we have

$$\Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} \left[ \Phi_i(\pi(\mathcal{V}_i)) = \mathsf{True} \,\big|\, \pi(\mathcal{V}') = \widetilde{\pi}(\mathcal{V}') \right] \geq (1 - \mathbf{e}2^{-k'})^{|\mathcal{C}''|} \geq \exp\left\{ -\frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\}.$$

Putting (5) and (6) back, we have

$$\Pr_{\pi(\mathcal{V}_i) \sim \{0,1\}^{\mathcal{V}_i}} [\Phi_i(\pi(\mathcal{V}_i)) = \mathsf{True}] \geq \exp\left\{ -\min\left\{ \frac{2k|\mathcal{C}_i|}{\eta D}, \frac{kn}{2^{4k}} \right\} - \frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\}$$

$$\geq \exp\left\{ -\min\left\{ \frac{2k|\mathcal{C}_i|}{\eta D} + \frac{|\mathcal{C}_i|}{kD}, \frac{kn}{2^{4k}} + \frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\} \right\}$$

$$\text{(since } \mathbf{e}2^{-k'} \cdot kD \leq 1\text{)}$$

$$\geq \exp\left\{ -\min\left\{ \frac{k^2|\mathcal{C}_i|}{D}, \frac{kn}{2^{4k}} + \frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\} \right\} \qquad \text{(since } 3/k \leq \eta \leq 1\text{)}$$

$$= \exp\left\{ -\min\left\{ \frac{\xi|\mathcal{C}_i|}{k^6(\alpha+1)}, \frac{kn}{2^{4k}} + \frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\} \right\},$$

$$\text{(since } D = k^8(\alpha+1)/\xi\text{)}$$

as desired. $\qquad\qquad \square$

**Corollary 5.5.** *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good and $\sigma$ satisfies Assumption 5.2. If $\mathbf{e}2^{-k'} \cdot kD \leq 1$, then* RejectionSampling$(\sigma, S)$ *runs in expected time*

$$\widetilde{O}\left( \sum_i |\mathcal{V}_i| \cdot \exp\left\{ \min\left\{ \frac{\xi|\mathcal{C}_i|}{k^6(\alpha+1)}, \frac{kn}{2^{4k}} + \frac{\mathbf{e}|\mathcal{C}_i|}{2^{k'}} \right\} \right\} \right),$$

*where each $\Phi_i = (\mathcal{V}_i, \mathcal{C}_i)$ is from Line 1 of* RejectionSampling$(\sigma, S)$.

## 5.1 Algorithms for the Atypical Setting

To give a sense of the bound in Corollary 5.5, we use it to analyze the atypical setting of Theorem 1.3 where either $\varepsilon$ or $\alpha$ is too small. Indeed, in these cases the naive rejection sampling algorithm is already highly efficient.

**Lemma 5.6** (Small Error Setting). *Assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is good and $\varepsilon \leq \exp\left\{ -n/2^{k/2} \right\}$. Then* RejectionSampling$(\star^{\mathcal{V}}, \mathcal{V})$ *runs in expected time $\widetilde{O}\left( (1/\varepsilon)^{\xi/k} \right)$ and has output distribution exactly $\mu$.*

*Proof.* By Fact 5.1, we only need to bound the expected runtime. By Proposition 3.2 and Lines 2 and 3 of `ConstructSep(V)`, we set $k' = (1 - 2\eta)k - 2$ in Corollary 5.5. Since $\eta = 15\log(k)/k$, $D = k^8(\alpha + 1)/\xi$, and $\alpha \leq \xi \cdot 2^{k/3}/k^{50}$ with $\xi \geq 2^{-k/8}$ and $k \geq 2^{20}$, we have

$$\mathbf{e}2^{-k'} \cdot kD = 4\mathbf{e}2^{-k} \cdot k^{39}(\alpha + 1)/\xi \leq 8\mathbf{e}k^{-11} \cdot 2^{-2k/3} \leq 1.$$

Then by Corollary 5.5, the expected runtime is upper bounded by

$$\widetilde{O}\left(\sum_i |\mathcal{V}_i| \cdot \exp\left\{\frac{kn}{2^{4k}} + \frac{4\mathbf{e}k^{30}|\mathcal{C}_i|}{2^k}\right\}\right) \leq \widetilde{O}\left(\sum_i |\mathcal{V}_i| \cdot \exp\left\{\frac{kn}{2^{4k}} + \frac{4\mathbf{e}n}{2^{2k/3}}\right\}\right)$$

$$\text{(since } |\mathcal{C}_i| \leq |\mathcal{C}| = \alpha n \leq n \cdot 2^{k/3}/k^{30})$$

$$\leq \widetilde{O}\left(\sum_i |\mathcal{V}_i| \cdot \exp\left\{\frac{n}{2k \cdot 2^{5k/8}}\right\}\right) \qquad \text{(since } k \geq 2^{20})$$

$$= \widetilde{O}\left(n \cdot \exp\left\{\frac{n}{2k \cdot 2^{5k/8}}\right\}\right) \leq \widetilde{O}\left(\exp\left\{\frac{n}{k \cdot 2^{5k/8}}\right\}\right)$$

$$\text{(since } n \geq 2^{\Omega(k)})$$

$$\leq \widetilde{O}\left((1/\varepsilon)^{\xi/k}\right) \quad \text{(since } \varepsilon \leq \exp\left\{-n/2^{k/2}\right\} \text{ and } \xi \geq 2^{-k/8})$$

as desired. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\square$

**Lemma 5.7** (Small Density Setting). *Assume* $(\Phi, k, \alpha, n, \xi, \eta, D)$ *is good and* $\alpha \leq 1/k^3$. *Then* `RejectionSampling($\star^{\mathcal{V}}, \mathcal{V}$)` *runs in expected time* $\widetilde{O}\left(n^{1+\xi/k}\right)$ *and has output distribution exactly* $\mu$.

*Proof.* Similar analysis as in the proof of Lemma 5.6. In addition, by Proposition 3.5 with $\ell = \ln n$, we have $\alpha^2 n^4 (\mathbf{e}k^2\alpha)^\ell < 1$ and thus the maximal connected component in $G_\Phi$ has size at most $\ln n$, i.e., each $\mathcal{C}_i$ in `RejectionSampling($\star^{\mathcal{V}}, \mathcal{V}$)` has size at most $\ln n$. Then by Corollary 5.5, the expected runtime is upper bounded by

$$\widetilde{O}\left(\sum_i |\mathcal{V}_i| \cdot \exp\left\{\frac{\xi\ln n}{k^6(\alpha + 1)}\right\}\right) \leq \widetilde{O}\left(n \cdot \exp\left\{\frac{\xi\ln n}{k}\right\}\right) = \widetilde{O}\left(n^{1+\xi/k}\right) \qquad\qquad \square$$

# 6 Algorithms for the Typical Setting

In this section, we present the sampling algorithm for the typical setting: $\alpha \geq 1/k^3$ and $\varepsilon \geq \exp\left\{-n/2^{k/2}\right\}$. We will conveniently assume our instance is nice (in particular, $\alpha \geq 1/k^3$), though some of the results also hold with weaker assumptions. From now on, unless specifically mentioned, we assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is nice and save the space of repeatedly putting this in the statements.

Our main algorithm is a modification of the ones in [HWY22]. Hence some of our notation and definitions will be similar to theirs, which we hope is easier to understand if the reader is already familiar with [HWY22].

Given a partial assignment $\sigma$ and $\mathcal{V}_{\mathsf{sep}}$ constructed above, we define $\mathcal{V}_{\mathsf{alive}}^\sigma$: For each $v \in \mathcal{V}$, $v \in \mathcal{V}_{\mathsf{alive}}^\sigma$ iff (i) $\sigma(v) = \star$ and $v \notin \mathcal{V}_{\mathsf{sep}}$, and (ii) for every clause $C \in \mathcal{C} \setminus \mathcal{C}_{\mathsf{sep}}$, either $C(\sigma) = \mathsf{True}$ or $|\mathsf{vbl}(C) \cap \Lambda(\sigma) \setminus (\mathcal{V}_{\mathsf{sep}} \cup \{v\})| \geq (2/3 - 2\eta)k$. Intuitively, $v \in \mathcal{V}_{\mathsf{alive}}^\sigma$ means after fixing $v$, each unsatisfied clause will still contain many unassigned variables, consistent with Assumption 5.2.

Now we present our `SolutionSampling($\Phi$)` algorithm in Algorithm 3 similar to [HWY22, Algorithm 4].

By dynamically maintaining and updating the size of each $\mathsf{vbl}(C) \cap \Lambda(\sigma) \setminus \mathcal{V}_{\mathsf{sep}}$, the total runtime for checking whether $v_i \in \mathcal{V}_{\mathsf{alive}}^\sigma$ is very efficient.

**Algorithm 3:** The `SolutionSampling` Algorithm for the Typical Setting

---

**Input:** A random $k$-CNF formula $\Phi = (\mathcal{V}, \mathcal{C})$

**Output:** A random assignment $\sigma$ distributed as $\mu$

**1** Obtain $\mathcal{V}_{\sf sep}, \mathcal{C}_{\sf sep} \leftarrow \texttt{ConstructSep}(\mathcal{V})$

**2** Initialize $\sigma \leftarrow \star^{\mathcal{V}}$

**3 foreach** $i = 1$ **to** $n$ **do**

**4** $\quad$ **if** $v_i \in \mathcal{V}^{\sigma}_{\sf alive}$ **then** Update $\sigma(v_i) \leftarrow \texttt{MarginSample}(\sigma, v_i)$

$\quad$ **end**

**5** $\sigma \leftarrow \texttt{RejectionSampling}(\sigma, \Lambda(\sigma))$

**6 return** $\sigma$

---

**Fact 6.1.** *The runtime of all the checking $v_i \in \mathcal{V}^{\sigma}_{\sf alive}$ combined is $\widetilde{O}(n)$.*

Assumption 5.2 will be preserved with $k' = (2/3 - 2\eta)k$ throughout if we only update the alive variables. This is indeed the case in Algorithm 3 and we highlight it as the following formal statements. For future referencing, we explicitly write Assumption 5.2 with $k' = (2/3 - 2\eta)k$ as Assumption 6.2.

**Assumption 6.2.** $\mathcal{V}_{\sf sep} \subseteq \Lambda(\sigma)$ and for every clause $C \in \mathcal{C} \setminus \mathcal{C}_{\sf sep}$, either $C(\sigma) = {\sf True}$ or $|{\sf vbl}(C) \cap \Lambda(\sigma) \setminus \mathcal{V}_{\sf sep}| \geq (2/3 - 2\eta)k$.

**Fact 6.3.** *Assume we construct a partial assignment $\sigma$ by starting with $\sigma = \star^{\mathcal{V}}$ and repeatedly fixing variables in $\mathcal{V}^{\sigma}_{\sf alive}$. Then $\sigma$ satisfies Assumption 6.2.*

*Proof.* We prove by induction. The base case $\sigma = \star^{\mathcal{V}}$ trivially holds since for any $C \in \mathcal{C} \setminus \mathcal{C}_{\sf sep}$, we have $|{\sf vbl}(C)| \geq k - 2$ by Proposition 3.2 and $|{\sf vbl}(C) \cap \mathcal{V}_{\sf sep}| \leq 2\eta k$ by Line 3 of $\texttt{ConstructSep}(\mathcal{V})$.

For the inductive case, assume we fix $v \in \mathcal{V}^{\sigma}_{\sf alive}$ and obtain $\sigma'$. Then by the definition of $\mathcal{V}^{\sigma}_{\sf alive}$, we know $v \notin \mathcal{V}_{\sf sep}$, which means $\Lambda(\sigma') \supseteq \Lambda(\sigma) \setminus \{v\} \supseteq \mathcal{V}_{\sf sep}$ by induction hypothesis. On the other hand, for every clause $C \in \mathcal{C} \setminus \mathcal{C}_{\sf sep}$,

- if $C(\sigma) = {\sf True}$, then $C(\sigma') = {\sf True}$,

- otherwise, $|{\sf vbl}(C) \cap \Lambda(\sigma') \setminus \mathcal{V}_{\sf sep}| \geq |{\sf vbl}(C) \cap \Lambda(\sigma) \setminus (\mathcal{V}_{\sf sep} \cup \{v\})| \geq (2/3 - 2\eta)k$ since $v \in \mathcal{V}^{\sigma}_{\sf alive}$. $\qquad\square$

As a corollary of Lemma 5.3, we have good control for the marginal of every remaining variable outside $\mathcal{V}_{\sf sep}$, and in particular, any $v \in \mathcal{V}^{\sigma}_{\sf alive}$.

**Corollary 6.4** (Local Uniformity). *Assume $\sigma$ satisfies Assumption 6.2. Then $\Phi^{\sigma}$ is satisfiable. Moreover, for each $b \in \{0, 1\}$ and $v \in \Lambda(\sigma) \setminus \mathcal{V}_{\sf sep}$, we have*

$$\frac{1 - \delta}{2} \leq \mu^{\sigma}_v(b) \leq \frac{1 + \delta}{2},$$

*where $\delta = \xi/(k^{40}\alpha)$.*

*Proof.* Let $k' = (2/3 - 2\eta)k$. Since $\eta = 15\log(k)/k$, $D = k^8(\alpha + 1)/\xi$, and $\alpha \leq \xi \cdot 2^{k/3}/k^{50}$ with $\xi \geq 2^{-k/8}$ and $k \geq 2^{20}$, we have

$$\mathbf{e}2^{-k'} \cdot kD = \mathbf{e}2^{-2k/3} \cdot k^{39}(\alpha + 1)/\xi \leq 2\mathbf{e}k^{-11} \leq 1.$$

Then by Lemma 5.3, we know $\Phi^{\sigma}$ is satisfiable and

$$\left| \mu^{\sigma}_v(b) - \frac{1}{2} \right| \leq \frac{1}{2} \cdot \mathbf{e}2^{-k'}D = \frac{1}{2} \cdot \mathbf{e}2^{-2k/3} \cdot k^{38}(\alpha + 1)/\xi$$

19

$$\leq \frac{1}{2} \cdot 2\mathbf{e}2^{-2k/3} \cdot k^{41} \cdot \alpha/\xi. \qquad \qquad \text{(since } \alpha \geq 1/k^3)$$

Since $\alpha \leq \xi \cdot 2^{k/3}/k^{50}$, we have

$$\frac{2\mathbf{e}2^{-2k/3} \cdot k^{41} \cdot \alpha/\xi}{\delta} = \frac{2\mathbf{e}k^{81} \cdot \alpha^2}{2^{2k/3} \cdot \xi^2} \leq \frac{2\mathbf{e}}{k^{19}} \leq 1$$

and thus $|\mu_v^\sigma(b) - 1/2| \leq \delta/2$. $\qquad \qquad \square$

Similarly, we have the following efficiency bound for the rejection sampling after replacing Assumption 5.2 with Assumption 6.2 in Corollary 5.5.

**Corollary 6.5.** *Assume $\sigma$ satisfies Assumption 6.2. Then* `RejectionSampling`$(\sigma, S)$ *runs in expected time*

$$\widetilde{O}\left(\sum_i |\mathcal{V}_i| \cdot \exp\left\{\frac{\xi|\mathcal{C}_i|}{k^6(\alpha+1)}\right\}\right),$$

*where each $\Phi_i = (\mathcal{V}_i, \mathcal{C}_i)$ is from Line 1 of* `RejectionSampling`$(\sigma, S)$.

For convenience, we will reserve $\delta = \xi/(k^{40}\alpha)$ as the local uniformity parameter from now on.

To obtain the correct marginal distribution for each Line 4, `MarginSample`$(\sigma, v)$ should sample from $\mu_v^\sigma$. By Corollary 6.4, this distribution is $\delta$-close to an unbiased coin. This inspires us to define the following distribution $\tau$ as a "lower bound" for any $\mu_v^\sigma$ that $v \in \mathcal{V}_{\mathsf{alive}}^\sigma$:

$$\tau = \begin{cases} 0 & \text{w.p} \quad (1-\delta)/2, \\ 1 & \text{w.p} \quad (1-\delta)/2, \\ \star & \text{w.p} \quad \delta. \end{cases}$$

As described in Algorithm 4, `MarginSample`$(\sigma, v)$ will first naively sample from $\tau$, and resample using `MarginOverflow`$(\sigma, v)$ if obtained $\star$ from $\tau$.

---

**Algorithm 4:** The `MarginSample` Algorithm

---
**Input:** $\sigma \in \{0, 1, \star, \star\}^\mathcal{V}$ and $v \in \mathcal{V}_{\mathsf{alive}}^\sigma$
**Output:** A binary random variable distributed as $\mu_v^\sigma$
**1** Sample $\sigma(v) \sim \tau$
**2** if $\sigma(v) = \star$ then return `MarginOverflow`$(\sigma, v)$
**3** else return $\sigma(v)$

---

Naturally, `MarginOverflow`$(\sigma, v)$ should complete $\tau$ into $\mu_v^\sigma$. Thus it should output a binary bit distributed proportional to $\mu_v^\sigma - \tau$, which we define as $\nu_v^\sigma$: For each $b \in \{0, 1\}$, we set

$$\nu_v^\sigma(b) = \frac{\mu_v^\sigma(b) - \tau(b)}{\tau(\star)} = \frac{\mu_v^\sigma(b) - (1-\delta)/2}{\delta}.$$

On the other hand, there exists a standard toolbox [NP05, Hub16, DHKN21], called *Bernoulli factory*, to provide samples from $\nu_v^\sigma$, which is a linear function of $\mu_v^\sigma$, using samples from $\mu_v^\sigma$. Here we use the statement in [HWY22]:

**Lemma 6.6** ([HWY22, Appendix A]). *There exists a Las Vegas algorithm* `BernoulliFactory()` *such that the following holds: Assume* $b_1, b_2, \ldots$ *are independent samples from* $\mu_v^\sigma$, *where* $\mu_v^\sigma$ *is unknown to the algorithm and* $\mu_v^\sigma(b) \geq \tau(b)$ *for* $b \in \{0, 1\}$. *Then* `BernoulliFactory`$(b_1, b_2, \ldots)$ *runs in expected time* $\widetilde{O}(1/\delta^2) = \widetilde{O}(1/\xi^2)$ *and has output distribution exactly* $\nu_v^\sigma$.

Samples from $\mu_v^\sigma$ can be provided by executing `RejectionSampling`$(\sigma, v)$, but simply doing so is just self-referencing: Why not let `MarginSample`$(\sigma, v)$ be `RejectionSampling`$(\sigma, v)$ in the first place?

The trick here is to recursively fix more variables in $\sigma$ and postpone the Bernoulli factory to the end. Hopefully at that point, most variables in $\sigma$ are fixed and $\Phi^\sigma$ can be factorized into small components, which makes the rejection sampling efficient. This will become rigorous as we describe `MarginOverflow`$(\sigma, v)$ shortly.

## 6.1 The Margin Overflow Algorithm and Truncation

To describe and analyze `MarginOverflow()`, we need the following notation to make rigorous our recursive sampling order: Let $\sigma$ be a partial assignment. For each $C \in \mathcal{C}$,

- $\mathcal{C}_\star^\sigma$: $C \in \mathcal{C}_\star^\sigma$ iff there exists some $v \in \mathsf{vbl}(C)$ that $\sigma(v) = \star$.
- $\mathcal{C}_{\mathsf{frozen}}^\sigma$: $C \in \mathcal{C}_{\mathsf{frozen}}^\sigma$ iff (i) $C(\sigma) \neq \mathsf{True}$ and $C \notin \mathcal{C}_{\mathsf{sep}}$, and (ii) $|\mathsf{vbl}(C) \cap \Lambda(\sigma) \setminus \mathcal{V}_{\mathsf{sep}}| < 1 + (2/3 - 2\eta)k$.
- $\mathcal{C}_{\mathsf{bad}}^\sigma$: $C \in \mathcal{C}_{\mathsf{bad}}^\sigma$ iff (i) $C(\sigma) \neq \mathsf{True}$ and $C \notin \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{sep}}$, and (ii) for any $v \in \mathsf{vbl}(C) \setminus \mathcal{V}_{\mathsf{sep}}$ with $\sigma(v) = \star$, there exists some $C' \in \mathcal{C}_{\mathsf{frozen}}^\sigma$ such that $v \in \mathsf{vbl}(C')$.

We remark that, though $\mathcal{C}_{\mathsf{frozen}}^\sigma \cap \mathcal{C}_{\mathsf{bad}}^\sigma = \emptyset$, it is possible that $\mathcal{C}_\star^\sigma \cap \mathcal{C}_{\mathsf{frozen}}^\sigma \neq \emptyset$ and $\mathcal{C}_\star^\sigma \cap \mathcal{C}_{\mathsf{bad}}^\sigma \neq \emptyset$. The definition of $\mathcal{C}_{\mathsf{frozen}}^\sigma$ is a direct opposite of $\mathcal{V}_{\mathsf{alive}}^\sigma$, while $\mathcal{C}_{\mathsf{bad}}^\sigma$ intermediately violates $\mathcal{V}_{\mathsf{alive}}^\sigma$ due to $\mathcal{C}_{\mathsf{frozen}}^\sigma$. This is formalized in the following fact.

**Fact 6.7.** *For any* $C \in \mathcal{C}_{\mathsf{sep}} \cup \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{bad}}^\sigma$, *we have* $\mathsf{vbl}(C) \cap \mathcal{V}_{\mathsf{alive}}^\sigma = \emptyset$. *As a consequence, we have* $\mathcal{C}_\star^\sigma \subseteq \mathcal{C}_\star^{\sigma'}$, $\mathcal{C}_{\mathsf{frozen}}^\sigma \subseteq \mathcal{C}_{\mathsf{frozen}}^{\sigma'}$, *and* $\mathcal{C}_{\mathsf{bad}}^\sigma \subseteq \mathcal{C}_{\mathsf{bad}}^{\sigma'}$ *if* $\sigma'$ *extends* $\sigma$ *by fixing some variable in* $\mathcal{V}_{\mathsf{alive}}^\sigma$.

*Proof.* Recall that $v \in \mathcal{V}_{\mathsf{alive}}^\sigma$ iff (a) $\sigma(v) = \star$, (b) $v \notin \mathcal{V}_{\mathsf{sep}}$, and (c) for every clause $C \in \mathcal{C} \setminus \mathcal{C}_{\mathsf{sep}}$ and $C(\sigma) \neq \mathsf{True}$, $|\mathsf{vbl}(C) \cap \Lambda(\sigma) \setminus (\mathcal{V}_{\mathsf{sep}} \cup \{v\})| \geq (2/3 - 2\eta)k$.

Now assume $v \in \mathsf{vbl}(C) \cap \mathcal{V}_{\mathsf{alive}}^\sigma$.

- If $C \in \mathcal{C}_{\mathsf{sep}}$, then by the definition of $\mathcal{V}_{\mathsf{sep}}$, we have $v \in \mathcal{V}_{\mathsf{sep}}$ and contradict to (b).
- If $C \in \mathcal{C}_{\mathsf{frozen}}^\sigma$, then by the definition of $\mathcal{C}_{\mathsf{frozen}}^\sigma$, we have $|\mathsf{vbl}(C) \cap \Lambda(\sigma) \setminus (\mathcal{V}_{\mathsf{sep}} \cup \{v\})| < (2/3 - 2\eta)k$ and contradict to (c).
- If $C \in \mathcal{C}_{\mathsf{bad}}^\sigma$, then by (a), we know $\sigma(v) = \star$. Then by the definition of $\mathcal{C}_{\mathsf{bad}}^\sigma$, we have $v \in \mathsf{vbl}(C')$ for some $C' \in \mathcal{C}_{\mathsf{frozen}}^\sigma$ and contradict to the last item with $C$ replaced by $C'$. $\square$

To preserve Assumption 6.2 and by Fact 6.3, we can only afford to sample variables in $\mathcal{V}_{\mathsf{alive}}^\sigma$. By Fact 6.7, this means we need to avoid $\mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{bad}}^\sigma \cup \mathcal{C}_{\mathsf{sep}}$.

On the other hand, the marginal distribution of a variable depends on all the variables and clauses connected to it. This motivates us to define, for each $v$ with $\sigma(v) = \star$, the bad interior $\mathcal{C}_{\mathsf{int}}^\sigma(v) \subseteq \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{bad}}^\sigma \cup \mathcal{C}_{\mathsf{sep}}$, which contains variables connected to $v$ (but unfortunately none of them is alive). Formally, we put $C \in \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{bad}}^\sigma \cup \mathcal{C}_{\mathsf{sep}}$ into $\mathcal{C}_{\mathsf{int}}^\sigma(v)$ iff either $v \in \mathsf{vbl}(C)$ or there exists some $C' \in \mathcal{C}_{\mathsf{int}}^\sigma(v)$ that $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\sigma) \neq \emptyset$.

21

To obtain alive variables to sample from, we need to take one step further to form the current component $\mathcal{C}_{\mathsf{con}}^\sigma(v)$. Formally, we put $C \in \mathcal{C}$ into $\mathcal{C}_{\mathsf{con}}^\sigma(v)$ iff $C \in \mathcal{C}_{\mathsf{int}}^\sigma(v)$, or $v \in \mathsf{vbl}(C)$, or there exists some $C' \in \mathcal{C}_{\mathsf{int}}^\sigma(v)$ that $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\sigma) \neq \emptyset$.

Then we take the union of the current components of all the $\star$'s, since we care about the marginals of these variables. Define $\mathcal{C}_{\mathsf{con}}^\sigma$ to be the union of $\mathcal{C}_{\mathsf{con}}^\sigma(v)$ for all $v$ with $\sigma(v) = \star$. Let $\mathcal{V}_{\mathsf{con}}^\sigma = \bigcup_{C \in \mathcal{C}_{\mathsf{con}}^\sigma} \mathsf{vbl}(C)$. Then we define

$$\mathsf{NextVar}(\sigma) = \begin{cases} v_i \in \mathcal{V}_{\mathsf{alive}}^\sigma \cap \mathcal{V}_{\mathsf{con}}^\sigma \text{ with smallest } i & \text{if } \mathcal{V}_{\mathsf{alive}}^\sigma \cap \mathcal{V}_{\mathsf{con}}^\sigma \neq \emptyset, \\ \bot & \text{otherwise,} \end{cases}$$

which will be the function for selecting the next variable to perform marginal sampling.

Now we give the pseudo-code of $\mathtt{MarginOverflow}(\sigma, v)$. We remark that Lines 3 and 4 is equivalent to calling $\mathtt{MarginSample}(\sigma, u)$. To avoid confusion from subroutines calling each other, we choose to expand it out as the current presentation.

---

**Algorithm 5:** The $\mathtt{MarginOverflow}$ Algorithm

---

**Input:** $\sigma \in \{0, 1, \star, \star\}^{\mathcal{V}}$ and $v \in \mathcal{V}$ with $\sigma(v) = \star$ and $v \in \mathcal{V}_{\mathsf{alive}}^{\overline\sigma}$ where $\overline\sigma$ equals $\sigma$ except $\overline\sigma(v) = \star$

**Output:** A binary random variable distributed as $\nu_v^\sigma$

**1** Let $u \leftarrow \mathsf{NextVar}(\sigma)$

**2 if** $u \neq \bot$ **then**

**3**      Sample $\sigma(u) \sim \tau$

**4**      **if** $\sigma(u) = \star$ **then** Update $\sigma(u) \leftarrow \mathtt{MarginOverflow}(\sigma, u)$

**5**      **return** $\mathtt{MarginOverflow}(\sigma, v)$

    **else**

**6**      **return** $\mathtt{BernoulliFactory}(b_1, b_2, \ldots)$ where $b_1, b_2, \ldots$ are independent samples provided by executing $\mathtt{RejectionSampling}(\sigma, v)$

    **end**

---

By pre-processing the maximal connected components in $\mathcal{C}_{\mathsf{sep}}$, we can dynamically maintain $\mathcal{C}_\star^\sigma, \mathcal{C}_{\mathsf{frozen}}^\sigma, \mathcal{C}_{\mathsf{bad}}^\sigma, \mathcal{V}_{\mathsf{alive}}^\sigma$ and their connectivity relation with $\mathcal{C}_{\mathsf{sep}}$. Then we can dynamically update $\mathcal{C}_{\mathsf{con}}^\sigma$ and $\mathcal{V}_{\mathsf{con}}^\sigma$ by joining the new connected components.

Therefore each computation of $\mathsf{NextVar}()$ can be done in worst case time $\mathsf{poly}(k, d)$, where $d = d(\Phi)$ is the maximum variable degree of $\Phi$. By Proposition 3.7, we obtain the following bound.

**Fact 6.8.** *With $\widetilde{O}(n)$ pre-processing time, the runtime of each $\mathsf{NextVar}()$ is $\widetilde{O}(1)$.*

Note that whenever $u$ from Line 1 is not $\bot$, we have $u \in \mathcal{V}_{\mathsf{alive}}^\sigma$. Therefore by Fact 6.3, Assumption 6.2 is preserved throughout the algorithm. This provides us the following correctness guarantee, the proof of which is almost identical to the inductive proof of [HWY22, Theorem 5.5].

**Lemma 6.9.** *Assume $\sigma$ satisfies Assumption 6.2. Then $\mathtt{MarginOverflow}(\sigma, v)$ terminates almost surely and has output distribution exactly $\nu_v^\sigma$.*

*Proof.* Observe that each deeper recursion will have the value of $u$ changed from $\star$ to $0/1/\star$. Therefore the number of $\star$'s in $\sigma$ is decreasing and thus the recursion ends eventually.

Now we prove the statement by induction on $\sigma$. The base case corresponds to the leaf of the recursion, where $u = \bot$ from Line 1. By Corollary 6.4, we know $\Phi^\sigma$ is satisfiable. Thus by Fact 5.1,

`RejectionSampling(σ, v)` terminates almost surely and has output distribution exactly $\mu_v^\sigma$. Now by Corollary 6.4, $\mu_v^\sigma$ is lower bounded by $\tau$. Therefore by Lemma 6.6, `BernoulliFactory(b₁, b₂, ...)` has output distribution exactly $\nu_v^\sigma$. This proves the base case.

For the inductive case that $u \neq \perp$, let $\sigma_0, \sigma_1, \sigma_\star$ equal $\sigma$ except $\sigma_0(u) = 0/1/\star$ respectively. Then by induction hypothesis, `MarginOverflow(σ_⋆, u)` returns a bit distributed as

$$\nu_{\sigma_\star}^u(b) = \frac{\mu_\sigma^u(b) - \tau(b)}{\tau(\star)} \quad \text{for } b \in \{0, 1\}.$$

Let $\sigma'$ be the updated $\sigma$ upon reaching Line 5. Thus $\sigma'$ equals $\sigma$ except

$$\mathbf{Pr}[\sigma'(u) = b] = \tau(b) + \tau(\star) \cdot \nu_{\sigma_\star}^u(b) = \mu_\sigma^u(b) \quad \text{for } b \in \{0, 1\}. \tag{7}$$

By induction hypothesis again, Line 5 terminates almost surely and obtains distribution $\nu_{\sigma'}^v$ where

$$
\begin{aligned}
\nu_{\sigma'}^v(b) &= \mathbf{Pr}[\sigma'(u) = 0] \cdot \nu_{\sigma_0}^v(b) + \mathbf{Pr}[\sigma'(u) = 1] \cdot \nu_{\sigma_1}^v(b) \\
&= \mu_\sigma^u(0) \cdot \nu_{\sigma_0}^v(b) + \mu_\sigma^u(1) \cdot \nu_{\sigma_1}^v(b) && \text{(by (7))} \\
&= \frac{1}{\tau(\star)} \cdot \left( \mu_\sigma^u(0) \cdot \mu_{\sigma_0}^v(b) + \mu_\sigma^u(1) \cdot \mu_{\sigma_1}^v(b) - \tau(b) \right) && \text{(by the definition of } \nu\text{)} \\
&= \frac{1}{\tau(\star)} \cdot \left( \mu_\sigma^v(b) - \tau(b) \right) = \nu_\sigma^v(b)
\end{aligned}
$$

for $b \in \{0, 1\}$ as desired. □

As an immediate corollary, we obtain the correctness of `MarginSample(σ, v)`.

**Corollary 6.10.** *Assume $\sigma$ satisfies Assumption 6.2. Then* `MarginSample(σ, v)` *terminates almost surely and has output distribution exactly $\mu_v^\sigma$.*

Then by the chain rule of conditional probability, Fact 6.3, and Fact 5.1, we obtain the correctness of our main algorithm.

**Corollary 6.11.** `SolutionSampling(Φ)` *terminates almost surely and has output distribution exactly $\mu$.*

Ideally, we only need to bound the expected runtime of each `MarginOverflow(σ, v)` and the final rejection sampling in `SolutionSampling(Φ)`; then we obtain the runtime of the whole algorithm. This will actually be a *perfect* sampler that outputs an uniform solution exactly, and is indeed the case for the standard $k$-CNFs in the local lemma regime [HWY22]. But the issue here is that $\Phi$ is random, and its structural properties break down when we analyze components of large size. Therefore to ensure that we have good structural properties at hand, we will have to halt when the component goes beyond a certain size. Fortunately, we are able to show that this truncation happens with small probability, and thus only incur small deviation in the total variation distance.

To give some intuition about the truncation, we analyze the efficiency of the leaf recursion of `MarginOverflow(σ, v)`. Recall our definition of $\mathcal{V}_{\mathsf{con}}^\sigma$ and $\mathcal{C}_{\mathsf{con}}^\sigma$ at the beginning of this subsection.

**Lemma 6.12.** *Assume $\sigma$ satisfies Assumption 6.2. If $\mathsf{NextVar}(\sigma) = \perp$, then* `MarginOverflow(σ, v)` *runs in expected runtime*

$$\widetilde{O}\left( \frac{|\mathcal{C}_{\mathsf{con}}^\sigma|}{\xi^2} \cdot \exp\left\{ \frac{\xi \cdot |\mathcal{C}_{\mathsf{con}}^\sigma|}{k^6(\alpha + 1)} \right\} \right).$$

*Proof.* Let $\Phi' = (\mathcal{V}', \mathcal{C}')$ be the maximal connected component in $\Phi^\sigma$ intersecting $v$. Then by Corollary 6.5, the expected runtime of $\texttt{RejectionSampling}(\sigma, v)$ is bounded by

$$\widetilde{O}\left(|\mathcal{V}'| \cdot \exp\left\{\frac{\xi \cdot |\mathcal{C}'|}{k^6(\alpha+1)}\right\}\right) = \widetilde{O}\left(|\mathcal{C}'| \cdot \exp\left\{\frac{\xi \cdot |\mathcal{C}'|}{k^6(\alpha+1)}\right\}\right),$$

where we use the fact that $|\mathcal{V}'| \leq k \cdot |\mathcal{C}'| = \widetilde{O}(|\mathcal{C}'|)$. Thus by Lemma 6.6 and the analysis of Lemma 6.9, the expected runtime of $\texttt{MarginOverflow}(\sigma, v)$ is

$$\widetilde{O}\left(\frac{|\mathcal{C}'|}{\xi^2} \cdot \exp\left\{\frac{\xi \cdot |\mathcal{C}'|}{k^6(\alpha+1)}\right\}\right).$$

Now it suffices to show $\mathcal{C}' \subseteq \mathcal{C}_{\mathsf{con}}^\sigma$.

Note that $\mathcal{C}'$ can be constructed as follows: Starting with $\mathcal{C}' = \emptyset$, we repeatedly put $C \in \mathcal{C}$ into $\mathcal{C}'$ if $\mathcal{C}(\sigma) \neq \mathsf{True}$ and, either (1) $v \in \mathsf{vbl}(C)$ or (2) $\mathsf{vbl}(C) \cap \mathsf{vbl}(C') \cap \Lambda(\sigma) \neq \emptyset$ for some $C' \in \mathcal{C}'$. Assume towards contradiction that $C$ is the first clause included in $\mathcal{C}'$ but not in $\mathcal{C}_{\mathsf{con}}^\sigma$.

- If $C$ satisfies condition (1), we know $C \in \mathcal{C}_{\mathsf{con}}^\sigma(v) \subseteq \mathcal{C}_{\mathsf{con}}^\sigma$ since $\sigma(v) = \bigstar$. A contradiction.
- Otherwise, $C$ satisfies condition (2). Since $C'$ is included in both $\mathcal{C}'$ and $\mathcal{C}_{\mathsf{con}}^\sigma$, there exists some $v'$ such that $\sigma(v') = \bigstar$ and $C' \in \mathcal{C}_{\mathsf{con}}^\sigma(v') \subseteq \mathcal{C}_{\mathsf{con}}^\sigma$. Then we have the following cases:
  - If $C' \in \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{bad}}^\sigma \cup \mathcal{C}_{\mathsf{sep}}$, then $C' \in \mathcal{C}_{\mathsf{int}}^\sigma(v')$. Thus $C \in \mathcal{C}_{\mathsf{con}}^\sigma(v') \subseteq \mathcal{C}_{\mathsf{con}}^\sigma$. A contradiction.
  - If $C' \in \mathcal{C}_\bigstar^\sigma$, then there exists some $v'' \in \mathsf{vbl}(C')$ such that $\sigma(v'') = \bigstar$. Then $C \in \mathcal{C}_{\mathsf{con}}^\sigma(v'') \subseteq \mathcal{C}_{\mathsf{con}}^\sigma$. A contradiction.
  - Otherwise, $C' \notin \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{bad}}^\sigma \cup \mathcal{C}_{\mathsf{sep}} \cup \mathcal{C}_\bigstar^\sigma$. Since $\mathsf{NextVar}(\sigma) = \bot$ and $C' \in \mathcal{C}_{\mathsf{con}}^\sigma$, we have $\mathsf{vbl}(C') \cap \mathcal{V}_{\mathsf{alive}}^\sigma = \emptyset$. Note that $\mathsf{vbl}(C')$ has no $\bigstar$. Thus for any $u \in \mathsf{vbl}(C') \cap \Lambda(\sigma) \setminus \mathcal{V}_{\mathsf{sep}}$, there exists some $C'' \in \mathcal{C} \setminus \mathcal{C}_{\mathsf{sep}}$ such that $C''(\sigma) \neq \mathsf{True}$ and

    $$|\mathsf{vbl}(C'') \cap \Lambda(\sigma) \setminus (\mathcal{V}_{\mathsf{sep}} \cup \{u\})| < (2/3 - 2\eta)k.$$

    These $C''$'s satisfy $|\mathsf{vbl}(C'') \cap \Lambda(\sigma) \setminus \mathcal{V}_{\mathsf{sep}}| < 1 + (2/3 - 2\eta)k$ and are thus in $\mathcal{C}_{\mathsf{frozen}}^\sigma$. This, together with $C'(\sigma) \neq \mathsf{True}$ and $\mathcal{C}' \notin \mathcal{C}_{\mathsf{frozen}}^\sigma \cup \mathcal{C}_{\mathsf{sep}}$, implies that $C' \in \mathcal{C}_{\mathsf{bad}}^\sigma$. A contradiction. $\qquad\blacksquare$

Similar to Lemma 6.12 and by Corollary 6.5, the efficiency of the final rejection sampling boils down to the size of the remaining components in $\Phi^\sigma$ where $\sigma$ is the partial assignment on Line 5 of $\texttt{SolutionSampling}(\Phi)$. Guided by these intuition, we will keep track of the size of $\mathcal{C}_{\mathsf{con}}^\sigma$ and truncate the program if it gets too large. In addition, we will halt the program if some component in $\Phi^\sigma$ is large upon the final rejection sampling.

Let $s \geq 1$ be the truncation parameter to be optimized later. We formalize our actual algorithms in Algorithm 6 and highlight the place where truncation happens. For convenience, we overload $\texttt{SolutionSampling}()$, $\texttt{MarginSample}()$, and $\texttt{MarginOverflow}()$ with the addition parameter $s$, and they reduce to the original version if $s = +\infty$.

Similarly as Fact 6.8, checking components' sizes can be done efficiency.

**Fact 6.13.** *With $\widetilde{O}(n)$ pre-processing time, the runtime of Line 5 of $\texttt{SolutionSampling}(\Phi, s)$ and Line 1 of $\texttt{MarginOverflow}(\sigma, v, s)$ is $\widetilde{O}(1)$.*

An immediate corollary of Lemma 6.12 is the following efficiency guarantee for the leaf recursion of $\texttt{MarginOverflow}(\sigma, v, s)$.

---
**Algorithm 6:** The Actual Algorithms
---

**Procedure** SolutionSampling($\Phi, s$):

1    Obtain $\mathcal{V}_{\mathsf{sep}}, \mathcal{C}_{\mathsf{sep}} \leftarrow$ ConstructSep($\mathcal{V}$)

2    Initialize $\sigma \leftarrow \star^{\mathcal{V}}$

3    **foreach** $i = 1$ **to** $n$ **do**

4      | **if** $v_i \in \mathcal{V}^{\sigma}_{\mathsf{alive}}$ **then** Update $\sigma(v_i) \leftarrow$ MarginSample($\sigma, v_i$)

     **end**

5    **if** *some connected component in* $\Phi^{\sigma}$ *has* $> s$ *clauses* **then Halt**      /* Truncation */

6    $\sigma \leftarrow$ RejectionSampling($\sigma, \Lambda(\sigma)$)

7    **return** $\sigma$

   **end**

**Procedure** MarginSample($\sigma, v, s$):

1    Sample $\sigma(v) \sim \tau$

2    **if** $\sigma(v) = \bigstar$ **then return** MarginOverflow($\sigma, v, s$)

3    **else return** $\sigma(v)$

   **end**

**Procedure** MarginOverflow($\sigma, v, s$):

1    **if** $|\mathcal{C}^{\sigma}_{\mathsf{con}}| > s$ **then Halt**                       /* Truncation */

2    Let $u \leftarrow$ NextVar($\sigma$)

3    **if** $u \neq \perp$ **then**

4      | Sample $\sigma(u) \sim \tau$

5      | **if** $\sigma(u) = \bigstar$ **then** Update $\sigma(u) \leftarrow$ MarginOverflow($\sigma, u, s$)

6      | **return** MarginOverflow($\sigma, v, s$)

     **else**

7      | **return** BernoulliFactory($b_1, b_2, \ldots$) where $b_1, b_2, \ldots$ are independent samples
        provided by executing RejectionSampling($\sigma, v$)

     **end**

   **end**

---

**Corollary 6.14.** *If $\sigma$ satisfies Assumption 6.2 and* NextVar($\sigma$) $= \perp$, *then the expected runtime of* MarginOverflow($\sigma, v, s$) *is*

$$\widetilde{O}\left(\frac{s}{\xi^2} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right).$$

The runtime of the final rejection sampling is also controlled by the truncation parameter and Corollary 6.5.

**Corollary 6.15.** *Assume $\sigma$ satisfies Assumption 6.2. Then* RejectionSampling($\sigma, \Lambda(\sigma)$) *on Line 6 of* SolutionSampling($\Phi, s$) *runs in expected time*

$$\widetilde{O}\left(n \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right),$$

*where each $\mathcal{C}_i$ is from Line 1 of* RejectionSampling($\sigma, \Lambda(\sigma)$).

In addition, since the difference only comes from the truncation, Corollary 6.11 allows us to bound the distance between algorithm's output and a uniform solution of $\Phi$ in terms of the probability that the program halts (i.e., truncation happens).

**Corollary 6.16.** `SolutionSampling(Φ, s)` *terminates almost surely and has output distribution* $p_{\mathsf{halt}}(\Phi, s)$*-close to* $\mu$ *in the total variation distance, where*

$$p_{\mathsf{halt}}(\Phi, s) = \mathbf{Pr}\left[\textit{truncation happens during the algorithm}\right].$$

## 6.2   The Recursive Cost Tree and the Simulation Tree

Now we turn to the most technical part: The analysis of the efficiency and $p_{\mathsf{halt}}(\Phi, s)$.

To this end, we use the notion of the recursive cost tree and the simulation tree similar to [HWY22]. The former captures the execution of a single `MarginOverflow(σ, v, s)`, and the latter represents the whole execution of `SolutionSampling(Φ, s)`.

**Definition 6.17** (Recursive Cost Tree). Let $\sigma$ be a partial assignment satisfying Assumption 6.2. We define the *recursive cost tree* for $\sigma$ as $\mathcal{T}_\sigma$. Here $\mathcal{T}_\sigma$ is a rooted tree with nodes labeled by distinct[8] partial assignments $\pi$ and edges labeled by values $\rho$ in $[0, 1]$ as follows:

- The root of $\mathcal{T}_\sigma$ is $\sigma$ and its depth is defined to be 0.
- For $i = 0, 1, \ldots$, let $\pi \in \mathcal{T}_\sigma$ be a node of depth $i$.

  If $|\mathcal{C}_{\mathsf{con}}^\pi| > s$, then we leave $\pi$ as a *recursing truncated* leaf node.

  Otherwise, let $u = \mathsf{NextVar}(\pi)$ and we proceed as follows:

  - If $u = \bot$, then we leave $\pi$ as a *Bernoulli* leaf node.
  - Otherwise, let $\pi_0, \pi_1, \pi_\star$ equal $\pi$ except that we fix $u$ to $0, 1, \star$ respectively. Then we append $\pi_0, \pi_1, \pi_\star$ as the child nodes of $\pi$ and label the edges by

$$\rho(\pi \to \pi_0) = \mu_u^\pi(0), \quad \rho(\pi \to \pi_1) = \mu_u^\pi(1), \quad \rho(\pi \to \pi_\star) = \delta.$$

The edge values reflect the `MarginOverflow(σ, v, +∞)` recursion without truncation, and it is an overestimate for the `MarginOverflow(σ, v, s)`. In addition, $\mathcal{T}_\sigma$ stops either at a Bernoulli leaf node, which corresponds to a regular leaf recursion and is ready for Bernoulli factory on Line 7, or at a recursing truncated leaf node, which corresponds to a truncation on Line 1 during the recursion.

We remark that the edge value only depends on the partial assignments of the endpoints. This is why we can use a single symbol $\rho$ without confusion.

**Remark 6.18.** Let $\sigma$ be a partial assignment satisfying Assumption 6.2 where $\sigma(v) = \star$ and the rest values are $0/1/\star$. We show a one-to-one correspondence between nodes in $\mathcal{T}_\sigma$ and the execution of `MarginOverflow(σ, v, s)`.

The starting point `MarginOverflow(σ, v, s)` corresponds to the root of $\mathcal{T}_\sigma$. Recall the algorithm description from Algorithm 6. Assume we just enter `MarginOverflow(π, w, s)`, which by induction corresponds to the node $\pi \in \mathcal{T}_\sigma$. Then after checking if $|\mathcal{C}_{\mathsf{con}}^\pi| > s$ (i.e., if $\pi$ is a recursing truncated leaf node) on Line 1, we will compute $u = \mathsf{NextVar}(\pi)$ and perform the Bernoulli factory if $u = \bot$, i.e., $\pi$ is a Bernoulli leaf node as designed. If $|\mathcal{C}_{\mathsf{con}}^\pi| \le s$ and $u \ne \bot$, the algorithm will update the assignment of $u$. Then we have $\tau(\star) = \delta$ probability of executing `MarginOverflow(π_⋆, u, s)` which means in $\mathcal{T}_\sigma$ proceeding to the child node $\pi_\star$. Afterwards, we will run `MarginOverflow(π_b, w, s)` for $b \in \{0, 1\}$, corresponding to the child node $\pi_b \in \mathcal{T}_\sigma$.

The definition of $\rho(\pi \to \pi_\star)$ is already explained above. For $b \in \{0, 1\}$, the probability of visiting $\pi_b$ is upper bounded by the corresponding probability with no truncation, i.e., setting $s = +\infty$ in `MarginOverflow(π_⋆, u, s)`, which, by Lemma 6.9, is exactly $\mu_u^\pi(b) = \rho(\pi \to \pi_b)$.

---

[8]The nodes are distinct by the definition, where the partial assignments of the child nodes of $\pi$ fix the value of $u = \mathsf{NextVar}(\pi)$ from $\star$ to $0/1/\star$.

To study the runtime of the whole `SolutionSampling`$(\Phi, s)$, we define the following simulation tree on top of recursive cost trees.

**Definition 6.19** (Simulation Tree). We define the *simulation tree* as $\mathcal{T}_{\mathsf{sim}}$. Here $\mathcal{T}_{\mathsf{sim}}$ is a rooted tree with nodes labeled by distinct[9] partial assignments $\pi$ and edges labeled by values $\rho$ in $[0, 1]$ as follows:

- The root of $\mathcal{T}_{\mathsf{sim}}$ is $\star^{\mathcal{V}}$ and its depth is defined to be 0.
- For $i = 0, 1, \ldots$, let $\pi \in \mathcal{T}_{\mathsf{sim}}$ be a node of depth $i$.
    - If $\pi$ has a $\bigstar$, then we say $\pi$ is a *recursing* node and we append $\mathcal{T}_\pi$ here.
    - Otherwise, let $u$ be the variable in $\mathcal{V}^\pi_{\mathsf{alive}}$ with minimal index:
        * If $u$ does not exist and each connected component in $\Phi^\pi$ has at most $s$ clauses, then we leave $\pi$ as a *sampling* leaf node.
        * If $u$ does not exist and some connected component in $\Phi^\pi$ has $> s$ clauses, then we leave $\pi$ as a *sampling truncated* leaf node.
        * Otherwise $u$ exists. Let $\pi_0, \pi_1, \pi_\bigstar$ equal $\pi$ except that we fix $u$ to $0, 1, \bigstar$ respectively. Then we append $\pi_0, \pi_1, \pi_\bigstar$ as the child nodes of $\pi$ and label the edges by

$$\rho(\pi \to \pi_0) = \mu_u^\pi(0), \quad \rho(\pi \to \pi_1) = \mu_u^\pi(1), \quad \rho(\pi \to \pi_\bigstar) = \delta.$$

Intuitively corresponding to `SolutionSampling`$(\Phi, s)$, a recursing node means that we are about to do `MarginOverflow()` inside a `MarginSample()` on Line 4, a sampling leaf node means that we now perform the final rejection sampling on Line 6, and a sampling truncated leaf node is analogous to the one in recursive cost tree that truncation happens on Line 5.

We remark that the edge value $\rho$ is indeed consistent with the one in the definition of the recursive cost tree, as both of them refer to the probability of the one-step update of the partial assignments of the endpoints of the edge without truncation. Thus we use the same symbol.

**Remark 6.20.** Similar to the recursive cost tree, there is a one-to-one correspondence between nodes in $\mathcal{T}_{\mathsf{sim}}$ and the execution of `SolutionSampling`$(\Phi, s)$. Let $\sigma$ be the partial assignment that `SolutionSampling`$(\Phi, s)$ maintains.

At the beginning, $\sigma = \star^{\mathcal{V}}$ and it is the root of $\mathcal{T}_{\mathsf{sim}}$. Each time we update $\sigma(v_i)$ for $v_i \in \mathcal{V}^\sigma_{\mathsf{alive}}$ on Line 4, this $v_i = u$ has the minimal index in $\mathcal{V}^\sigma_{\mathsf{alive}}$ since the for-loop on Line 3 goes in the ascending order. Then, based on the outcome of $b \leftarrow$ `MarginSample`$(\sigma, v_i, s)$, we update $\sigma$ to $\sigma_b, b \in \{0, 1\}$. Recall that `MarginSample`$(\sigma, v_i, s)$ may call `MarginOverflow`$(\sigma_\bigstar, v_i, s)$. Together with $\sigma_0, \sigma_1$, these $\sigma_b$'s are presented in $\mathcal{T}_{\mathsf{sim}}$ as the child nodes of $\sigma$, where $\sigma_\bigstar$ is a recursing node and we append $\mathcal{T}_{\sigma_\bigstar}$ and follow the correspondence in Remark 6.18.

To see the edge values, the probability of obtaining $\sigma_\bigstar$ is precisely $\tau(\bigstar) = \delta = \rho(\sigma \to \sigma_\bigstar)$. For $b \in \{0, 1\}$, the probability of visiting $\sigma_b$ is upper bounded by the corresponding probability when we ignore truncation, which in turn is exactly $\mu_{v_i}^\sigma(b) = \rho(\sigma \to \sigma_b)$ by Corollary 6.10.

Finally on Line 5, we reach a partial assignment $\sigma$ ready for the final rejection sampling Line 6. Depending on the components' sizes in $\Phi^\sigma$, it gives a sampling (truncated) leaf node.

By the correspondence above, we see that Assumption 6.2 is always preserved.

**Fact 6.21.** *Assumption 6.2 holds for any node in* $\mathcal{T}_{\mathsf{sim}}$.

---

[9] This is also clear from the definition of the simulation tree and expanding the construction of the recursive cost tree. In general, the partial assignments of the child nodes fix the variable from $\star$ to $0/1/\bigstar$.

For convenience, we define the following quantities:

- For a node $\pi$ in $\mathcal{T}_{\mathsf{sim}}$, $\rho(\pi)$ denotes the product of the edge values from the root of $\mathcal{T}_{\mathsf{sim}}$ to $\pi$.
- $\mathcal{N}_{\mathsf{rec}}$ denotes the set of recursing nodes of $\mathcal{T}_{\mathsf{sim}}$, and define $d_{\mathsf{rec}} = \max_{\sigma \in \mathcal{N}_{\mathsf{rec}}} \mathsf{depth}(\mathcal{T}_\sigma)$ to be the maximal depth of the recursive cost trees encountered.
- $\mathcal{N}_{\mathsf{samp\text{-}trunc}}$ denotes the set of sampling truncated leaf nodes of $\mathcal{T}_{\mathsf{sim}}$, corresponding to Line 5 of $\texttt{SolutionSampling}(\Phi, s)$.
- $\mathcal{N}_{\mathsf{rec\text{-}trunc}}$ denotes the set of recursing truncated leaf nodes of $\mathcal{T}_{\mathsf{sim}}$, corresponding to Line 1 of $\texttt{MarginOverflow}(\sigma, v, s)$.
- $\mathcal{N}_{\mathsf{trunc}} = \mathcal{N}_{\mathsf{samp\text{-}trunc}} \cup \mathcal{N}_{\mathsf{rec\text{-}trunc}}$ denotes the set of all truncated leaf nodes.

At this point, we can bound $p_{\mathsf{halt}}(\Phi, s)$ using the leaf nodes' information of $\mathcal{T}_{\mathsf{sim}}$.

**Lemma 6.22.** $p_{\mathsf{halt}}(\Phi, s) \leq \sum_{\pi \in \mathcal{N}_{\mathsf{trunc}}} \rho(\pi)$.

*Proof.* By Remark 6.20, we have a one-to-one correspondence between the truncated leaf nodes in $\mathcal{T}_{\mathsf{sim}}$ and the place where truncation happens during $\texttt{SolutionSampling}(\Phi, s)$. In addition, for any partial assignment $\pi$, $\rho(\pi)$ upper bounds the probability that the algorithm visits $\pi$. Therefore by the definition of $\mathcal{N}_{\mathsf{trunc}}$ and $p_{\mathsf{halt}}(\Phi, s)$ from Corollary 6.16, we have

$$p_{\mathsf{halt}}(\Phi, s) = \mathbf{Pr}\left[\text{reaching some node in } \mathcal{N}_{\mathsf{trunc}} \text{ during the algorithm}\right] \leq \sum_{\pi \in \mathcal{N}_{\mathsf{trunc}}} \rho(\pi). \qquad \square$$

The runtime can also be analyzed similarly.

**Lemma 6.23.** $\texttt{SolutionSampling}(\Phi, s)$ *runs in expected time*

$$\widetilde{O}\left(n \cdot (1+\delta)^{d_{\mathsf{rec}}} \cdot \frac{s}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right).$$

*Proof.* Recall the description of $\texttt{SolutionSampling}(\Phi, s)$ from Algorithm 6. The runtime of Lines 1 and 2 is $\widetilde{O}(n)$ by Fact 4.1. Lines 5 and 6 runs in expected time $\widetilde{O}\left(n \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right)$ by Fact 6.13 and Corollary 6.15. Now it remains to bound the runtime of Lines 3 and 4.

Firstly checking condition on Line 4 takes $\widetilde{O}(n)$ time in total by Fact 6.1. Now assume we call $\texttt{MarginSample}(\sigma, v_i, s)$ on Line 4. Let $\sigma_\star$ equal $\sigma$ except we fix $v_i$ to $\star$. By Remark 6.20, $\sigma_\star$ is a recursing node where we append the recursive cost tree $\mathcal{T}_{\sigma_\star}$ for $\texttt{MarginOverflow}(\sigma_\star, v_i, s)$. The expected runtime of $\texttt{MarginOverflow}(\sigma_\star, v_i, s)$ has two parts:

(i) Visiting partial assignments $\pi$, checking $|\mathcal{C}_{\mathsf{con}}^\pi|$, calculating $\mathsf{NextVar}(\pi)$, and sampling from $\tau$.

(ii) Performing Bernoulli factory on leaf recursions if not truncated.

Let $\rho'(\pi)$ be the product of the edge values from the root of $\mathcal{T}_{\sigma_\star}$ to $\pi$. By the correspondence described in Remark 6.18, the probability of visiting a partial assignment $\pi \in \mathcal{T}_{\sigma_\star}$ conditioned on starting at $\sigma_\star$ is upper bounded by $\rho'(\pi)$. Without loss of generality, we expand $\mathcal{T}_{\sigma_\star}$ to a complete ternary tree where the parent-to-child edge weights are $\zeta, 1 - \zeta, \delta$ respectively for some $\zeta \in [0, 1]$. This is consistent with the existing edge values $\rho$, where $\zeta = \mu_u^\pi(0)$ for node $\pi$ and $u = \mathsf{NextVar}(\pi)$. At this point, we have

$$\mathbb{E}[\text{runtime for (i)}] \leq \sum_{\pi \in \mathcal{T}_{\sigma_\star}} \rho'(\pi) \cdot \widetilde{O}(1) \leq \sum_{d=0}^{d_{\mathsf{rec}}} (1+\delta)^d \cdot \widetilde{O}(1) = \widetilde{O}\left(\frac{(1+\delta)^{d_{\mathsf{rec}}}}{\delta}\right).$$

28

By Corollary 6.14, we can bound the runtime of (ii) similarly

$$\mathbb{E}[\text{runtime for (ii)}] \leq \sum_{\pi \in \mathcal{T}_{\sigma_\star} \text{ is a Bernoulli leaf node}} \rho'(\pi) \cdot \widetilde{O}\left(\frac{s}{\xi^2} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right)$$

$$\leq (1+\delta)^{d_{\text{rec}}} \cdot \widetilde{O}\left(\frac{s}{\xi^2} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right),$$

where we use the fact that Bernoulli factory happens only on leaf nodes. Since we only have $\rho(\sigma \to \sigma_\star) = \delta$ probability of executing $\texttt{MarginOverflow}(\sigma_\star, v_i, s)$, we have

$$\mathbb{E}[\text{runtime of } \texttt{MarginSample}(\sigma, v_i)] = \widetilde{O}(1) + \delta \cdot (\mathbb{E}[\text{runtime for (i)}] + \mathbb{E}[\text{runtime for (ii)}])$$

$$\leq \widetilde{O}\left((1+\delta)^{d_{\text{rec}}} \cdot \left(1 + \frac{\delta s}{\xi^2} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right)\right)$$

$$\leq \widetilde{O}\left((1+\delta)^{d_{\text{rec}}} \cdot \left(1 + \frac{s}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right)\right)$$

$$= \widetilde{O}\left((1+\delta)^{d_{\text{rec}}} \cdot \frac{s}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right), \quad (\text{since } s \geq 1 \geq \xi)$$

where we use $\alpha \geq 1/k^3$ and $\delta = \xi/(k^{40}\alpha) \leq \widetilde{O}(\xi)$ in the third step. Hence

$$\mathbb{E}[\text{runtime of Lines 3 and 4}] \leq \widetilde{O}\left(n \cdot (1+\delta)^{d_{\text{rec}}} \cdot \frac{s}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right).$$

Putting everything together, we have

$$\mathbb{E}[\text{total runtime}] \leq \widetilde{O}\left(n + n \cdot (1+\delta)^{d_{\text{rec}}} \cdot \frac{s}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\} + n \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right)$$

$$= \widetilde{O}\left(n \cdot (1+\delta)^{d_{\text{rec}}} \cdot \frac{s}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right). \qquad \square$$

# 7 Truncation Analysis

Given Lemma 6.22, Corollary 6.16, and Lemma 6.23, we need to carefully select the truncation parameter $s$ such that both $d_{\text{rec}}$ and $p_{\text{halt}}(\Phi, s)$ can be bounded. The goal of this section is to establish such relations and prove the following formal statements. We will still assume $(\Phi, k, \alpha, n, \xi, \eta, D)$ is nice and omit it from all the statements.

**Lemma 7.1.** $d_{\text{rec}} \leq s \cdot k + 1$.

**Lemma 7.2.** *Assume* $6k^4\alpha \log(n) < s \leq n/2^{5k/\log(k)}$. *Then*

$$p_{\text{halt}}(\Phi, s) \leq n^{10}(1+\delta)^{d_{\text{rec}}+1} \cdot k^{-s/(6k^4\alpha)}.$$

Let $\sigma$ be a partial assignment. For convenience, we recall the definitions:

- $v \in \mathcal{V}_{\text{alive}}^\sigma$ iff (i) $\sigma(v) = \text{☆}$ and $v \notin \mathcal{V}_{\text{sep}}$, and (ii) for every clause $C \in \mathcal{C} \backslash \mathcal{C}_{\text{sep}}$, either $C(\sigma) = \text{True}$ or $|\text{vbl}(C) \cap \Lambda(\sigma) \backslash (\mathcal{V}_{\text{sep}} \cup \{v\})| \geq (2/3 - 2\eta)k$.

- $C \in \mathcal{C}_\star^\sigma$ iff there exists some $v \in \text{vbl}(C)$ that $\sigma(v) = \text{★}$.

- $C \in \mathcal{C}_{\text{frozen}}^\sigma$ iff (i) $C(\sigma) \neq \text{True}$ and $C \notin \mathcal{C}_{\text{sep}}$, and (ii) $|\text{vbl}(C) \cap \Lambda(\sigma) \backslash \mathcal{V}_{\text{sep}}| < 1 + (2/3 - 2\eta)k$.

- $C \in \mathcal{C}_{\mathsf{bad}}^{\sigma}$ iff (i) $C(\sigma) \neq \mathsf{True}$ and $C \notin \mathcal{C}_{\mathsf{frozen}}^{\sigma} \cup \mathcal{C}_{\mathsf{sep}}$, and (ii) for any $v \in \mathsf{vbl}(C) \setminus \mathcal{V}_{\mathsf{sep}}$ with $\sigma(v) = \star$, there exists some $C' \in \mathcal{C}_{\mathsf{frozen}}^{\sigma}$ such that $v \in \mathsf{vbl}(C')$.

- $C \in \mathcal{C}_{\mathsf{int}}^{\sigma}(v)$ iff (i) $C \in \mathcal{C}_{\mathsf{frozen}}^{\sigma} \cup \mathcal{C}_{\mathsf{bad}}^{\sigma} \cup \mathcal{C}_{\mathsf{sep}}$, and (ii) either $v \in \mathsf{vbl}(C)$ or there exists some $C' \in \mathcal{C}_{\mathsf{int}}^{\sigma}(v)$ that $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\sigma) \neq \emptyset$.

- $C \in \mathcal{C}_{\mathsf{con}}^{\sigma}(v)$ iff $C \in \mathcal{C}_{\mathsf{int}}^{\sigma}(v)$, or $v \in \mathsf{vbl}(C)$, or there exists some $C' \in \mathcal{C}_{\mathsf{int}}^{\sigma}(v)$ that $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\sigma) \neq \emptyset$.

- $\mathcal{C}_{\mathsf{con}}^{\sigma}$ is the union of $\mathcal{C}_{\mathsf{con}}^{\sigma}(v)$ for all $v$ with $\sigma(v) = \star$.

We will prove Lemma 7.1 in Subsection 7.1. Then we construct witnesses for truncated nodes in Subsection 7.2 and prove Lemma 7.2 in Subsection 7.3.

## 7.1 Size-to-Depth Reduction

We start by relating $s$ and $d_{\mathsf{rec}}$, and show that small truncation size implies small depth in the recursive cost trees. To this end, we will use $|\mathcal{C}_{\mathsf{con}}^{\pi}|$ as an intermediate measure for partial assignments $\pi$ in recursive cost trees. Indeed, $|\mathcal{C}_{\mathsf{con}}^{\pi}|$ is upper bounded by $s$ by truncation, and we only need to lower bound it in terms of $d_{\mathsf{rec}}$.

We start by proving the connectivity, which reduces to the following technical lemma showing that $\mathcal{C}_{\mathsf{int}}^{\pi}$ and $\mathcal{C}_{\mathsf{con}}^{\pi}$ are increasing in $\pi$.

**Lemma 7.3.** *Let $\pi$ and $\pi'$ be partial assignments. Assume $\pi'$ extends $\pi$ by fixing some variable in $\mathcal{V}_{\mathsf{alive}}^{\pi}$. Then $\mathcal{C}_{\mathsf{int}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$ and $\mathcal{C}_{\mathsf{con}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{con}}^{\pi'}(v)$ hold for any $v$ with $\pi(v) = \star$.*

*Proof.* We first show $\mathcal{C}_{\mathsf{int}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$. Let $C \in \mathcal{C}_{\mathsf{int}}^{\pi}(v)$ and we verify the conditions for $C \in \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$:

- Condition (i). By the condition (i) for $C \in \mathcal{C}_{\mathsf{int}}^{\pi}(v)$, we have $C \in \mathcal{C}_{\mathsf{frozen}}^{\pi} \cup \mathcal{C}_{\mathsf{bad}}^{\pi} \cup \mathcal{C}_{\mathsf{sep}}$. Then $C \in \mathcal{C}_{\mathsf{frozen}}^{\pi'} \cup \mathcal{C}_{\mathsf{bad}}^{\pi'} \cup \mathcal{C}_{\mathsf{sep}}$ since $\mathcal{C}_{\mathsf{frozen}}^{\pi} \subseteq \mathcal{C}_{\mathsf{frozen}}^{\pi'}, \mathcal{C}_{\mathsf{bad}}^{\pi} \subseteq \mathcal{C}_{\mathsf{bad}}^{\pi'}$ by Fact 6.7.

- Condition (ii). We have two cases based on the condition (ii) for $C \in \mathcal{C}_{\mathsf{int}}^{\pi}(v)$:

  - If $v \in \mathsf{vbl}(C)$, then the same reason holds for $C \in \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$.
  - Otherwise, there exists some $C' \in \mathcal{C}_{\mathsf{int}}^{\pi}(v)$ that $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\pi) \neq \emptyset$. Now note that $\pi'$ extends $\pi$ on a variable in $\mathcal{V}_{\mathsf{alive}}^{\pi}$, which, by condition (i) and Fact 6.7, is not contained in $C'$. Thus $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\pi') = \mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\pi) \neq \emptyset$, which means the condition (ii) here holds due to the same $C'$.

Now we prove $\mathcal{C}_{\mathsf{con}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{con}}^{\pi'}(v)$ with similar arguments. Let $C \in \mathcal{C}_{\mathsf{con}}^{\pi}(v)$ and we verify $C \in \mathcal{C}_{\mathsf{con}}^{\pi'}(v)$:

- If $C \in \mathcal{C}_{\mathsf{int}}^{\pi}(v)$, then $C \in \mathcal{C}_{\mathsf{int}}^{\pi'}(v) \subseteq \mathcal{C}_{\mathsf{con}}^{\pi'}(v)$ since $\mathcal{C}_{\mathsf{int}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$.

- If $v \in \mathsf{vbl}(C)$, then $C \in \mathcal{C}_{\mathsf{con}}^{\pi'}(v)$ by the same reason.

- Otherwise, there exists some $C' \in \mathcal{C}_{\mathsf{int}}^{\pi}(v)$ that $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\pi) \neq \emptyset$. Note that we have $C' \in \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$ since $\mathcal{C}_{\mathsf{int}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$. As $\pi'$ differs from $\pi$ on a variable in $\mathcal{V}_{\mathsf{alive}}^{\pi}$, by Fact 6.7, this variable is not in $C' \in \mathcal{C}_{\mathsf{con}}^{\pi}(v) \subseteq \mathcal{C}_{\mathsf{frozen}}^{\pi} \cup \mathcal{C}_{\mathsf{bad}}^{\pi} \cup \mathcal{C}_{\mathsf{sep}}$. Thus $\mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\pi') = \mathsf{vbl}(C') \cap \mathsf{vbl}(C) \cap \Lambda(\pi) \neq \emptyset$, which, combined with $C' \in \mathcal{C}_{\mathsf{int}}^{\pi'}(v)$, implies $C \in \mathcal{C}_{\mathsf{con}}^{\pi'}(v)$. $\square$

As a result, we can lower bound $|\mathcal{C}_{\mathsf{con}}^{\pi}|$ by the depth of $\pi$ in $\mathcal{T}_{\sigma}$.

**Corollary 7.4.** *Let $\sigma \in \mathcal{N}_{\mathsf{rec}}$ and $\pi \in \mathcal{T}_{\sigma}$. Then $|\mathcal{C}_{\mathsf{con}}^{\pi}| \geq \mathsf{depth}(\pi, \mathcal{T}_{\sigma})/k$ where $\mathsf{depth}(\pi, \mathcal{T}_{\sigma})$ is the depth of $\pi$ in $\mathcal{T}_{\sigma}$.*

30

*Proof.* Let $L = \mathsf{depth}(\pi, \mathcal{T}_\sigma)$. Along the path from $\sigma$ to $\pi$, we fix $L$ distinct variables $v_1, v_2, \ldots, v_L$. Let partial assignments $\pi_0, \pi_1, \ldots, \pi_L$ be the evolution of the process, i.e., $\pi_0 = \sigma$, $\pi_L = \pi$, and each $\pi_i$ extends $\pi_{i-1}$ by fixing $v_i$. Since $v_i = \mathsf{NextVar}(\pi_{i-1}) \in \mathcal{V}_{\mathsf{con}}^{\pi_{i-1}}$, there exists $C_i \in \mathcal{C}_{\mathsf{con}}^{\pi_{i-1}}$ such that $v_i \in \mathsf{vbl}(C_i)$. By Lemma 7.3, these $C_i$'s remain in $\mathcal{C}_{\mathsf{con}}^\pi$. Thus $|\mathcal{C}_{\mathsf{con}}^\pi|$ is at least the number of distinct clauses in $C_1, C_2, \ldots, C_L$, which, in turn, is at least $L/k$. □

Now Lemma 7.1 follows immediately.

*Proof of Lemma 7.1.* Recall that $d_{\mathsf{rec}} = \max_{\sigma \in \mathcal{N}_{\mathsf{rec}}} \mathsf{depth}(\mathcal{T}_\sigma) = \max_{\sigma \in \mathcal{N}_{\mathsf{rec}}, \pi \in \mathcal{T}_\sigma} \mathsf{depth}(\pi, \mathcal{T}_\sigma)$. Let $\sigma$ and $\pi$ achieve $\mathsf{depth}(\pi, \mathcal{T}_\sigma) = d_{\mathsf{rec}}$. If $\pi$ is a Bernoulli leaf node, then $|\mathcal{C}_{\mathsf{con}}^\pi| \leq s$. By Corollary 7.4, we have $|\mathcal{C}_{\mathsf{con}}^\pi| \geq d_{\mathsf{rec}}/k$ and thus $d_{\mathsf{rec}} \leq s \cdot k$. Now assume $\pi \in \mathcal{N}_{\mathsf{rec\text{-}trunc}}$ is a recursing truncated leaf node.

If $\pi = \sigma$, then $d_{\mathsf{rec}} = 0$ trivially. Otherwise let $\pi'$ be the parent node of $\pi$. Then $\mathsf{depth}(\pi', \mathcal{T}_\sigma) = d_{\mathsf{rec}} - 1$ and we have $|\mathcal{C}_{\mathsf{con}}^{\pi'}| \geq (d_{\mathsf{rec}} - 1)/k$ by Corollary 7.4. Since $\pi'$ is not truncated, we also have $|\mathcal{C}_{\mathsf{con}}^{\pi'}| \leq s$, which implies $d_{\mathsf{rec}} \leq s \cdot k + 1$. □

## 7.2 Witness for Truncation

To establish Lemma 7.2, we will construct succinct witnesses for truncated nodes $\mathcal{N}_{\mathsf{trunc}}$. Then in Subsection 7.3, we will enumerate all possible witnesses and apply a union bound to show that with high probability none of them appears. Though we have two types of truncation $\mathcal{N}_{\mathsf{rec\text{-}trunc}}$ and $\mathcal{N}_{\mathsf{samp\text{-}trunc}}$, the witness construction is similar.

Let $\pi$ be a partial assignment triggering truncation. In a nutshell, the witness will consist of many connected clauses where most of the clauses are either frozen (i.e., in $\mathcal{C}_{\mathsf{frozen}}^\pi$) or contains $\star$ (i.e., in $\mathcal{C}_\star^\pi$). The former case, together with the locally sparse properties of $\Phi$, indicates that many variables in $\pi$ are fixed towards the bad direction that does not satisfy the clauses. The latter case should also be rare since, by local uniformity, each $\star$ appears with probability $\delta \ll 1$.

### Truncation inside the Margin Overflow

We start with the recursing truncated nodes $\pi \in \mathcal{N}_{\mathsf{rec\text{-}trunc}}$, which corresponds to partial assignments $\sigma \in \mathcal{N}_{\mathsf{rec}}$ and $\pi \in \mathcal{T}_\sigma \cap \mathcal{N}_{\mathsf{trunc}}$. We want to zoom in to the frozen clauses $\mathcal{C}_{\mathsf{frozen}}^\pi$ of $\pi$ since each clause there is still not satisfied and yet many variables within are fixed (to the unsatisfying direction).

However, $\mathcal{C}_{\mathsf{frozen}}^\pi$ alone may not be connected and we cannot afford the enumeration. Therefore, we put in $\mathcal{C}_{\mathsf{bad}}^\pi$ and $\mathcal{C}_{\mathsf{sep}}$. These clauses do not contain many fixed variables (indeed by definition), but at least they are also not satisfied and is close to clauses in $\mathcal{C}_{\mathsf{frozen}}^\pi$. Thus we still have control for the variables within.

Unfortunately, at this point we still cannot guarantee large connected components. At best, we will only have connected components $\mathcal{C}_{\mathsf{int}}^\pi(v_i)$ where $v_1, v_2, \ldots, v_t$ are the $\star$'s in $\pi$; and these are not necessarily connected to each other. Indeed, the definition of $\mathcal{C}_{\mathsf{con}}^\pi$ involves taking one step further from $\mathcal{C}_{\mathsf{int}}^\pi(v_i)$; only after that it will be truncated due to exceeding size $s$.

The final thing we can do is to incorporate clauses in $\mathcal{C}_\star^\pi$, which is still acceptable since we have control for the probability that we encounter any fixed $\star$. This is indeed the case here: By Lemma 7.3, we connect $\mathcal{C}_{\mathsf{int}}^\pi(v_i)$'s by including edges from $\mathcal{C}_\star^\pi$. Put differently, each $v_i$ is contained in $\mathcal{C}_\star^\pi(v_j)$ for some previous $v_j$.

**Lemma 7.5.** *Let $\sigma \in \mathcal{N}_{\mathsf{rec}}$ and $\pi \in \mathcal{T}_\sigma$. Let $v_1, v_2, \ldots, v_t$ be the $\star$'s in $\pi$ in the order of the path from $\sigma$ to $\pi$.[10] Then for any $i \geq 2$, there exists some $j < i$ and $C \in \mathcal{C}_{\mathsf{con}}^\pi(v_j)$ such that $v_i \in \mathsf{vbl}(C)$.*

---

[10]That is, $v_1$ is the unique $\star$ in $\sigma$ and $v_t$ is the last variable fixed to $\star$ before reaching $\pi$.

*Proof.* Let $\pi' \in \mathcal{T}_\sigma$ be the ancestor of $\pi$ that fixes $v_i$ to $\star$. That is, $\mathsf{NextVar}(\pi') = v_i$. By the definition of $\mathsf{NextVar}()$, there exists $C \in \mathcal{C}_{\mathsf{con}}^{\pi'}$ that $v_i \in \mathsf{vbl}(C)$. Since the $\star$'s in $\pi'$ are $v_1, \ldots, v_{i-1}$, we have $\mathcal{C}_{\mathsf{con}}^{\pi'} = \bigcup_{j<i} \mathcal{C}_{\mathsf{con}}^{\pi'}(v_j)$. Therefore, there exists $j < i$ such that $C \in \mathcal{C}_{\mathsf{con}}^{\pi'}(v_j)$. Now by Lemma 7.3, we know $C$ remains in $\mathcal{C}_{\mathsf{con}}^{\pi}(v_j)$ as desired, since $\pi$ is obtained from $\pi'$ by repeatedly fixing alive variables. $\square$

As a result, we can connect $\mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$'s efficiently in a spanning tree fashion using edges in $\mathcal{C}_{\star}^{\pi}$.

**Corollary 7.6.** *Let $\sigma \in \mathcal{N}_{\mathsf{rec}}$ and $\pi \in \mathcal{T}_\sigma$. Let $v_1, v_2, \ldots, v_t$ be the $\star$'s in $\pi$ in the order of the path from $\sigma$ to $\pi$. Then there exists $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \subseteq \mathcal{C}_{\star}^{\pi}$ disjoint from $\bigcup_i \mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$ such that the following holds:*

- *$\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \bigcup_i \mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$ is connected in $G_\Phi$, and $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$ covers $v_1, \ldots, v_t$.*
- *For any $\mathcal{C}' \subseteq \mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$, we have $|\{v \in \mathcal{V}' \mid \pi(v) = \star\}| \geq |\mathcal{C}'|$ where $\mathcal{V}' = \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)$.*

*Proof.* We construct $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$ by inspecting $\mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$ sequentially. At first, $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} = \emptyset$ and $i = 1$.

By definition, each $\mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$ is connected by itself. If at some point $i > 1$, $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \bigcup_{j \leq i} \mathcal{C}_{\mathsf{int}}^{\pi}(v_j)$ is not connected. By Lemma 7.5, there exists $j < i$ and $C \in \mathcal{C}_{\mathsf{con}}^{\pi}(v_j)$ such that $v_i \in \mathsf{vbl}(C)$ and we put this $C$ into $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$. Note that $C \in \mathcal{C}_{\star}^{\pi}$ since $v_i \in \mathsf{vbl}(C)$. On the other hand, before including $C$, $v_i$ is not covered in $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$ since otherwise $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \bigcup_{j \leq i} \mathcal{C}_{\mathsf{int}}^{\pi}(v_j)$ is already connected.

Therefore, after this process, $\mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$'s are connected by $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$ which covers all the $\star$'s. In addition, the second item holds since each newly included clause brings in at least one $\star$ distinct from all previous ones. $\square$

Define $\mathcal{C}_{\mathsf{int}}^{\pi} = \bigcup_i \mathcal{C}_{\mathsf{int}}^{\pi}(v_i)$. At this point, the witness is already in shape: $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}$. Indeed, every clause in $\mathcal{C}_{\mathsf{int}}^{\pi}$ is not satisfied by $\pi$, and $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$ contains all the $\star$'s. Now we need to show that the size of $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}$ scales with the size of $\mathcal{C}_{\mathsf{con}}^{\pi}$, which is in turn lower bounded by $s$ upon truncation.

**Lemma 7.7.** *Let $\sigma \in \mathcal{N}_{\mathsf{rec}}$ and $\pi \in \mathcal{T}_\sigma$. Then*

$$6k^4\alpha \cdot \max\left\{|\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}|, \log(n)\right\} \geq \min\left\{|\mathcal{C}_{\mathsf{con}}^{\pi}|, n/2^{2k/\log(k)}\right\}.$$

*Moreover, if $\pi \in \mathcal{N}_{\mathsf{trunc}}$ and $6k^4\alpha \log(n) < s \leq n/2^{2k/\log(k)}$, then*

$$|\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}| \geq \frac{s}{6k^4\alpha}.$$

*Proof.* Let $\mathcal{V}' = \bigcup_{C \in \mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}} \mathsf{vbl}(C)$. Then $\mathcal{V}'$ is connected in $H_\Phi$ since $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}$ is connected in $G_\Phi$ by Corollary 7.6. Thus by Proposition 3.6,

$$\mathcal{V}'' := \left\{v \mid v \in \mathcal{V}' \text{ or } v \text{ is adjacent to } \mathcal{V}'\right\} \leq 3k^4\alpha \cdot \max\left\{|\mathcal{V}'|, \lfloor k \log(n) \rfloor\right\}$$
$$\leq 3k^5\alpha \cdot \max\left\{|\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi} \cup \mathcal{C}_{\mathsf{int}}^{\pi}|, \log(n)\right\}.$$

By the definition of $\mathcal{C}_{\mathsf{con}}^{\pi}$, every clause in $\mathcal{C}_{\mathsf{con}}^{\pi}$ is contained in $\mathcal{C}_{\mathsf{int}}^{\pi}$, or is connected to some clause in $\mathcal{C}_{\mathsf{int}}^{\pi}$, or contains some $\star$ in $\pi$ and thus is connected to some clause $\mathcal{C}_{\star\text{-}\mathsf{int}}^{\pi}$ by Corollary 7.6. Therefore $\mathcal{V}''$ is the support of $\mathcal{C}_{\mathsf{con}}^{\pi}$.

Now let $\mathcal{C}'' \subseteq \mathcal{C}_{\mathsf{con}}^{\pi}$ be arbitrary and has size $\min\left\{|\mathcal{C}_{\mathsf{con}}^{\pi}|, n/2^{2k/\log(k)}\right\}$. Then $\mathcal{V}''$ is also the support of $\mathcal{C}''$. By Item 2 of Proposition 3.3 and $\eta \leq 1$, we have

$$|\mathcal{V}''| \geq \left|\bigcup_{C \in \mathcal{C}''} \mathsf{vbl}(C)\right| \geq |\mathcal{C}''| \cdot k/2 = \min\left\{|\mathcal{C}_{\mathsf{con}}^{\pi}|, n/2^{2k/\log(k)}\right\} \cdot k/2,$$

32

which completes the proof for the first half.

For the second half, notice that $\pi \in \mathcal{N}_{\mathsf{trunc}}$ additionally implies $|\mathcal{C}_{\mathsf{con}}^\pi| > s$. Therefore

$$6k^4\alpha \cdot \max\left\{|\mathcal{C}_{\star\text{-int}}^\pi \cup \mathcal{C}_{\mathsf{int}}^\pi|, \log(n)\right\} \geq \min\left\{s, n/2^{2k/\log(k)}\right\} = s.$$

Since $s > 6k^4\alpha\log(n)$, we must take the former inside the max, which gives the desired bound. $\quad\square$

Now that we have a relatively large witness. The next step for us is to show that $\mathcal{C}_{\mathsf{int}}^\pi$ contains many fixed variables in $\pi$ using the locally sparse properties of $\Phi$.

The caveat here is that, most of the structural properties in Section 3 hold only when we don't have *too many* clauses, whereas it is possible that $\mathcal{C}_{\star\text{-int}}^\pi \cup \mathcal{C}_{\mathsf{int}}^\pi$ exceeds this threshold. Though $\mathcal{C}_{\star\text{-int}}^\pi \cup \mathcal{C}_{\mathsf{int}}^\pi$ is a subset of $\mathcal{C}_{\mathsf{con}}^\pi$ and we truncate once $|\mathcal{C}_{\mathsf{con}}^\pi| > s$, it is not guaranteed that the size increase of $\mathcal{C}_{\mathsf{con}}^\pi$ is smooth that we have a reasonable *upper* bound upon truncation.

To circumvent this issue, we introduce a pruning process on $\mathcal{C}_{\star\text{-int}}^\pi \cup \mathcal{C}_{\mathsf{int}}^\pi$ to obtain the actual witness of size not to large while maintaining some key properties useful later.

**Lemma 7.8.** *Let $\sigma \in \mathcal{N}_{\mathsf{rec}}$ and $\pi \in \mathcal{T}_\sigma$. There exist $\overline{\mathcal{C}}_{\star\text{-int}}^\pi \subseteq \mathcal{C}_{\star\text{-int}}^\pi$ and $\overline{\mathcal{C}}_{\mathsf{int}}^\pi \subseteq \mathcal{C}_{\mathsf{int}}^\pi$ such that the following holds:*

1. *If $|\mathcal{C}_{\star\text{-int}}^\pi \cup \mathcal{C}_{\mathsf{int}}^\pi| \leq n/2^{4k/\log(k)}$, then $\overline{\mathcal{C}}_{\star\text{-int}}^\pi = \mathcal{C}_{\star\text{-int}}^\pi$ and $\overline{\mathcal{C}}_{\mathsf{int}}^\pi = \mathcal{C}_{\mathsf{int}}^\pi$.*
   *Otherwise we have $n/2^{5k/\log(k)} \leq |\overline{\mathcal{C}}_{\star\text{-int}}^\pi \cup \overline{\mathcal{C}}_{\mathsf{int}}^\pi| \leq n/2^{4k/\log(k)}$.*

2. *$\overline{\mathcal{C}}_{\star\text{-int}}^\pi \cup \overline{\mathcal{C}}_{\mathsf{int}}^\pi$ is connected in $G_\Phi$, and $\overline{\mathcal{C}}_{\star\text{-int}}^\pi$ covers at least $|\overline{\mathcal{C}}_{\star\text{-int}}^\pi|$ many $\star$'s.*

3. *For any $C \in \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{bad}}^\pi$ and $v \in \mathsf{vbl}(C)$ with $\pi(v) = \star\!\!\!\!\!\star$, there exists some $C' \in \overline{\mathcal{C}}_{\mathsf{int}}^\pi$ such that $v \in \mathsf{vbl}(C')$ and $C' \in \mathcal{C}_{\mathsf{frozen}}^\pi \cup \mathcal{C}_{\mathsf{sep}}$.*

4. *For any $C \in \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{bad}}^\pi$ and $v \in \mathsf{vbl}(C)$ with $\pi(v) = \star$, there exists some $C' \in \overline{\mathcal{C}}_{\star\text{-int}}^\pi$ such that $v \in \mathsf{vbl}(C')$.*

*Proof.* We start with $\overline{\mathcal{C}}_{\star\text{-int}}^\pi = \mathcal{C}_{\star\text{-int}}^\pi, \overline{\mathcal{C}}_{\mathsf{int}}^\pi = \mathcal{C}_{\mathsf{int}}^\pi$ then perform pruning iteratively. At the beginning, Items 2 and 4 follow from Corollary 7.6, and Item 3 holds due to the definition of $\mathcal{C}_{\mathsf{int}}^\pi, \mathcal{C}_{\mathsf{bad}}^\pi$ and Fact 6.21.

If $|\overline{\mathcal{C}}_{\star\text{-int}}^\pi \cup \overline{\mathcal{C}}_{\mathsf{int}}^\pi| > n/2^{4k/\log(k)}$, then we have the following pruning cases:

- If there exists $\overline{C} \in \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{bad}}^\pi$, then let $\mathcal{S}_1, \mathcal{S}_2, \ldots, \mathcal{S}_t$ be the maximal connected components of $\overline{\mathcal{C}}_{\star\text{-int}}^\pi \cup \overline{\mathcal{C}}_{\mathsf{int}}^\pi$ in $G_\Phi$ after removing $\overline{C}$. Assume that $\mathcal{S}_1$ has the maximal size. Then we update

$$\overline{\mathcal{C}}_{\star\text{-int}}^\pi \leftarrow \mathcal{S}_1 \cap \overline{\mathcal{C}}_{\star\text{-int}}^\pi \quad \text{and} \quad \overline{\mathcal{C}}_{\mathsf{int}}^\pi \leftarrow \mathcal{S}_1 \cap \overline{\mathcal{C}}_{\mathsf{int}}^\pi.$$

- Otherwise $\overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{bad}}^\pi = \emptyset$. Then let $\overline{C} \in \overline{\mathcal{C}}_{\star\text{-int}}^\pi \cup \overline{\mathcal{C}}_{\mathsf{int}}^\pi$ be arbitrary such that removing it does not disconnect $\overline{\mathcal{C}}_{\star\text{-int}}^\pi \cup \overline{\mathcal{C}}_{\mathsf{int}}^\pi$ in $G_\Phi$, and we update

$$\overline{\mathcal{C}}_{\star\text{-int}}^\pi \leftarrow \overline{\mathcal{C}}_{\star\text{-int}}^\pi \setminus \{\overline{C}\} \quad \text{and} \quad \overline{\mathcal{C}}_{\mathsf{int}}^\pi \leftarrow \overline{\mathcal{C}}_{\star\text{-int}}^\pi \setminus \{\overline{C}\}.$$

Now we verify the conditions. The connectivity is trivially preserved, and the number of $\star$'s is always lower bounded by $|\overline{\mathcal{C}}_{\star\text{-int}}^\pi|$ due to $\overline{\mathcal{C}}_{\star\text{-int}}^\pi \subseteq \mathcal{C}_{\star\text{-int}}^\pi$ and Corollary 7.6. Thus Item 2 holds.

Items 3 and 4 is trivial for the second pruning case since $\overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{bad}}^\pi = \emptyset$ there. For the first pruning case, notice that $\mathcal{S}_1, \ldots, \mathcal{S}_t$ are disjoint from each other. Upon the update, for any $C \in \mathcal{S}_1 \cap \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{bad}}^\pi$ and $v \in \mathsf{vbl}(C)$ with $\pi(v) = \star\!\!\!\!\!\star$, its previous witness $C' \in \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap (\mathcal{C}_{\mathsf{frozen}}^\pi \cup \mathcal{C}_{\mathsf{sep}})$ is different from $\overline{C}$ and is connected to $C$ in $G_\Phi$. Thus $C' \in \mathcal{S}_1$ comes along and Item 3 holds. Similar argument holds for Item 4.

Finally we prove Item 1 when the iterative pruning stops. Note that each time we start with size larger than $n/2^{4k/\log(k)}$ and fall into one of the two pruning cases. The second case decreases the size by one, and thus if we stop afterwards, we have $|\overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \cup \overline{\mathcal{C}}^{\pi}_{\text{int}}| > n/2^{4k/\log(k)} - 1 > n/2^{5k/\log(k)}$. The first case removes a clause $\overline{C}$ and decompose the component into $t$ disjoint parts. Since $\overline{C}$ contains at most $k$ literals and the $t$ parts are connected by $\overline{C}$, we know $t \leq k$. Thus if we stop after this case, by averaging argument we have

$$|\overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \cup \overline{\mathcal{C}}^{\pi}_{\text{int}}| \geq \frac{1}{t}\sum_{i=1}^{t}|\mathcal{S}_i| \geq \frac{n/2^{4k\log(k)/k} - 1}{k} \geq n/2^{5k/\log(k)}. \qquad \square$$

Define $\mathcal{W}^{\pi} = \overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \cup \overline{\mathcal{C}}^{\pi}_{\text{int}}$ as our witness for $\pi$. Items 3 and 4 of Lemma 7.8 show that the unassigned variables in bad clauses of $\mathcal{W}^{\pi}$ are also contained as frozen ones (i.e., in $\mathcal{W}^{\pi} \cap \mathcal{C}^{\pi}_{\text{frozen}}$), or as separators (i.e., in $\mathcal{W}^{\pi} \cap \mathcal{C}_{\text{sep}}$), or by some clause in $\overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \subseteq \mathcal{W}^{\pi}$. This allows us to leverage the structural properties of $\Phi$ to show the following "saturation" result, which intuitively says that most clauses in the witness are the frozen ones or contain $\star$'s.

**Lemma 7.9.** *Let $\sigma \in \mathcal{N}_{\text{rec}}$ and $\pi \in \mathcal{T}_{\sigma}$. If $|\mathcal{W}^{\pi}| \geq \log(n)$, then*

$$|\overline{\mathcal{C}}^{\pi}_{\star\text{-int}}| + |\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}^{\pi}_{\text{frozen}}| \geq (1 - 5\eta) \cdot |\mathcal{W}^{\pi}|.$$

*Proof.* Since $\overline{\mathcal{C}}^{\pi}_{\text{int}} \subseteq \mathcal{C}^{\pi}_{\text{int}}$ only contains clauses from $\mathcal{C}^{\pi}_{\text{frozen}}, \mathcal{C}^{\pi}_{\text{bad}}, \mathcal{C}_{\text{sep}}$ which are disjoint, we expand

$$|\overline{\mathcal{C}}^{\pi}_{\text{int}}| = |\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}^{\pi}_{\text{frozen}}| + |\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}^{\pi}_{\text{bad}}| + |\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}_{\text{sep}}|. \tag{8}$$

By Corollary 4.7 and assuming $|\mathcal{W}^{\pi}| \geq \log(n)$, we have

$$|\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}_{\text{sep}}| \leq |\mathcal{W}^{\pi} \cap \mathcal{C}_{\text{sep}}| \leq \frac{1+\eta}{k} \cdot |\mathcal{W}^{\pi}| \leq \eta|\mathcal{W}^{\pi}|, \tag{9}$$

where we use the fact that $\eta = 15\log(k)/k \geq 1/(k-1)$.

Define $\mathcal{C}_1 = \overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \cup (\overline{\mathcal{C}}^{\pi}_{\text{int}} \setminus \mathcal{C}^{\pi}_{\text{bad}})$ and $\mathcal{C}_2 = \overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}^{\pi}_{\text{bad}}$. Let $\mathcal{V}_1 = \bigcup_{C \in \mathcal{C}_1} \mathsf{vbl}(C)$ and $\mathcal{V}_2 = \bigcup_{C \in \mathcal{C}_2} \mathsf{vbl}(C)$. Then for any $v \in \mathcal{V}_2$, we have the following cases:

- If $\pi(v) = \star$, by Item 3 of Lemma 7.8, $v$ is covered in $\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap (\mathcal{C}^{\pi}_{\text{frozen}} \cup \mathcal{C}_{\text{sep}})$ and thus is in $\mathcal{V}_1$.
- If $\pi(v) = \bigstar$, by Item 4 of Lemma 7.8, $v$ is covered in $\overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \subseteq \mathcal{C}_1$ and thus is in $\mathcal{V}_1$.
- Otherwise $\pi(v) \in \{0,1\}$. Since $C(\pi) \neq \mathsf{True}$ and $\pi$ satisfies Assumption 6.2 by Fact 6.21, the number of options for $v$ is

$$|\mathsf{vbl}(C) \setminus \Lambda(\pi)| = |\mathsf{vbl}(C)| - |\mathsf{vbl}(C) \cap \Lambda(\pi)| \leq |\mathsf{vbl}(C)| - |\mathsf{vbl}(C) \cap \Lambda(\pi) \setminus \mathcal{V}_{\text{sep}}| \leq k - (2/3 - 2\eta)k.$$

Thus

$$|\mathcal{V}_1 \cup \mathcal{V}_2| = |\mathcal{V}_1| + |\mathcal{V}_2 \setminus \mathcal{V}_1| \leq k|\mathcal{C}_1| + (1/3 + 2\eta)k \cdot |\mathcal{C}_2|. \tag{10}$$

On the other hand, since $\overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \cap \overline{\mathcal{C}}^{\pi}_{\text{int}} = \emptyset$ and $\mathcal{W}^{\pi} = \overline{\mathcal{C}}^{\pi}_{\star\text{-int}} \cup \overline{\mathcal{C}}^{\pi}_{\text{int}}$, we have $\mathcal{C}_1 \cap \mathcal{C}_2 = \emptyset$ and $\mathcal{C}_1 \cup \mathcal{C}_2 = \mathcal{W}^{\pi}$ of size at most $n/2^{4k/\log(k)}$ by Item 1 of Lemma 7.8. Thus by Item 2 of Proposition 3.3, we have

$$|\mathcal{V}_1 \cup \mathcal{V}_2| \geq \frac{k|\mathcal{C}_1 \cup \mathcal{C}_2|}{1+\eta} = \frac{k|\mathcal{C}_1| + k|\mathcal{C}_2|}{1+\eta}. \tag{11}$$

Combining (10) and (11), we have

$$|\overline{\mathcal{C}}^{\pi}_{\text{int}} \cap \mathcal{C}^{\pi}_{\text{bad}}| = |\mathcal{C}_2| \leq \frac{\eta|\mathcal{C}_1|}{\frac{2}{3} - \frac{7\eta}{3} - 2\eta^2} \leq 4\eta|\mathcal{C}_1| \leq 4\eta|\mathcal{W}^{\pi}|, \tag{12}$$

34

where we use the fact that $\eta \leq 1/9$. Finally we obtain

$$
\begin{aligned}
|\overline{\mathcal{C}}^\pi_{\star\text{-int}}| + |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}| &= |\overline{\mathcal{C}}^\pi_{\star\text{-int}}| + |\overline{\mathcal{C}}^\pi_{\text{int}}| - |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{bad}}| - |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}_{\text{sep}}| && \text{(by (8))} \\
&= |\mathcal{W}^\pi| - |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{bad}}| - |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}_{\text{sep}}| \\
&\geq |\mathcal{W}^\pi| \cdot (1 - 5\eta) && \text{(by (9) and (12))}
\end{aligned}
$$

as desired. $\qquad\square$

The clauses in $\overline{\mathcal{C}}^\pi_{\star\text{-int}}$ contribute at least $|\overline{\mathcal{C}}^\pi_{\star\text{-int}}|$ $\star$'s. To complement, we show that the clauses in $\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}$ contribute almost maximal amount of variables that are fixed towards the unsatisfying direction. Indeed, each clause can access at most $(1/3 + 2\eta)k$ variables due to Assumption 6.2, and clauses can overlap on many variables. Nevertheless, we use the structural properties of $\Phi$ to show that the clauses in $\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}$ will achieve this extremal ratio.

**Lemma 7.10.** *Let $\sigma \in \mathcal{N}_{\text{rec}}$ and $\pi \in \mathcal{T}_\sigma$. Let*

$$
\mathcal{V}' = \bigcup_{C \in \overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}} \{v \in \text{vbl}(C) \mid \pi(v) \in \{0, 1, \star\}\}
$$

*be the set of accessed variables contained in $\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}$. Then $|\mathcal{V}'| \geq |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}| \cdot (1 - 4\eta)k/3$.*

*Proof.* The number of accessed variables in $C \in \mathcal{C}^\pi_{\text{frozen}}$ is

$$
\begin{aligned}
|\text{vbl}(C) \setminus \{v \in \text{vbl}(C) \mid \pi(v) = \star\}| &= |\text{vbl}(C) \setminus \{v \in \text{vbl}(C) \cap \Lambda(\pi) \mid \pi(v) = \star\}| \\
&\geq |\text{vbl}(C)| - |\text{vbl}(C) \cap \mathcal{V}_{\text{sep}}| - |\text{vbl}(C) \cap \Lambda(\pi) \setminus \mathcal{V}_{\text{sep}}| \\
&\geq (k - 2) - 2\eta k - (1 + (2/3 - 2\eta)k) \\
&= k/3 - 3,
\end{aligned}
$$

where we use Proposition 3.2, Algorithm 1, and the definition of $\mathcal{C}^\pi_{\text{frozen}}$ for the third line. Thus $|\text{vbl}(C) \cap \mathcal{V}'| \geq k/3 - 3 \geq (1 - \eta)k/3$.

Note that $\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}} \subseteq \mathcal{W}^\pi$ has size at most $n/2^{4k/\log(k)}$ by Item 1 of Lemma 7.8. By Proposition 3.4 with $b = (1 - \eta)/3$, we have

$$
|\mathcal{V}'| \geq |\mathcal{C}^{\sigma'}_{\text{int}} \cap \mathcal{C}^{\sigma'}_{\text{frozen}}| \cdot (1 - 4\eta)k/3. \qquad\square
$$

**Truncation before the Final Rejection Sampling**

Now we turn to the sampling truncated nodes $\pi \in \mathcal{N}_{\text{samp-trunc}}$, which corresponds to partial assignments $\pi \in \mathcal{N}_{\text{trunc}}$ that do not contain any $\star$. In this case we truncate if some connected component in $\Phi^\pi$ has more than $s$ clauses.

Let $\Phi' = (\mathcal{V}', \mathcal{C}')$ be the maximal connected component in $\Phi^\pi$ of size $|\mathcal{C}'| > s$. To keep notation consistent, we start with $\mathcal{C}^\pi_{\star\text{-int}}$ and $\mathcal{C}^\pi_{\text{int}}$. Here, the construction is simple: We set $\mathcal{C}^\pi_{\star\text{-int}} = \emptyset$ and $\mathcal{C}^\pi_{\text{int}} = \mathcal{C}'$. Now, comparing with Lemma 7.7, we now have a simpler and better lower bound:

$$
|\mathcal{C}^\pi_{\star\text{-int}} \cup \mathcal{C}^\pi_{\text{int}}| = |\mathcal{C}'| \geq s.
$$

To deal with the same trouble of $|\mathcal{C}^\pi_{\star\text{-int}} \cup \mathcal{C}^\pi_{\text{int}}|$ exceeding the threshold for structural properties, we perform the pruning in Lemma 7.8.

**Claim 7.11.** Lemma 7.8 works for $\pi \in \mathcal{N}_{\text{samp-trunc}}$ as well.

*Proof.* We only need to verify Item 3 of Lemma 7.8 for the starting case $\overline{\mathcal{C}}^\pi_{\star\text{-int}} = \mathcal{C}^\pi_{\star\text{-int}} = \emptyset, \overline{\mathcal{C}}^\pi_{\text{int}} = \mathcal{C}^\pi_{\text{int}} = \mathcal{C}'$, and the rest follows the proof of Lemma 7.8 identically.

Let $C \in \mathcal{C}^\pi_{\text{bad}}$ and $v \in \mathsf{vbl}(C)$ with $\pi(v) = \star$ be arbitrary. According to Definition 6.19, we have $\mathcal{V}^\pi_{\text{alive}} = \emptyset$ and thus $v \notin \mathcal{V}^\pi_{\text{alive}}$. Recall the definition of $\mathcal{V}^\pi_{\text{alive}}$, and we have the following cases:

- If $v \in \mathcal{V}_{\text{sep}}$, then there exists $C' \in \mathcal{C}_{\text{sep}}$ such that $v \in \mathsf{vbl}(C')$ as well. Then $C' \in \mathcal{C}'$ since $\mathcal{C}'$ is maximally connected and $C'$ is not satisfied by $\pi$ due to Fact 6.21 and Assumption 6.2.

- Otherwise, there exists $C' \in \mathcal{C} \setminus \mathcal{C}_{\text{sep}}$ such that $|\mathsf{vbl}(C') \cap \Lambda(\pi) \setminus (\mathcal{V}_{\text{sep}} \cup \{v\})| < (2/3 - 2\eta)k$ and $C'(\pi) \neq \mathsf{True}$. If $v \notin \mathsf{vbl}(C')$, then $|\mathsf{vbl}(C') \cap \Lambda(\pi) \setminus \mathcal{V}_{\text{sep}}| < (2/3 - 2\eta)k$ and thus violating Assumption 6.2 and Fact 6.21. Therefore $v \in \mathsf{vbl}(C')$ and $|\mathsf{vbl}(C') \cap \Lambda(\pi) \setminus \mathcal{V}_{\text{sep}}| < 1 + (2/3 - 2\eta)k$. This means $C' \in \mathcal{C}^\pi_{\text{frozen}}$ and $C'$ is connected to $C$ in $G_\Phi$, which implies $C' \in \mathcal{C}'$ as $\mathcal{C}'$ is maximally connected. $\square$

After the pruning, we define our witness $\mathcal{W}^\pi = \overline{\mathcal{C}}^\pi_{\star\text{-int}} \cup \overline{\mathcal{C}}^\pi_{\text{int}}$ analogously. Then Lemma 7.9 and Lemma 7.10 can be proved identically due to Claim 7.11.

### Summarizing Properties of the Witness

Finally we summarize the properties of the witness.

**Corollary 7.12.** *Assume $6k^4\alpha \log(n) < s \leq n/2^{5k/\log(k)}$. Let $\pi \in \mathcal{N}_{\text{trunc}}$. Then we have witness $\mathcal{W}^\pi = \overline{\mathcal{C}}^\pi_{\star\text{-int}} \cup \overline{\mathcal{C}}^\pi_{\text{int}} \subseteq \mathcal{C}$ such that the following holds:*

1. *$s/(6k^4\alpha) \leq |\mathcal{W}^\pi| \leq n/2^{4k/\log(k)}$, $\mathcal{W}^\pi$ is connected in $G_\Phi$, and $\overline{\mathcal{C}}^\pi_{\star\text{-int}}, \overline{\mathcal{C}}^\pi_{\text{int}}$ are disjoint.*
2. *$\overline{\mathcal{C}}^\pi_{\star\text{-int}}$ covers at least $|\overline{\mathcal{C}}^\pi_{\star\text{-int}}|$ many $\star$'s in $\pi$.*
3. *$|\overline{\mathcal{C}}^\pi_{\star\text{-int}}| + |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}| \geq (1 - 5\eta) \cdot |\mathcal{W}^\pi|$.*
4. *$\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}$ accesses at least $|\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}| \cdot (1 - 4\eta)k/3$ distinct variables in $\pi$.*

*Proof.* We verify for $\pi \in \mathcal{N}_{\text{rec-trunc}}$ and the argument for $\mathcal{N}_{\text{samp-trunc}}$ is similar due to the discussion above. By Lemma 7.7, we get $|\mathcal{C}^\pi_{\star\text{-int}} \cup \mathcal{C}^\pi_{\text{int}}| \geq s/(6k^4\alpha)$. Then we perform the pruning in Lemma 7.8. Since $s \leq n/2^{5k/\log(k)}$ and $\alpha \geq 1/k^3$, the obtained witness $\mathcal{W}^\pi = \overline{\mathcal{C}}^\pi_{\star\text{-int}} \cup \overline{\mathcal{C}}^\pi_{\text{int}}$ has size at least $s/(6k^4\alpha)$, which is at least $\log(n)$ as $s > 6k^4\alpha \log(n)$. The other properties follow directly from Lemma 7.8, Lemma 7.9, and Lemma 7.10. $\square$

## 7.3 Refutation of Witnesses

We now show that the number of possible truncation witnesses is small and the algorithm visits any one of them in small probability. These two combined establishes Lemma 7.2 by a union bound.

To better describe the witness, we will provide side information on $\mathcal{W}^\pi$ via the following augmentation. For technical issue, we need to provide the location $z$ for the first generated $\star$ of $\pi$ in addition to $\mathcal{W}^\pi$. This $\star$ may not be covered in $\mathcal{W}^\pi$ due to the pruning process Lemma 7.8.

**Definition 7.13** (Witness Augmentation). For $\pi \in \mathcal{N}_{\text{trunc}}$, we augment $\mathcal{W}^\pi = \overline{\mathcal{C}}^\pi_{\star\text{-int}} \cup \overline{\mathcal{C}}^\pi_{\text{int}}$ to $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$ as follows:

- $\ell = |\mathcal{W}^\pi|$, $q = |\overline{\mathcal{C}}^\pi_{\text{int}} \cap \mathcal{C}^\pi_{\text{frozen}}|$, and $r$ equals the number of $\star$'s contained in $\overline{\mathcal{C}}^\pi_{\star\text{-int}}$.
- $\mathcal{Q} = \overline{\mathcal{C}}^\pi_{\text{int}}$, $\mathcal{R} = \overline{\mathcal{C}}^\pi_{\star\text{-int}}$, and $f$ indicates the locations of the $r$ $\star$'s in $\overline{\mathcal{C}}^\pi_{\star\text{-int}}$.
- $z$ is the first generated[11] $\star$ of $\pi$ in $\mathcal{T}_{\text{sim}}$ (and set $z = \bot$ if $\pi$ has no $\star$).

---

[11]Formally, if $\pi \in \mathcal{T}_\sigma$ for some $\sigma \in \mathcal{N}_{\text{rec}}$, then $z$ is the unique $\star$ in $\sigma$.

By Corollary 7.12, a large witness is guaranteed to exist if we set $s$ suitably large.

**Corollary 7.14.** *If $6k^4\alpha\log(n) < s \le n/2^{5k/\log(k)}$, then $\ell \ge s/(6k^4\alpha)$ in any witness augmentation.*

In the reverse direction, we can count the number of possible witness augmentations satisfying properties in Corollary 7.12.

**Lemma 7.15.** *Assume $6k^4\alpha\log(n) < s \le n/2^{5k/\log(k)}$. For any fixed $\ell, q, r$, there are at most $n^7(k^3\alpha)^\ell k^{2r}$ possible $\mathcal{Q}, \mathcal{R}, f, z$ such that $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$ is an augmentation for a witness satisfying properties in Corollary 7.12. Moreover, $\ell \ge s/(6k^4\alpha)$ and $q + r \ge (1 - 5\eta)\ell$.*

*Proof.* Assume $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$ is the augmentation for $\mathcal{W}^\pi, \pi \in \mathcal{N}_{\mathsf{trunc}}$. Then $|\mathcal{Q} \cup \mathcal{R}| = |\mathcal{W}^\pi| = \ell$. Since $\mathcal{W}^\pi$ is connected in $G_\Phi$ by Item 1 of Corollary 7.12, $\mathcal{Q} \cup \mathcal{R}$ has at most $m \cdot \alpha^2 n^4 (\mathbf{e}k^2\alpha)^\ell$ possibilities by Proposition 3.5. Then we enumerate $\mathcal{R}$. By Item 2 of Corollary 7.12, $|\mathcal{R}| \le r$, and thus $\mathcal{R}$ has $\binom{\ell}{\le r}$ possibilities given $\mathcal{Q} \cup \mathcal{R}$. Now we list all possible $f$. Since $\mathcal{R}$ is fixed and contains at most $k|\mathcal{R}| \le kr$ distinct variables, we know that $f$ has at most $\binom{kr}{r}$ probabilities. Finally $z \in \mathcal{V} \cup \{\bot\}$ has $n + 1$ options. In all, the total count is upper bounded by

$$\alpha^2 n^4 m (\mathbf{e}k^2\alpha)^\ell \cdot \binom{\ell}{\le r} \cdot \binom{kr}{r} \cdot (n+1) \le n^7(k^3\alpha)^\ell k^{2r},$$

where we use the fact that $m = \alpha n$, $n \ge 2^{\Omega(k)}$, $\binom{\ell}{\le r} \le 2^\ell$, and $\binom{kr}{r} \le (\mathbf{e}k)^r \le k^{2r}$. The "moreover" part follows directly from Items 1 to 3 of Corollary 7.12. $\qquad\square$

Now we bound the probability of encountering any witness augmentation. The idea here is that the witness augmentation determines the unsatisfied clauses $\mathcal{Q}$ and the clauses $\mathcal{R}$ containing $\bigstar$'s. Then as we keep this in mind and simulate the execution of the algorithm using the simulation tree $\mathcal{T}_{\mathsf{sim}}$, whenever we need to fix a variable, we know it is fixed towards the unsatisfying direction if it appears in $\mathcal{Q}$, or it is fixed to $\bigstar$ if it is appears in $\mathcal{R}$ and is indicated so by side information $f$. For both cases, we have good probability bound by the edge values $\rho$ in $\mathcal{T}_{\mathsf{sim}}$.

**Lemma 7.16.** *Assume $6k^4\alpha\log(n) < s \le n/2^{5k/\log(k)}$. Let $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$ be a witness augmentation. Then*

$$\sum_{\substack{\pi \in \mathcal{N}_{\mathsf{trunc}} \\ \textit{augmented as } (\ell,q,r,\mathcal{Q},\mathcal{R},f,z)}} \rho(\pi) \le (2\delta)^r \left(2k^{20}2^{-k/3}\right)^q (1+\delta)^{d_{\mathsf{rec}}+1}.$$

*Proof.* We first mark edges in $\mathcal{T}_{\mathsf{sim}}$ leading to possible $\pi \in \mathcal{N}_{\mathsf{trunc}}$ augmented as $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$.

Let $\sigma \in \mathcal{T}_{\mathsf{sim}}$ be an internal node. Let $\sigma_0, \sigma_1, \sigma_\bigstar$ be its child nodes which extend $\sigma$ by fixing variable $v$ to $0, 1, \bigstar$ respectively. We classify $\sigma$ into one of the following types and mark its outgoing edges accordingly:

(i) If $v$ is identified by $f$ as a $\bigstar$ in $\mathcal{R}$, then we say $\sigma$ is T1 and mark the edge $\sigma \to \sigma_\bigstar$.

(ii) Else if $v \in \mathsf{vbl}(C)$ for some clause $C \in \mathcal{Q}$, then we say $\sigma$ is T2 and mark edges $\sigma \to \sigma_\bigstar, \sigma \to \sigma_b$, where $b \in \{0, 1\}$ is the unique value that does not satisfy $C$ if assigned to $v$.[12]

(iii) Else if $z = \bot$, then we say $\sigma$ is T3 and mark edges $\sigma \to \sigma_0, \sigma \to \sigma_1$.

(iv) Else if $v \ne z$ and $\sigma(z) = \bigwhitestar$, then we say $\sigma$ is T4 and mark edges $\sigma \to \sigma_0, \sigma \to \sigma_1$.

---

[12]Pedantically, if $v$ appears in $C$ as $v$, then $b = 0$; otherwise $v$ appears in $C$ as $\neg v$, then $b = 1$.

(v) Else, we say $\sigma$ is T5 and mark all edges $\sigma \to \sigma_0, \sigma \to \sigma_1, \sigma \to \sigma_\star$.

For correctness, we need to show that we do not miss any truncated leaf node. Assume towards contradiction that $\pi \in \mathcal{N}_{\mathsf{trunc}}$ is missed and $\mathcal{W}^\pi$ is augmented as $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$. Along the path from root to $\pi$, let $\sigma \in \mathcal{T}_{\mathsf{sim}}$ be the last node that the marked edges lead to. Define $\sigma_0, \sigma_1, \sigma_\star, v$ as above, and it means the edge $\sigma \to \sigma_{\pi(v)}$ is missed. Then we have the following case analysis:

- $\sigma$ is T1. This cannot happen since $f$ indicates $\pi(v) = \star$ and the edge is already marked.
- $\sigma$ is T2. Since $\mathcal{Q} = \overline{\mathcal{C}}_{\mathsf{int}}^\pi \subseteq \mathcal{C}_{\mathsf{int}}^\pi$, we have $C(\pi) \neq \mathsf{True}$ by the definition of $\mathcal{C}_{\mathsf{int}}^\pi$. Thus $\pi(v)$ equals $\star$ or the unique $b \in \{0, 1\}$ that does not satisfy $C$ if assigned to $v$. Since both edges $\sigma \to \sigma_\star, \sigma \to \sigma_b$ are already marked, this is a contradiction.
- $\sigma$ is T3. This cannot happen since $\pi(v) \neq \star$ by $z = \bot$, and the edge is already marked.
- $\sigma$ is T4. This means $\pi(v) = \star$ since otherwise the edge $\sigma \to \sigma_{\pi(v)}$ is already marked. By definition, $z$ is the first generated $\star$ of $\pi$. Then due to $v \neq z$ and $\pi(v) = \star$, $z$ is already visited before reaching $\sigma$ and updating $v$ to $\star$. This contradicts $\sigma(v) = \star$.
- $\sigma$ is T5. This cannot happen since all three edges are marked.

Let $\pi \in \mathcal{N}_{\mathsf{trunc}}$ be arbitrary and augmented as $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$. Now we count the number of T1/T2/T5 nodes on the path from the root to $\pi$ in $\mathcal{T}_{\mathsf{sim}}$. It is easy to see that there are exactly $r$ T1 nodes, corresponding to $\star$'s in $\mathcal{R} = \overline{\mathcal{C}}_{\star\text{-int}}^\pi$. The number of T2 nodes is the number of accessed variables in $\mathcal{Q} = \overline{\mathcal{C}}_{\mathsf{int}}^\pi \supseteq \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{frozen}}^\pi$, minus some T1 nodes that are also accessed in $\mathcal{Q}$. Therefore

$$
\begin{aligned}
\#\mathsf{T2} &\geq \#\text{accessed variables in } \overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{frozen}}^\pi - \#\mathsf{T1} \\
&\geq |\overline{\mathcal{C}}_{\mathsf{int}}^\pi \cap \mathcal{C}_{\mathsf{frozen}}^\pi| \cdot (1 - 4\eta)k/3 - r \qquad \text{(by Item 4 of Corollary 7.12)} \\
&= (1 - 4\eta)kq/3 - r.
\end{aligned}
$$

To bound the number of T5 nodes, we observe that it can only appear upon and after $z \neq \bot$ is access on the path. Since $z$ is the first generated $\star$ in $\pi$, this means, by the definition of $\mathcal{T}_{\mathsf{sim}}$, we enter a recursive cost tree after accessing $z$. Hence its number of upper bounded by one plus the depth of this particular recursive cost tree, which is in turn bounded by $1 + d_{\mathsf{rec}}$ by definition.

Let $\mathcal{T}$ be the sub-tree of $\mathcal{T}_{\mathsf{sim}}$ consisting of all paths from the root to truncated leaf nodes $\mathcal{N}_{\mathsf{trunc}}$ augmented as $(\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)$. By the argument above, all edges in $\mathcal{T}$ are marked. In addition, along any root-to-leaf path of $\mathcal{T}$, there are

$$
r \text{ T1 nodes}, \geq (1 - 4\eta)kq/3 - r \text{ T2 nodes}, \leq 1 + d_{\mathsf{rec}} \text{ T5 nodes}, \tag{13}
$$

and the rest are T3/4 nodes. In addition, by the definition of the types and edge weights $\rho$, we have

$$
\text{total weight of the outgoing marked edges of a } \begin{cases} \mathsf{T1} \text{ node} \\ \mathsf{T2} \text{ node} \\ \mathsf{T5} \text{ node} \\ \mathsf{T3/T4} \text{ node} \end{cases} \text{ is } \begin{cases} = \delta, \\ \leq (1 + 3\delta)/2, \\ = 1 + \delta, \\ = 1, \end{cases} \tag{14}
$$

where we use Corollary 6.4 for T2 nodes. As a result, we have

$$
\sum_{\substack{\pi \in \mathcal{N}_{\mathsf{trunc}} \\ \text{augmented as } (\ell, q, r, \mathcal{Q}, \mathcal{R}, f, z)}} \rho(\pi) = \sum_{\text{leaf } \pi \in \mathcal{T}} \rho(\pi) = \sum_{\substack{\text{root-to-leaf path } \mathcal{P} \text{ in } \mathcal{T} \\ \text{edge } e \text{ on } \mathcal{P}}} \rho(e)
$$

38

$$\leq \max_{\mathcal{P}} \delta^{\#\mathsf{T1}\in\mathcal{P}} \left(\frac{1+3\delta}{2}\right)^{\#\mathsf{T2}\in\mathcal{P}} (1+\delta)^{\#\mathsf{T5}\in\mathcal{P}} \qquad \text{(by (14))}$$

$$\leq \delta^r \left(\frac{1+3\delta}{2}\right)^{(1-4\eta)kq/3-r} (1+\delta)^{1+d_{\mathsf{rec}}} \qquad \text{(by (13))}$$

$$\leq (2\delta)^r \left(\left(\frac{1+3\delta}{2}\right)^{k/3} \cdot 2^{4\eta k/3}\right)^q (1+\delta)^{1+d_{\mathsf{rec}}}$$

$$= (2\delta)^r \left(\left(\frac{1+3\delta}{2}\right)^{k/3} \cdot k^{20}\right)^q (1+\delta)^{1+d_{\mathsf{rec}}} \quad \text{(since } \eta = 15\log(k)/k)$$

$$\leq (2\delta)^r \left(2^{-k/3} \cdot 2 \cdot k^{20}\right)^q (1+\delta)^{1+d_{\mathsf{rec}}},$$

where we use the fact $\delta = \xi/(k^{40}\alpha) \leq 1/k^{37}$ by $\alpha \geq 1/k^3$ and $\xi \leq 1$ for the last line. $\qquad \square$

Finally we are ready to prove Lemma 7.2 .

*Proof of Lemma 7.2.* By Lemma 6.22, it suffices to enumerate all possible witness augmentations using Lemma 7.15 and apply Lemma 7.16 for each fixed one:

$$p_{\mathsf{halt}}(\Phi, s) \leq \sum_{\substack{\text{possible } (\ell,q,r,\mathcal{Q},\mathcal{R},f,z)}} \sum_{\substack{\pi\in\mathcal{N}_{\mathsf{trunc}} \\ \text{augmented as } (\ell,q,r,\mathcal{Q},\mathcal{R},f,z)}} \rho(\pi)$$

$$\leq \sum_{\substack{\text{possible } (\ell,q,r,\mathcal{Q},\mathcal{R},f,z)}} (2\delta)^r \left(2k^{20}2^{-k/3}\right)^q (1+\delta)^{d_{\mathsf{rec}}+1} \qquad \text{(by Lemma 7.16)}$$

$$\leq \sum_{\substack{\ell,q,r \\ \ell\geq s/(6k^4\alpha) \\ q+r\geq(1-5\eta)\ell}} n^7(k^3\alpha)^\ell k^{2r} \cdot (2\delta)^r \left(2k^{20}2^{-k/3}\right)^q (1+\delta)^{d_{\mathsf{rec}}+1} \qquad \text{(by Lemma 7.15)}$$

$$\leq n^7(1+\delta)^{d_{\mathsf{rec}}+1} \sum_{\substack{\ell,q,r \\ \ell\geq s/(6k^4\alpha)}} (k^3\alpha)^\ell (2k^2\delta)^{(1-5\eta)\ell-q} \left(2k^{20}2^{-k/3}\right)^q \qquad \text{(since } 2k^2\delta \leq 1)$$

$$= n^7(1+\delta)^{d_{\mathsf{rec}}+1} \sum_{\substack{\ell,q,r \\ \ell\geq s/(6k^4\alpha)}} \left(\frac{2k^5\alpha \cdot \delta}{(2k^2\delta)^{5\eta}}\right)^\ell \left(\frac{k^{18}2^{-k/3}}{\delta}\right)^q$$

$$= n^7(1+\delta)^{d_{\mathsf{rec}}+1} \sum_{\substack{\ell,q,r \\ \ell\geq s/(6k^4\alpha)}} \left(\frac{2k^{-35}\xi}{(2k^{-38}\xi/\alpha)^{5\eta}}\right)^\ell \left(\frac{k^{58}2^{-k/3}\alpha}{\xi}\right)^q \qquad \text{(since } \delta = \xi/(k^{40}\alpha))$$

$$\leq n^7(1+\delta)^{d_{\mathsf{rec}}+1} \sum_{\substack{\ell,q,r \\ \ell\geq s/(6k^4\alpha)}} \left(\frac{2k^{-35}\xi}{(2^{-k/3})^{5\eta}}\right)^\ell k^{8q} \qquad \text{(since } \alpha \leq \xi \cdot 2^{k/3}/k^{50})$$

$$\leq n^7(1+\delta)^{d_{\mathsf{rec}}+1} \sum_{\substack{\ell,q,r \\ \ell\geq s/(6k^4\alpha)}} \left(\frac{2k^{-27}}{k^{-25}}\right)^\ell \qquad \text{(since } q \leq \ell, \xi \leq 1, \text{ and } \eta = 15\log(k)/k)$$

$$\leq n^{10}(1+\delta)^{d_{\mathsf{rec}}+1} \sum_{\substack{\ell \geq s/(6k^4\alpha)}} \left(2k^{-2}\right)^\ell \qquad \text{(since } r \leq n \text{ and } q \leq m = \alpha n)$$

$$\leq n^{10}(1+\delta)^{d_{\mathsf{rec}}+1} \cdot k^{-s/(6k^4\alpha)} \qquad\qquad\qquad \text{(since } k \geq 2^{20})$$

as desired. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

# 8 Putting Everything Together

Now we put everything together and prove Theorem 1.3. As we mentioned before, the final algorithm is a combination of two different ones for different ranges of parameters. The atypical setting refers to the case where $\alpha$ or $\varepsilon$ is very small, and it will be handled by the naive rejection sampling described in Section 5. The typical setting is the case where both $\alpha$ and $\varepsilon$ are reasonably large, then it will be handled by the more sophisticated $\texttt{SolutionSampling}(\Phi, s)$ presented in Section 6.

*Proof of Theorem 1.3.* If $\varepsilon \leq \exp\left\{-n/2^{k/2}\right\}$ or $\alpha \leq 1/k^3$, we run $\texttt{RejectionSampling}(\maltese^{\mathcal{V}}, \mathcal{V})$ for $\widetilde{O}\left((n/\varepsilon)^{1+\xi/k}\right)$ steps. By Markov's inequality and Lemmas 5.6 and 5.7, the probability of $\texttt{RejectionSampling}(\maltese^{\mathcal{V}}, \mathcal{V})$ not terminating within these number of steps is at most $\varepsilon$ when the instance is good, which, by Corollary 3.14, happens with probability $1 - o(1/n)$. In addition, the total variation distance between the output and $\mu$ is guaranteed to be at most $\varepsilon$ as desired.

If $\alpha \geq 1/k^3$ and $\varepsilon \geq \exp\left\{-n/2^{k/2}\right\}$, we run $\texttt{SolutionSampling}(\Phi, s)$ with

$$s = 6k^4\alpha\log(n/\varepsilon) = \widetilde{O}(1)$$

for $\widetilde{O}\left((n/\varepsilon)^{1+\xi/k}/\xi\right)$ steps. Now we prove the correctness assuming the input is a nice instance, which happens with probability $1 - o(1/n)$ by Corollary 3.14. By Lemma 6.23 and Lemma 7.1, the expected runtime of $\texttt{SolutionSampling}(\Phi, s)$ is bounded by

$$\widetilde{O}\left(\frac{n}{\xi} \cdot (1+\delta)^{s\cdot k+1} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right) \leq \widetilde{O}\left(\frac{n \cdot (n/\varepsilon)^{\xi/(2k)}}{\xi} \cdot \exp\left\{\frac{\xi \cdot s}{k^6(\alpha+1)}\right\}\right)$$

$$\leq \widetilde{O}\left(\frac{n \cdot (n/\varepsilon)^{\xi/(2k)}}{\xi} \cdot \exp\left\{\frac{6\xi \cdot \log(n/\varepsilon)}{k^2}\right\}\right)$$

$$\leq \widetilde{O}\left((n/\varepsilon)^{\xi/k} \cdot n/\xi\right),$$

where we use the bound

$$(1+\delta)^{s\cdot k} \leq \mathbf{e}^{s\cdot k\cdot\delta} = \exp\left\{\frac{6\xi \cdot \log(n/\varepsilon)}{k^{35}}\right\} \leq (n/\varepsilon)^{\xi/(2k)} \tag{15}$$

for the first inequality. Thus the probability of not terminating within the prescribed number of steps is at most $\varepsilon/2$ by Markov's inequality. Since $\exp\left\{-n/2^{k/2}\right\} \leq \varepsilon < 1$ and $\alpha \leq 2^{k/3}/k^{50}$, we have

$$6k^4\alpha\log(n) < s \leq 6k^4\alpha \cdot 2n/2^{k/2} \leq n/2^{k/6} \leq n/2^{5k/\log(k)}.$$

By Lemma 6.22 and Lemma 7.2, we have

$$p_{\mathsf{halt}}(\Phi, s) \leq n^{10}(1+\delta)^{s\cdot k+2} \cdot k^{-s/(6k^4\alpha)} \qquad\qquad \text{(by Lemma 7.1)}$$

$$\leq 4n^{11}/\varepsilon \cdot k^{-s/(6k^4\alpha)} \qquad\qquad\qquad \text{(by (15) and } \xi/(2k) \leq 1)$$

$$= 4n^{11}/\varepsilon \cdot k^{-\log(n/\varepsilon)}$$

$$\leq \varepsilon/2. \qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{(since } k \geq 2^{20})$$

Then by Corollary 6.16, the total variation distance of the output distribution and $\mu$ is at most $\varepsilon/2 + p_{\mathsf{halt}}(\Phi, s) \leq \varepsilon$ as desired, where the first $\varepsilon/2$ comes from algorithm not terminating within $\widetilde{O}\left((n/\varepsilon)^{1+\xi/k}/\xi\right)$ steps. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

# Acknowledgement

# References

[ACO08] Dimitris Achlioptas and Amin Coja-Oghlan. Algorithmic Barriers from Phase Transitions. In *Proceedings of the 49th Annual IEEE Symposium on Foundations of Computer Science*, pages 793–802. IEEE, oct 2008. 2

[AJ22] Konrad Anand and Mark Jerrum. Perfect sampling in infinite spin systems via strong spatial mixing. *SIAM Journal on Computing*, 51(4):1280–1295, 2022. 3, 5, 6

[AM02] D. Achlioptas and C. Moore. The Asymptotic Order of the Random $k$-SAT Threshold. In *The 43rd Annual IEEE Symposium on Foundations of Computer Science, 2002. Proceedings.*, pages 779–788. IEEE Comput. Soc, 2002. 1

[AP03] Dimitris Achlioptas and Yuval Peres. The threshold for random $k$-SAT is $2^k(\ln 2 - O(k))$. In *Proceedings of the thirty-fifth ACM symposium on Theory of computing - STOC '03*, page 223, New York, New York, USA, 2003. ACM Press. 1

[AZ08] John Ardelius and Lenka Zdeborová. Exhaustive enumeration unveils clustering and freezing in the random 3-satisfiability problem. *Physical Review E*, 78(4):040101, 2008. 5

[BGG⁺19] Ivona Bezáková, Andreas Galanis, Leslie A. Goldberg, Heng Guo, and Daniel Štefankovič. Approximation via correlation decay when strong spatial mixing fails. *SIAM J. Comput.*, 48(2):279–349, 2019. 3, 4

[BH22] Guy Bresler and Brice Huang. The algorithmic phase transition of random k-sat for low degree polynomials. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 298–309, 2022. 2

[Bón06] Miklós Bóna. *A walk through combinatorics: an introduction to enumeration and graph theory.* World Scientific, 2006. 46

[CF14] Amin Coja-Oghlan and Alan M. Frieze. Analyzing walksat on random formulas. *SIAM J. Comput.*, 43(4):1456–1485, 2014. 7, 11

[CMM23] Zongchen Chen, Nitya Mani, and Ankur Moitra. From algorithms to connectivity and back: finding a giant component in random k-sat. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 3437–3470. SIAM, 2023. 3, 6

[COHH17] A. Coja-Oghlan, A. Haqshenas, and S. Hetterich. `Walksat` Stalls Well Below Satisfiability. *SIAM Journal on Discrete Mathematics*, 31(2):1160–1173, jan 2017. 2

[Coj10] Amin Coja-Oghlan. A better algorithm for random k-sat. *SIAM J. Comput.*, 39(7):2823–2864, 2010. 1, 2, 5

[COP16] Amin Coja-Oghlan and Konstantinos Panagiotou. The asymptotic k-sat threshold. *Advances in Mathematics*, pages 985–1068, 2016. 1, 11

[DHKN21] Shaddin Dughmi, Jason Hartline, Robert D Kleinberg, and Rad Niazadeh. Bernoulli factories and black-box reductions in mechanism design. *Journal of the ACM (JACM)*, 68(2):1–30, 2021. 7, 20

[DSS22] Jian Ding, Allan Sly, and Nike Sun. Proof of the satisfiability conjecture for large $k$. *Annals of Mathematics*, 196(1):1–388, 2022. 1, 2, 5, 11

[EL73] Paul Erdős and László Lovász. Problems and results on 3-chromatic hypergraphs and some related questions. In *Colloquia Mathematica Societatis Janos Bolyai 10. Infinite and Finite Sets, Keszthely (Hungary)*. Citeseer, 1973. 5, 7, 10

[ES91] P Erdos and Joel Spencer. Lopsided lovsz local lemma and latin transversals. *Discrete Applied Mathematics*, 30(151-154):10–1016, 1991. 2

[FGYZ21] Weiming Feng, Heng Guo, Yitong Yin, and Chihao Zhang. Fast sampling and counting $k$-sat solutions in the local lemma regime. *Journal of the ACM (JACM)*, 68(6):1–42, 2021. 2, 3, 5

[FHY21] Weiming Feng, Kun He, and Yitong Yin. Sampling constraint satisfaction solutions in the local lemma regime. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 1565–1578, 2021. 2, 3, 5

[Fri99] Ehud Friedgut. Sharp thresholds of graph properties, and the $k$-sat problem. *J. Amer. Math. Soc.*, 12(4):1017–1054, 1999. With an appendix by Jean Bourgain. 1

[GGGH22] Andreas Galanis, Leslie Ann Goldberg, Heng Guo, and Andrés Herrera-Poyatos. Fast sampling of satisfying assignments from random k-sat. *CoRR*, abs/2206.15308, 2022. 3, 6

[GGGY21] Andreas Galanis, Leslie Ann Goldberg, Heng Guo, and Kuan Yang. Counting solutions to random CNF formulas. *SIAM J. Comput.*, 50(6):1701–1738, 2021. 2, 3, 4, 5, 6, 7, 11, 13, 14, 46

[GST16] Heidi Gebauer, Tibor Szabó, and Gábor Tardos. The local lemma is asymptotically tight for sat. *Journal of the ACM (JACM)*, 63(5):1–32, 2016. 2

[Het16] Samuel Hetterich. Analysing survey propagation guided decimationon random formulas. In *ICALP*, volume 55 of *LIPIcs*, pages 65:1–65:12. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2016. 2

[HSS11] Bernhard Haeupler, Barna Saha, and Aravind Srinivasan. New constructive aspects of the lovász local lemma. *J. ACM*, 58(6):28:1–28:28, 2011. 5, 7, 10

[HSW21] Kun He, Xiaoming Sun, and Kewen Wu. Perfect sampling for (atomic) lovász local lemma. *CoRR*, abs/2107.03932, 2021. 2, 3, 5

[HSZ19] Jonathan Hermon, Allan Sly, and Yumeng Zhang. Rapid mixing of hypergraph independent sets. *Random Struct. Algorithms*, 54(4):730–767, 2019. 4

[Hub16] Mark Huber. Nearly optimal bernoulli factories for linear functions. *Combinatorics, Probability and Computing*, 25(4):577–591, 2016. 7, 20

[HWY22] Kun He, Chunyang Wang, and Yitong Yin. Sampling lovász local lemma for general constraint satisfaction solutions in near-linear time. In *63rd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2022, Denver, CO, USA, October 31 - November 3, 2022*, pages 147–158. IEEE, 2022. 3, 5, 6, 7, 18, 20, 21, 22, 23, 26

[HWY23] Kun He, Chunyang Wang, and Yitong Yin. Deterministic counting lovász local lemma beyond linear programming. In Nikhil Bansal and Viswanath Nagarajan, editors, *Proceedings of the 2023 ACM-SIAM Symposium on Discrete Algorithms, SODA 2023, Florence, Italy, January 22-25, 2023*, pages 3388–3425. SIAM, 2023. 3, 5, 6

[JPV21] Vishesh Jain, Huy Tuan Pham, and Thuy Duong Vuong. On the sampling lovász local lemma for atomic constraint satisfaction problems. *CoRR*, abs/2102.08342, 2021. 2, 3, 5

[KKKS98] Lefteris M. Kirousis, Evangelos Kranakis, Danny Krizanc, and Yannis C. Stamatiou. Approximating the unsatisfiability threshold of random formulas. *Random Structures & Algorithms*, 12(3):253–269, 1998. 1, 11

[MMZ05] M. Mézard, T. Mora, and R. Zecchina. Clustering of solutions in the random satisfiability problem. *Phys. Rev. Lett.*, 94:197205, 2005. 1

[Moi19] Ankur Moitra. Approximate counting, the lovász local lemma, and inference in graphical models. *J. ACM*, 66(2):10:1–10:25, 2019. 2, 3, 5, 6

[MPZ02] M. Mézard, G. Parisi, and R. Zecchina. Analytic and algorithmic solution of random satisfiability problems. *Science*, 297(5582):812–815, 2002. 1

[MS07] Andrea Montanari and Devavrat Shah. Counting good truth assignments of random $k$-SAT formulae. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 2007)*, pages 1255–1264, jul 2007. 2

[MU17] Michael Mitzenmacher and Eli Upfal. *Probability and computing: Randomization and probabilistic techniques in algorithms and data analysis*. Cambridge university press, 2017. 48, 49

[NP05] Şerban Nacu and Yuval Peres. Fast simulation of new coins from old. *The Annals of Applied Probability*, 15(1A):93–115, 2005. 7, 20

[QWZ22] Guoliang Qiu, Yanheng Wang, and Chihao Zhang. A perfect sampler for hypergraph independent sets. In Mikolaj Bojanczyk, Emanuela Merelli, and David P. Woodruff, editors, *49th International Colloquium on Automata, Languages, and Programming, ICALP 2022, July 4-8, 2022, Paris, France*, volume 229 of *LIPIcs*, pages 103:1–103:16. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022. 4

[RS98] Martin Raab and Angelika Steger. "balls into bins" - A simple and tight analysis. In Michael Luby, José D. P. Rolim, and Maria J. Serna, editors, *Randomization and Approximation Techniques in Computer Science, Second International Workshop, RANDOM'98, Barcelona, Spain, October 8-10, 1998, Proceedings*, volume 1518 of *Lecture Notes in Computer Science*, pages 159–170. Springer, 1998. 11

# A    Proofs of the Structural Properties

*Proof of Proposition 3.2.* Assume $k \geq 3$. For each $C \in \mathcal{C}$, let $\mathcal{E}(C)$ be the event that $|\mathsf{vbl}(C)| \leq k-3$ and $1_{\mathcal{E}(C)} \in \{0, 1\}$ be the indicator of $\mathcal{E}(C)$. Then by union bound, we have

$$\mathbb{E}\left[1_{\mathcal{E}(C)}\right] = \mathbf{Pr}\left[\mathcal{E}(C)\right] \leq \binom{n}{k-3}\left(\frac{k-3}{n}\right)^k \leq \left(\frac{en}{k-3}\right)^{k-3}\left(\frac{k-3}{n}\right)^k = \frac{e^{k-3}(k-3)^3}{n^3}$$

$$\leq \frac{1}{\alpha n^{2.5}} = \frac{1}{mn^{1.5}}. \qquad \text{(assume } \alpha \leq 2^k \text{ and } n \geq 2^{\Omega(k)}\text{)}$$

Then by Markov's inequality, we have

$$\mathbf{Pr}\left[\exists \text{ such } \mathcal{E}(C)\right] = \mathbf{Pr}\left[\sum_{C \in \mathcal{C}} 1_{\mathcal{E}(C)} \geq 1\right] \leq \frac{1}{n^{1.5}} = o(1/n). \qquad \square$$

*Proof of Proposition 3.3.* We first prove Item 1. By Proposition 3.2, we have $|\mathsf{vbl}(C)| \geq k-2$ for all $C \in \mathcal{C}$ with probability $1 - o(1/n)$. Given this, we focus on the case $k-2 \leq |\mathcal{V}'| \leq n/2^{k/\log(k)}$.

Let $s$ be an integer that $k-2 \leq s \leq n/2^{k/\log(k)}$. Define $t = \lceil(1+\eta)s/k\rceil$. For any fixed subset $X$ of variables of size $s$ and subset $Y$ of clauses of size $t$, we have

$$\mathbf{Pr}\left[\mathsf{vbl}(C) \subseteq X, \forall C \in Y\right] = \left(\frac{s}{n}\right)^{k \cdot t}.$$

Then by enumerating all possible $Y$, we have

$$\mathbf{Pr}\left[\exists Y, \mathsf{vbl}(C) \subseteq X, \forall C \in Y\right] \leq \binom{m}{t}\left(\frac{s}{n}\right)^{kt} \leq \left(\frac{e\alpha n}{t}\right)^t\left(\frac{s}{n}\right)^{kt} \qquad \text{(since } m = \alpha n\text{)}$$

$$\leq \left(\frac{e\alpha n}{s/k}\right)^t\left(\frac{s}{n}\right)^{kt} = (ek\alpha)^t \cdot \left(\frac{s}{n}\right)^{(k-1)t} \qquad \text{(since } t \geq s/k\text{)}$$

$$\leq 2^{4 \cdot ((1+\eta)s+k)} \cdot \left(\frac{s}{n}\right)^{(k-1)t} \qquad \text{(since } t \leq (1+\eta)s/k+1 \text{ and } \alpha \leq 2^k\text{)}$$

$$\leq 2^{4 \cdot ((1+\eta)s+k)} \cdot \left(\frac{s}{n}\right)^{(1-1/k)(1+\eta)s}. \qquad \text{(since } t \geq (1+\eta)s/k\text{)}$$

Thus by union bound over all possible $X$, we have

$$\mathbf{Pr}\left[\exists \text{ such } X, Y\right] \leq \sum_{s=k-2}^{\lfloor n/2^{k/\log(k)}\rfloor} \binom{n}{s} \cdot 2^{4 \cdot ((1+\eta)s+k)} \cdot \left(\frac{s}{n}\right)^{(1-1/k)(1+\eta)s}$$

$$\leq \sum_{s=k-2}^{\lfloor n/2^{k/\log(k)}\rfloor} \left(\frac{4n}{s}\right)^s \cdot 2^{4 \cdot ((1+\eta)s+k)} \cdot \left(\frac{s}{n}\right)^{(1-1/k)(1+\eta)s}$$

$$= 2^{4k} \sum_{s=k-2}^{\lfloor n/2^{k/\log(k)}\rfloor} \left(2^{4\eta+6}\left(\frac{s}{n}\right)^{\eta-\frac{1}{k}-\frac{\eta}{k}}\right)^s \leq 2^{4k} \sum_{s=k-2}^{\lfloor n/2^{k/\log(k)}\rfloor} \left(2^{6(\eta+1)}\left(\frac{s}{n}\right)^{\eta/2}\right)^s$$

$$\text{(assume } \eta - \tfrac{1}{k} - \tfrac{\eta}{k} \geq \eta/2\text{)}$$

$$\leq 2^{4k} \sum_{s=k-2}^{\lfloor \ln^2 n\rfloor} \left(2^{6(\eta+1)}\left(\frac{\ln^2 n}{n}\right)^{\eta/2}\right)^s + 2^{4k} \sum_{s=\lfloor \ln^2 n\rfloor+1}^{\lfloor n/2^{k/\log(k)}\rfloor} \left(2^{6(\eta+1)}2^{-\frac{\eta k}{2\log(k)}}\right)^s$$

$$\leq 2^{4k} n^{-\frac{\eta(k-2)}{4}} \sum_{s=k-2}^{\lfloor \ln^2 n \rfloor} \left( 2^{6(\eta+1)} \left( \frac{\ln^2 n}{\sqrt{n}} \right)^{\eta/2} \right)^s + 2^{4k} \sum_{s=\lfloor \ln^2 n \rfloor + 1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( 2^{6(\eta+1)} 2^{-\frac{\eta k}{2\log(k)}} \right)^s$$

$$\leq 2^{4k} n^{-2} \sum_{s=k-2}^{\lfloor \ln^2 n \rfloor} 2^{-s} + 2^{4k} \sum_{s=\lfloor \ln^2 n \rfloor + 1}^{\lfloor n/2^{k/\log(k)} \rfloor} 2^{-s}$$

$$\text{(assume } \eta(k-2) \geq 8, \ n \geq 2^{\Omega(1+1/\eta)}, \text{ and } \frac{\eta k}{2\log(k)} \geq 6\eta + 7)$$

$$= o(1/n). \hspace{5cm} \text{(assume } n \geq 2^{\Omega(k)})$$

Finally we note that if $\alpha \leq 2^k$, $k/\log(k) \geq 14(1 + 1/\eta)$, and $n \geq 2^{\Omega(k)}$, then all the assumptions above are satisfied.

Now we turn to Item 2. Fix an arbitrary $\mathcal{C}' \subset \mathcal{C}$ with $|\mathcal{C}'| \leq n/2^{2k/\log(k)}$. Let $\mathcal{V}' = \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)$ which satisfies $|\mathcal{V}'| \leq k|\mathcal{C}'| \leq n/2^{k/\log(k)}$. Then by Item 1, we have

$$|\mathcal{C}'| \leq \left| \{ C \in \mathcal{C} \mid \mathsf{vbl}(C) \subseteq \mathcal{V}' \} \right| \leq (1+\eta)|\mathcal{V}'|/k,$$

which implies $\left| \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C) \right| = |\mathcal{V}'| \geq k|\mathcal{C}'|/(1+\eta)$. $\hspace{1cm}$ □

To prove Proposition 3.4, we will need the following technical lemma.

**Proposition A.1.** *Let $\eta = \eta(k) \in (0,1)$ be a parameter. Assume $\alpha \leq 2^k$, $\frac{k}{\log(k)} \geq \frac{5}{\eta}$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n)$ over the random $\Phi$, the following holds: For every $\mathcal{V}' \subset \mathcal{V}$ with $|\mathcal{V}'| \leq n/2^{k/\log(k)}$, we have*

$$\left| \{ C \in \mathcal{C} \mid |\mathsf{vbl}(C) \cap \mathcal{V}'| \geq \eta k \} \right| \leq k|\mathcal{V}'|.$$

*Proof.* Let $s \leq n/2^{k/\log(k)}$ be an integer. For any fixed subset $X$ of variables of size $s$ and subset $Y$ of clauses of size $ks$, we have

$$\mathbf{Pr}\left[ |\mathsf{vbl}(C) \cap X| \geq \eta k, \forall C \in Y \right] \leq \left( \binom{k}{\lceil \eta k \rceil} \cdot \left( \frac{s}{n} \right)^{\lceil \eta k \rceil} \right)^{ks} \leq \left( 2^k \cdot \left( \frac{s}{n} \right)^{\eta k} \right)^{ks} = \left( \frac{2^{1/\eta} \cdot s}{n} \right)^{\eta k^2 s},$$

where $\binom{k}{\eta k}$ chooses the (first) $\eta k$ locations in $C$ that use variables from $X$. Thus by union bound, we have

$$\mathbf{Pr}\left[ \exists \text{ such } X, Y \right] \leq \sum_{s=1}^{\lfloor n/2^{k/\log(k)} \rfloor} \binom{n}{s} \binom{m}{ks} \cdot \left( \frac{2^{1/\eta} \cdot s}{n} \right)^{\eta k^2 s}$$

$$\leq \sum_{s=1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( \frac{en}{s} \right)^s \cdot \left( \frac{em}{ks} \right)^{ks} \cdot \left( \frac{2^{1/\eta} \cdot s}{n} \right)^{\eta k^2 s}$$

$$= \sum_{s=1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( \frac{e^{k+1} 2^{k^2} \alpha^k s^{\eta k^2 - k - 1}}{k^k n^{\eta k^2 - k - 1}} \right)^s \hspace{2cm} \text{(since } m = \alpha n)$$

$$\leq \sum_{s=1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( \frac{2^{2k^2} s^{\eta k^2 - k - 1}}{n^{\eta k^2 - k - 1}} \right)^s \hspace{2cm} \text{(since } \alpha \leq 2^k \text{ and assume } e^{k+1} \leq k^k)$$

45

$$\leq \sum_{s=1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( \frac{2^4 \cdot s^\eta}{n^\eta} \right)^{k^2 s/2} \qquad \text{(assume } \eta k^2 - k - 1 \geq \eta k^2/2)$$

$$\leq \sum_{s=1}^{\lfloor \ln^2 n \rfloor} \left( \frac{2^4 \cdot \ln^{2\eta} n}{n^\eta} \right)^{k^2 s/2} + \sum_{s=\lfloor \ln^2 n \rfloor + 1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( \frac{2^4}{2^{\eta k/\log(k)}} \right)^{k^2 s/2}$$

$$\leq n^{-\eta k^2/4} \sum_{s=1}^{\lfloor \ln^2 n \rfloor} \left( \frac{2^{22} \cdot \ln^{2\eta} n}{n^{\eta/2}} \right)^{k^2 s/2} + \sum_{s=\lfloor \ln^2 n \rfloor + 1}^{\lfloor n/2^{k/\log(k)} \rfloor} \left( \frac{2^4}{2^{\eta k/\log(k)}} \right)^{k^2 s/2}$$

$$\leq n^{-2} \sum_{s=1}^{\lfloor \ln^2 n \rfloor} 2^{-k^2 s/2} + \sum_{s=\lfloor \ln^2 n \rfloor + 1}^{\lfloor n/2^{k/\log(k)} \rfloor} 2^{-k^2 s/2}$$

$$\text{(assume } \eta k^2 \geq 8, \ n \geq 2^{\Omega(1/\eta)}, \text{ and } \eta k/\log(k) \geq 5)$$

$$= o(1/n).$$

Finally we note that if $k/\log(k) \geq 5/\eta$, $\alpha \leq 2^k$, and $n \geq 2^{\Omega(k)}$, then all the assumptions above are satisfied. $\qquad \square$

Now we proceed to the proof of Proposition 3.4.

*Proof of Proposition 3.4.* We assume $\Phi$ satisfies the properties in Proposition A.1 and Proposition 3.3, which by union bound happens with probability $1 - o(1/n)$. We also assume $b \leq 1$ since otherwise the statement trivially holds.

Fix an arbitrary $\mathcal{V}' \subset \mathcal{V}$ with $|\mathcal{V}'| \leq n/2^{3k/\log(k)}$. Let $\mathcal{C}' = \{C \in \mathcal{C} \mid |\mathsf{vbl}(C) \cap \mathcal{V}'| \geq bk\}$. Then

$$\left| \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C) \setminus \mathcal{V}' \right| \leq (1-b)k \cdot |\mathcal{C}'|.$$

By Proposition A.1, we know $|\mathcal{C}'| \leq k|\mathcal{V}'| \leq n/2^{2k/\log(k)}$. Then by Item 2 of Proposition 3.3, we have

$$\frac{k|\mathcal{C}'|}{1+\eta} \leq \left| \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C) \right| = |\mathcal{V}'| + \left| \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C) \setminus \mathcal{V}' \right| \leq |\mathcal{V}'| + (1-b)k \cdot |\mathcal{C}'|,$$

which implies $|\mathcal{V}'| \geq k|\mathcal{C}'| \cdot \left( \frac{1}{1+\eta} - (1-b) \right) \geq (b-\eta)k \cdot |\mathcal{C}'|$. $\qquad \square$

**Lemma A.2** ([GGGY21, Lemma 8.5]). *For any labeled tree $T$ on a subset of $\mathcal{C}$, the probability that $T$ is a sub-graph of $G_\Phi$ is at most $(k^2/n)^{|V(T)|-1}$ where $V(T)$ is the number of nodes of $T$.*

*Proof of Proposition 3.5.* Let $C \in \mathcal{C}$ be arbitrary and let $U \subseteq \mathcal{C}$ be a size-$\ell$ set of clauses containing $C$. For any fixed labeled spanning tree on $U$, by Lemma A.2 it appears in $G_\Phi$ with probability at most $(k^2/n)^{\ell-1}$. Meanwhile by standard result (See e.g., [Bón06]), there are $\ell^{\ell-2}$ many possible $U$. Thus by union bound, we have

$$\mathbf{Pr}\left[ G_\Phi[U] \text{ is connected} \right] \leq \ell^{\ell-2} (k^2/n)^{\ell-1}.$$

Now let $Z_{\ell,C}$ be the number of connected sets of clauses with size $\ell$ containing $C$. Then

$$\mathbb{E}\left[ Z_{\ell,C} \right] = \sum_{U \subseteq \mathcal{C}: C \in U, |U| = \ell} \mathbf{Pr}\left[ G_\Phi[U] \text{ is connected} \right]$$

46

$$\leq \binom{m-1}{\ell-1} \cdot \ell^{\ell-2} \left(\frac{k^2}{n}\right)^{\ell-1} \leq \left(\frac{\mathbf{e}(m-1)}{\ell-1}\right)^{\ell-1} \cdot \ell^{\ell-2} \left(\frac{k^2}{n}\right)^{\ell-1}$$

$$\leq \left(\frac{\mathbf{e}mk^2 \cdot \ell^{\frac{\ell-2}{\ell-1}}}{n \cdot (\ell-1)}\right)^{\ell-1} \leq (\mathbf{e}k^2\alpha)^{\ell-1}. \qquad (\text{since } \ell^{\ell-2} \leq (\ell-1)^{\ell-1} \text{ and } m = \alpha n)$$

Then by Markov's inequality, we have

$$\mathbf{Pr}\left[Z_{\ell,C} \geq \alpha^2 n^4 (\mathbf{e}k^2\alpha)^{\ell-1}\right] \leq \alpha^{-2} n^{-4} = n^{-2}m^{-2}.$$

Finally, by union bound, we have

$$\mathbf{Pr}\left[\exists \text{ such } Z_{\ell,C} \geq \alpha^2 n^3 (\mathbf{e}k^2\alpha)^{\ell-1}\right] \leq m^2 \cdot n^{-2}m^{-2} = 1/n^2 = o(1/n). \qquad \square$$

*Proof of Proposition 3.6.* Define $\widetilde{\mathcal{V}} = |\{v \in \mathcal{V} \mid v \in \mathcal{V}' \text{ or } v \text{ is adjacent to } \mathcal{V}'\}|$. Let

$$\mathcal{C}' = \left\{C \in \mathcal{C} \mid \mathsf{vbl}(C) \cap \mathcal{V}' \neq \emptyset\right\}.$$

Since $|\widetilde{\mathcal{V}}| \leq k|\mathcal{C}'|$, it suffices to bound $|\mathcal{C}'| \leq 3k^3\alpha \max\left\{|\mathcal{V}'|, \lfloor k\log(n)\rfloor\right\}$.

We first focus on the case $|\mathcal{V}'| \geq \lfloor k\log(n)\rfloor$. Since $H_\Phi[\mathcal{V}']$ is connected, there exists some $\mathcal{C}'' \subseteq \mathcal{C}'$ such that $|\mathcal{V}'|/k \leq |\mathcal{C}''| \leq |\mathcal{V}'|$ and $\mathcal{V}'$ is connected in $H_\Phi$ using $\mathcal{C}''$. In particular, we have

$$|\mathcal{C}''| \geq \lfloor k\log(n)\rfloor/k \geq \log(n) - 1.$$

Let $\widetilde{\mathcal{C}} = \mathcal{C}' \setminus \mathcal{C}''$. Since $k^3\alpha \geq 1$, it suffices to bound $|\mathcal{C}'| \leq 2k^3\alpha|\mathcal{V}'| + |\mathcal{V}'|$. Then plugging in $|\mathcal{C}'| = |\widetilde{\mathcal{C}}| + |\mathcal{C}''|$ and $|\mathcal{C}''| \leq |\mathcal{V}'|$, it suffices to prove $|\widetilde{\mathcal{C}}| \leq 2k^3\alpha|\mathcal{V}'|$. Now for any fixed $\mathcal{C}'', \mathcal{V}', \widetilde{\mathcal{C}}$ satisfying:

- $|\mathcal{C}''| \geq \log(n) - 1$, $|\mathcal{V}'| \geq |\mathcal{C}''|$, $|\widetilde{\mathcal{C}}| \geq 2k^3\alpha|\mathcal{V}'|$, and $\mathcal{C}'' \cap \widetilde{\mathcal{C}} = \emptyset$.
- $G_\Phi[\mathcal{C}'']$ is connected, $\mathcal{V}' \subseteq \bigcup_{C \in \mathcal{C}''} \mathsf{vbl}(C)$, and $\mathsf{vbl}(\widetilde{C}) \cap \mathcal{V}' \neq \emptyset$ holds for all $\widetilde{C} \in \widetilde{\mathcal{C}}$.

Let $s_1 = |\mathcal{C}''|$, $s_2 = |\mathcal{V}'|$, and $s_3 = |\widetilde{\mathcal{C}}|$. We now define the following events:

- $\mathcal{E}(\mathcal{C}'', \mathcal{V}', \widetilde{\mathcal{C}})$ is the event that "$\mathcal{C}'', \mathcal{V}', \widetilde{\mathcal{C}}$ satisfy the conditions above".
- $\mathcal{E}(\mathcal{C}'')$ is the event that "$G_\Phi[\mathcal{C}'']$ is connected".
- $\mathcal{E}(\mathcal{V}', \widetilde{\mathcal{C}})$ is the event that "$\mathsf{vbl}(\widetilde{C}) \cap \mathcal{V}' \neq \emptyset$ holds for all $\widetilde{C} \in \widetilde{\mathcal{C}}$".

By union bounding over all $s_1^{s_1-2}$ labeled spanning trees over $\mathcal{C}''$ and using Lemma A.2, we have

$$\mathbf{Pr}\left[\mathcal{E}(\mathcal{C}'')\right] \leq s_1^{s_1-2}\left(\frac{k^2}{n}\right)^{s_1-1}.$$

Since $\mathcal{C}'' \cap \widetilde{\mathcal{C}} = \emptyset$, by independence we have

$$\mathbf{Pr}\left[\mathcal{E}(\mathcal{V}', \widetilde{\mathcal{C}}) \,\middle|\, \mathcal{E}(\mathcal{C}'')\right] = \mathbf{Pr}\left[\mathcal{E}(\mathcal{V}', \widetilde{\mathcal{C}})\right] \leq \left(k \cdot \frac{s_2}{n}\right)^{s_3}.$$

Hence

$$\mathbf{Pr}\left[\mathcal{E}(\mathcal{C}'', \mathcal{V}', \widetilde{\mathcal{C}})\right] \leq \mathbf{Pr}\left[\mathcal{E}(\mathcal{C}'') \wedge \mathcal{E}(\mathcal{V}', \widetilde{\mathcal{C}})\right] \leq s_1^{s_1-2}\left(\frac{k^2}{n}\right)^{s_1-1}\left(\frac{ks_2}{n}\right)^{s_3}.$$

Thus by union bound, we have

$$\mathbf{Pr}\left[\exists \text{ such } \mathcal{E}(\mathcal{C}'',\mathcal{V}',\widetilde{\mathcal{C}})\right] \leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq s_1}\sum_{s_3\geq 2k^3\alpha\cdot s_2}\binom{m}{s_1}\binom{ks_1}{s_2}\binom{m}{s_3}\cdot s_1^{s_1-2}\left(\frac{k^2}{n}\right)^{s_1-1}\left(\frac{ks_2}{n}\right)^{s_3}$$

$$\left(\binom{ks_1}{s_2} \text{ comes from } \mathcal{V}' \subseteq \bigcup_{C\in\mathcal{C}''}\mathsf{vbl}(C)\right)$$

$$\leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq s_1}\sum_{s_3\geq 2k^3\alpha\cdot s_2}\frac{n}{k^2 s_1^2}\left(\mathbf{e}k^2\alpha\right)^{s_1}\left(\frac{\mathbf{e}ks_1}{s_2}\right)^{s_2}\left(\frac{\mathbf{e}k\alpha s_2}{s_3}\right)^{s_3}$$

$$\text{(since } m=\alpha n)$$

$$\leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq s_1}\frac{n}{k^2 s_1^2}\left(\mathbf{e}k^2\alpha\right)^{s_1}\left(\mathbf{e}k\right)^{s_2}\sum_{s_3\geq 2k^3\alpha\cdot s_2}\left(\frac{\mathbf{e}}{2k^2}\right)^{s_3} \quad \text{(since } s_2\geq s_1)$$

$$\leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq s_1}\frac{2n}{k^2 s_1^2}\left(\mathbf{e}k^2\alpha\right)^{s_1}\left(\mathbf{e}k\right)^{s_2}\left(\frac{\mathbf{e}}{2k^2}\right)^{2k^3\alpha\cdot s_2} \quad \text{(assume } k\geq 2)$$

$$\leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq s_1}\frac{2n}{k^2 s_1^2}\left(\mathbf{e}k^2\alpha\right)^{s_1}\left(\mathbf{e}k\right)^{s_2}\left(\frac{\mathbf{e}}{2k^2}\right)^{k^3\alpha\cdot s_1}\left(\frac{\mathbf{e}}{2k^2}\right)^{k^3\alpha\cdot s_2}$$

$$\text{(since } s_2\geq s_1)$$

$$= \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq s_1}\frac{2n}{k^2 s_1^2}\left(\frac{\mathbf{e}k^2\alpha}{(2k^2/\mathbf{e})^{k^3\alpha}}\right)^{s_1}\left(\frac{\mathbf{e}k}{(2k^2/\mathbf{e})^{k^3\alpha}}\right)^{s_2}$$

$$\leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq 1}\frac{2n\cdot 8^{-s_1-s_2}}{k^2 s_1^2}$$

$$\left(\text{assume } \frac{\mathbf{e}k^2\alpha}{(2k^2/\mathbf{e})^{k^3\alpha}}\leq\frac{1}{8} \text{ and } \frac{\mathbf{e}k}{(2k^2/\mathbf{e})^{k^3\alpha}}\leq\frac{1}{8}\right)$$

$$\leq \sum_{s_1\geq\log(n)-1}\sum_{s_2\geq 1}\frac{2n\cdot 8^{-s_1-s_2}}{(\log(n)-1)^2}\leq\sum_{s_1\geq\log(n)-1}\frac{32n\cdot 8^{-s_1}}{(\log(n)-1)^2}$$

$$\leq \frac{4n}{n^2(\log(n)-1)^2}=o(1/n).$$

Now we analyze the assumptions. Define $t=k^3\alpha$. Then the calculation above demands $k\geq 2$ and

$$\frac{t}{k}\left(\frac{\mathbf{e}}{2k^2}\right)^t\leq\frac{1}{8\mathbf{e}} \quad\text{and}\quad k\left(\frac{\mathbf{e}}{2k^2}\right)^t\leq\frac{1}{8\mathbf{e}}.$$

Thus it suffices to assume $k^3\alpha=t\geq 1$ and $k\geq 30$.

Now we turn to the case $|\mathcal{V}'|<\lfloor k\log(n)\rfloor$. If $|\widetilde{\mathcal{V}}|<\lfloor k\log(n)\rfloor$, then we are done since $\alpha\geq 1/k^3$. Otherwise consider an arbitrary connected $\widehat{\mathcal{V}}\supset\mathcal{V}'$ such that $|\widehat{\mathcal{V}}|=\lfloor k\log(n)\rfloor$. Then by applying the previous argument on $\widehat{\mathcal{V}}$, we have

$$|\widetilde{\mathcal{V}}|\leq\left|\left\{v\in\mathcal{V}\,\middle|\,v\in\widehat{\mathcal{V}}\text{ or }v\text{ is adjacent to }\widehat{\mathcal{V}}\right\}\right|\leq 3k^4\alpha|\widehat{\mathcal{V}}|=3k^4\alpha\cdot\lfloor k\log(n)\rfloor. \qquad\square$$

*Proof of Proposition 3.7.* The degrees of the variables in $\Phi$ distribute as a balls-and-bins experiment with $km$ balls and $n$ bins. Let $D_1,\ldots,D_n\sim\mathsf{Poi}(k\alpha)$ be $n$ independent Poisson random variables with parameter $k\alpha$. Then the degrees of the variables in $\Phi$ has the same distribution as $\{D_1,\ldots,D_n\}$ conditioned on the event $\mathcal{E}$ that $\sum_{i=1}^n D_i=km$ [MU17, Chapter 5.4]. Note that $\sum_{i=1}^n D_i$ is a Poisson random variable with parameter $k\alpha n=km$. Thus

$$\mathbf{Pr}\left[\mathcal{E}\right]=\mathbf{e}^{-km}\cdot\frac{(km)^{km}}{(km)!}\geq\frac{1}{\sqrt{2\pi km}}=\frac{1}{\sqrt{2\pi k\alpha n}}.$$

48

Let $D = 4k\alpha + 6\log(n)$. For any fixed $i \in [n]$, we have

$$\mathbf{Pr}\left[D_i \geq D\right] = \mathbf{Pr}\left[\mathsf{Poi}(k\alpha) \geq D\right] \leq \frac{\mathbf{e}^{-k\alpha}(\mathbf{e}k\alpha)^D}{D^D} \qquad \text{(by [MU17, Theorem 5.4])}$$

$$\leq \mathbf{e}^{-k\alpha}(\mathbf{e}/4)^D \leq \mathbf{e}^{-k\alpha} \cdot 2^{-D/2} \qquad \text{(since } D \geq 4k\alpha)$$

$$\leq \mathbf{e}^{-k\alpha} \cdot n^{-3}. \qquad \text{(since } D \geq 6\log(n))$$

Define $U = \{i \in [n] \mid D_i \geq D\}$. Then

$$\mathbf{Pr}\left[\exists v \in \mathcal{V}, \deg_{\mathcal{C}}(v) \geq D\right] = \mathbf{Pr}\left[|U| \geq 1 \mid \mathcal{E}\right] \leq \frac{\mathbf{Pr}\left[|U| \geq 1\right]}{\mathbf{Pr}\left[\mathcal{E}\right]}$$

$$\leq \sqrt{2\pi k\alpha n} \cdot n \cdot \mathbf{Pr}\left[D_i \geq D\right] \qquad \text{(by Markov's inequality)}$$

$$\leq \sqrt{2\pi k\alpha n} \cdot \mathbf{e}^{-k\alpha} \cdot n^{-2}$$

$$= O(1/n^{1.5}) = o(1/n). \qquad \square$$

*Proof of Proposition 3.8.* The calculation is similar to the proof of Proposition 3.7.

Let $D_1, \ldots, D_n \sim \mathsf{Poi}(k\alpha)$ be $n$ independent Poisson random variables with parameter $k\alpha$. Then the degrees of the variables in $\Phi$ has the same distribution as $\{D_1, \ldots, D_n\}$ conditioned on the event $\mathcal{E}$ that $\sum_{i=1}^{n} D_i = km$. Note that $\sum_{i=1}^{n} D_i$ is a Poisson random variable with parameter $k\alpha n = km$. Thus

$$\mathbf{Pr}\left[\mathcal{E}\right] = \mathbf{e}^{-km} \cdot \frac{(km)^{km}}{(km)!} \geq \frac{1}{\sqrt{2\pi km}} = \frac{1}{\sqrt{2\pi k\alpha n}}.$$

For any fixed $i \in [n]$, we have

$$\mathbf{Pr}\left[D_i \geq D\right] = \mathbf{Pr}\left[\mathsf{Poi}(k\alpha) \geq D\right] \leq \frac{\mathbf{e}^{-k\alpha}(\mathbf{e}k\alpha)^D}{D^D} \qquad \text{(by [MU17, Theorem 5.4])}$$

$$\leq \mathbf{e}^{-k\alpha}(\mathbf{e}/8)^D \leq (\mathbf{e}/8)^D \qquad \text{(assume } D \geq 8k\alpha)$$

$$\leq 2^{-4k-1}. \qquad \text{(assume } D \geq 8k)$$

Define $U = \{i \in [n] \mid D_i \geq D\}$. Then by Chernoff-Hoeffding bound, we have

$$\mathbf{Pr}\left[|U| \geq n/2^{4k}\right] \leq \mathbf{Pr}\left[|U| - \mathbb{E}[|U|] \geq n/2^{4k+1}\right] \leq \mathbf{e}^{-n/2^{4k+1}}.$$

Thus

$$\mathbf{Pr}\left[|\{v \in \mathcal{V} \mid \deg_{\mathcal{C}}(v) \geq D\}| \geq n/2^{4k}\right]$$

$$= \mathbf{Pr}\left[|U| \geq n/2^{4k} \,\Big|\, \mathcal{E}\right] \leq \frac{\mathbf{Pr}\left[|U| \geq n/2^{4k}\right]}{\mathbf{Pr}\left[\mathcal{E}\right]}$$

$$\leq \sqrt{2\pi k\alpha n} \cdot \mathbf{e}^{-n/2^{4k+1}}$$

$$= o(1/n). \qquad \text{(assume } n \geq 2^{\Omega(k)} \text{ and } \alpha \leq 2^k)$$

Finally we note that if $k \geq 2$, $\alpha \leq 2^k$, $D \geq 8k(\alpha + 1)$, and $n \geq 2^{\Omega(k)}$, then all the assumptions above are satisfied. $\square$

*Proof of Proposition 3.9.* Let $\widetilde{\mathcal{V}} = |\{v \in \mathcal{V}' \mid \deg_{\mathcal{C}}(v) \geq D\}|$ and $\widetilde{\mathcal{C}} = \left\{C \in \mathcal{C} \,\Big|\, \mathsf{vbl}(C) \cap \widetilde{\mathcal{V}} \neq \emptyset\right\}$. To lower bound $|\widetilde{\mathcal{C}}|$, we perform a double counting for the size of $\left\{(v, C) \,\Big|\, v \in \widetilde{\mathcal{V}}, C \in \widetilde{\mathcal{C}}\right\}$, which is lower bounded by $D \cdot |\widetilde{\mathcal{V}}|$ and upper bounded by $k \cdot |\widetilde{\mathcal{C}}|$. Therefore we have $|\widetilde{\mathcal{C}}| \geq D|\widetilde{\mathcal{V}}|/k$.

By Proposition 3.8, we have $|\widetilde{\mathcal{V}}| \leq n/2^{4k}$ with probability $1 - o(1/n)$. Since $D \leq 2^{2k}$, we have $|\widetilde{\mathcal{C}}| \leq D|\widetilde{\mathcal{V}}| \leq Dn/2^{4k} \leq n/2^{2k/\log(k)}$. By Item 2 of Proposition 3.3 with $\eta = 1$, we have

$$\left| \bigcup_{C \in \mathcal{C}: \mathsf{vbl}(C) \cap \mathcal{V}' \neq \emptyset} \mathsf{vbl}(C) \right| \geq \left| \bigcup_{C \in \widetilde{\mathcal{C}}} \mathsf{vbl}(C) \right| \geq \frac{k|\widetilde{\mathcal{C}}|}{2} \geq D|\widetilde{\mathcal{V}}|/2$$

with probability $1 - o(1/n)$. On the other hand, by Proposition 3.6 we have

$$\left| \bigcup_{C \in \mathcal{C}: \mathsf{vbl}(C) \cap \mathcal{V}' \neq \emptyset} \mathsf{vbl}(C) \right| \leq 3k^4 \alpha \cdot \max\left\{ |\mathcal{V}'|, k \log(n) \right\} \leq 3k^5 \alpha \cdot |\mathcal{V}'|$$

with probability $1 - o(1/n)$, where we use the bound $\max\left\{ |\mathcal{V}'|, k \log(n) \right\} \leq k|\mathcal{V}'|$ as $|\mathcal{V}'| \geq \log(n)$. Rearranging and assuming $D \geq 6k^7 \alpha$, we have

$$|\widetilde{\mathcal{V}}| \leq |\mathcal{V}'| \cdot \frac{3k^5 \alpha}{D/2} \leq |\mathcal{V}'|/k^2.$$

Finally we note that if $k \geq 2^{10}$, $6k^7(\alpha + 1) \leq D \leq 2^{2k}$, $1/k^3 \leq \alpha \leq 2^k$, and $n \geq 2^{\Omega(k)}$, then all the assumptions used above are satisfied. $\qquad\square$

Proposition 3.10 is a simple union bound of the following lemma.

**Lemma A.3.** *Let $\varepsilon = \varepsilon(k, n)$ be a parameter satisfying $1/n \leq \varepsilon \leq 2^{-2.5k}$. Assume $k \geq 12$, $\alpha \leq 2^k$, and $n \geq 2^{\Omega(k)}$. Then with probability $1 - o(1/n^3)$ over the random $\Phi$, the following holds: Fix an arbitrary $\mathcal{C}' \subseteq \mathcal{C}$ with $|\mathcal{C}'| \leq \varepsilon n$. Let $C_{i_1}, \ldots, C_{i_\ell} \in \mathcal{C} \setminus \mathcal{C}'$ be clauses with distinct indices. For each $s \in [\ell]$, define $\mathcal{V}_s = \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C) \cup \bigcup_{j=1}^{s-1} \mathsf{vbl}(C_{i_j})$. If $|\mathsf{vbl}(C_{i_s}) \cap \mathcal{V}_s| \geq 6$ holds for all $s \in [\ell]$, then $\ell \leq \varepsilon n$.*

*Proof.* Assume $\mathcal{C}'$ and $C_{i_1}, \ldots, C_{i_\ell}$ violates the statement. By discarding redundant clauses from $C_{i_1}, \ldots, C_{i_\ell}$, we assume $\ell = \lfloor \varepsilon n \rfloor + 1$. Now, as long as $|\mathcal{C}'| < \lfloor \varepsilon n \rfloor$ and $\mathcal{C} \setminus \{\mathcal{C}' \cup \{C_{i_1}, \ldots, C_{i_\ell}\}\}$ is not empty, we can enlarge $\mathcal{C}'$ by including new clauses and the statement is still violated. Therefore we assume $|\mathcal{C}'| = \min\{\lfloor \varepsilon n \rfloor, m - \ell\} = \min\{\ell - 1, m - \ell\}$.

Note that the sets $Y = \bigcup_{j=1}^{\ell} \mathsf{vbl}(C_{i_j}) \setminus \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)$ and $\mathcal{C}'$ have the following properties:

- $|Y| = \sum_{s=1}^{\ell} |\mathsf{vbl}(C_{i_s})| - |\mathsf{vbl}(C_{i_s}) \cap \mathcal{V}_s| \leq (k - 6)\ell$.

  This is because each $C_{i_s}$ intersects $\mathcal{V}_s$ with at least 6 variables.

- There exists $\widetilde{\mathcal{C}} \subset \mathcal{C} \setminus \mathcal{C}'$ with $|\widetilde{\mathcal{C}}| = \ell$ such that $\mathsf{vbl}(\widetilde{C}) \subseteq Y \cup \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)$ for all $\widetilde{C} \in \widetilde{\mathcal{C}}$.

  This is because we can pick $\widetilde{\mathcal{C}} = \{C_{i_1}, \ldots, C_{i_\ell}\}$.

Now for any fixed $\mathcal{C}', Y, \widetilde{\mathcal{C}}$ satisfying $|\mathcal{C}'| = \min\{\ell - 1, m - \ell\}$, $|\widetilde{\mathcal{C}}| = \ell$, and $|Y| = t \leq (k - 6)\ell$. We define event $\mathcal{E}(\mathcal{C}', Y, \widetilde{\mathcal{C}})$ to be "$\mathsf{vbl}(\widetilde{C}) \subseteq Y \cup \bigcup_{C \in \mathcal{C}'} \mathsf{vbl}(C)$ for all $\widetilde{C} \in \widetilde{\mathcal{C}}$". Then

$$\mathbf{Pr}\left[ \mathcal{E}(\mathcal{C}', Y, \widetilde{\mathcal{C}}) \right] \leq \left( \frac{k|\mathcal{C}'| + |Y|}{n} \right)^{k|\widetilde{\mathcal{C}}|} \leq \left( \frac{k(\ell - 1) + (k - 6)\ell}{n} \right)^{k\ell} \leq (4k\varepsilon)^{k\ell},$$

where the last inequality is due to $\ell \leq \varepsilon n + 1 \leq 2\varepsilon n$. Therefore by union bound, we have

$$\mathbf{Pr}\left[ \exists \text{ such } \mathcal{E}(\mathcal{C}', Y, \widetilde{\mathcal{C}}) \right] \leq \sum_{t=0}^{(k-6)\ell} \binom{m}{\min\{\ell - 1, m - \ell\}}^2 \cdot \binom{n}{t} \cdot (4k\varepsilon)^{k\ell}.$$

50

Note that $(k-6)\ell \le (k-6)(\varepsilon n+1) \le 2k\varepsilon n \le n/2$ assuming $\varepsilon \le 1/(4k)$. Thus $\binom{n}{t} \le \binom{n}{(k-6)\ell} \le \left(\frac{en}{(k-6)\ell}\right)^{(k-6)\ell}$. Also both $\binom{m}{\ell-1}$ and $\binom{m}{m-\ell}$ are upper bounded by $\left(\frac{em}{\ell-1}\right)^{\ell} = \left(\frac{e\alpha n}{\ell-1}\right)^{\ell}$. Then we have

$$\mathbf{Pr}\left[\exists \text{ such } \mathcal{E}(\mathcal{C}', Y, \widetilde{\mathcal{C}})\right] \le n \cdot \left(\frac{e\alpha n}{\ell-1}\right)^{2\ell} \cdot \left(\frac{en}{(k-6)\ell}\right)^{(k-6)\ell} \cdot (4k\varepsilon)^{k\ell}$$

$$\le n \cdot \left(\frac{e^{k-4} \cdot 2^{4k} \cdot n^{k-4} \cdot k^k \cdot \varepsilon^k}{(\ell-1)^2 \cdot \ell^{k-6} \cdot (k-6)^{k-6}}\right)^{\ell} \qquad \text{(since } m = \alpha n \text{ and } \alpha \le 2^k\text{)}$$

$$\le n \cdot \left(\frac{e^{k-4} \cdot 2^{4k} \cdot n^2 \cdot \varepsilon^6 \cdot k^k}{(\ell-1)^2 (k-6)^{k-6}}\right)^{\ell} \qquad \text{(since } \ell = \lfloor \varepsilon n \rfloor + 1 \ge \varepsilon n\text{)}$$

$$\le n \cdot \left(\frac{e^{k-4} \cdot 2^{4k+2} \cdot \varepsilon^4 \cdot k^k}{(k-6)^{k-6}}\right)^{\ell} \qquad \text{(since } \ell - 1 = \lfloor \varepsilon n \rfloor \ge \varepsilon n/2\text{)}$$

$$\le n \cdot \left(e^{k-4} \cdot 2^{4k+14} \cdot \varepsilon^4 \cdot k^6\right)^{\ell}$$

$$\text{(since } \tfrac{k^k}{(k-6)^{k-6}} \le (4(k-6))^6 \le (4k)^6 \text{ for } k \ge 12\text{)}$$

$$\le n \cdot \left(2^{10k-1} \cdot \varepsilon^4\right)^{\ell} =: \widetilde{p}. \qquad \text{(since } k \ge 12\text{)}$$

Now we have two cases:

- If $\varepsilon n \ge 5\log(n)$, then assuming $2^{10k-1} \cdot \varepsilon^4 \le 1/2$, we have

$$\widetilde{p} \le n \cdot (1/2)^{\ell} \le n \cdot (1/2)^{\varepsilon n} = o(1/n^3).$$

- Otherwise $\varepsilon \le 5\log(n)/n$. Then assuming $n \ge 2^{\Omega(k)}$, we have $2^{10k-1}\varepsilon^4 = o(1/n^3)$. Now since $\varepsilon n \ge 1$, we have $\ell \ge 2$ and $\widetilde{p} \le n \cdot o(1/n^3)^2 = o(1/n^3)$.

Finally we note that if $k \ge 12$, $\alpha \le 2^k$, $1/n \le \varepsilon \le 2^{-2.5k}$, and $n \ge 2^{\Omega(k)}$, then all the assumptions above are satisfied. $\qquad\square$

Now we put explicit parameters into Lemma A.3 to prove Proposition 3.10.

*Proof of Proposition 3.10.* For each $z \in [n/2^{4k}]$, let $\mathcal{E}_z$ be the event that there exists some $\mathcal{C}' \subset \mathcal{C}$ with $|\mathcal{C}'| = z$ that violates the desired property. Now we apply Lemma A.3 with $\varepsilon = z/n$. Notice that if $k \ge 12$ and $n \ge 2^{\Omega(k)}$, then all the assumptions in Lemma A.3 are satisfied. Thus $\mathbf{Pr}[\mathcal{E}_z] = o(1/n^3)$. Then the corollary follows immediately by union bound over all possible $z$ and assuming $n \ge 2^{\Omega(k)}$. $\qquad\square$