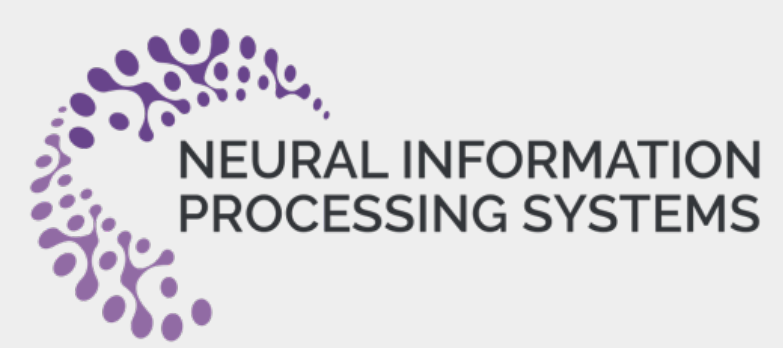


# A Private Approximation of the 2nd-Moment Matrix of Any Subsamplable Input



Bar Mahpud Or Sheffet

Faculty of Engineering  
Bar-Ilan University  
Israel



## Introduction

Estimating the second moment matrix is central to statistics and machine learning. We focus on achieving differential privacy guarantees even for worst-case inputs under a new assumption: **subsamplability**.

## Subsamplability Assumption

A dataset  $X$  is  $(m, \alpha, \beta)$ -subsamplable if a random subsample of  $m' \geq m$  points yields a spectral approximation of the second moment:

$$\Pr[(1 - \alpha)\Sigma \preceq \hat{\Sigma} \preceq (1 + \alpha)\Sigma] \geq 1 - \beta.$$

This property enables robust DP estimation with only minimal assumptions.

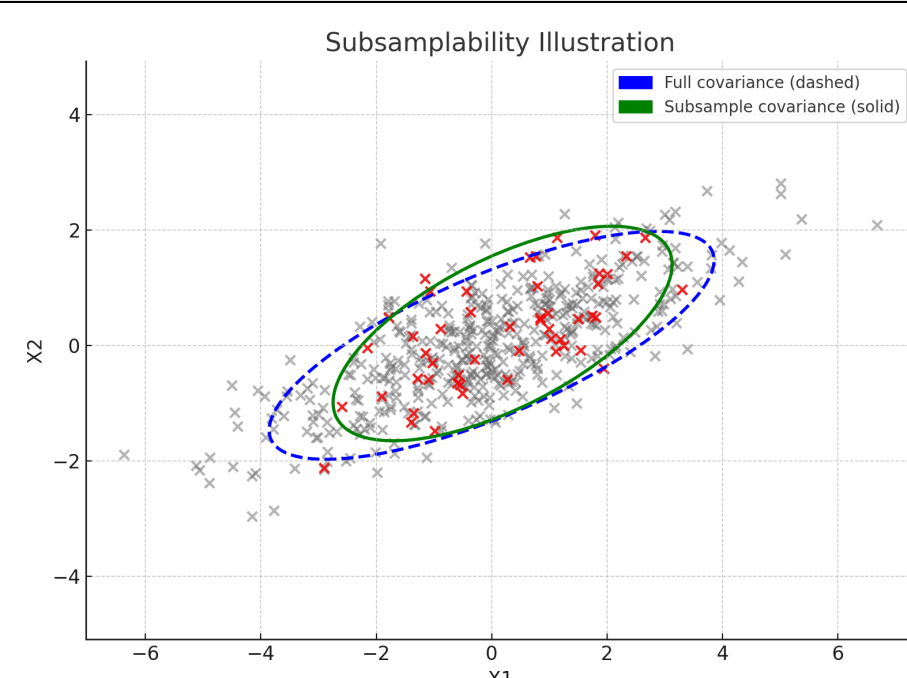


Figure 1. Subsamplability Illustration

## Our Contributions

We Present a new algorithm that achieve strong privacy-utility trade-offs even for worst-case inputs under subsamplability assumptions on the data.

Building upon *subsamplability*, we give a recursive algorithmic framework similar to Kamath et al. [3] that abides zero-Concentrated Differential Privacy (zCDP) while preserving w.h.p the accuracy of the second moment estimation upto an arbitrary factor of  $(1 \pm \gamma)$ .

We then show how to apply our algorithm to approximate the second moment matrix of a distribution  $\mathcal{D}$ , even when a noticeable *fraction* of the input are outliers.

## Guarantees

Let  $P_{\text{tail}} = \{x \in X : \exists u \in \mathbb{R}^d : \langle x, u \rangle^2 > m(1 + \alpha) \cdot \frac{1}{n} \sum_{x \in X} \langle x, u \rangle^2\}$ .

If the input is  $(m, \alpha, \beta)$ -subsamplable with  $\beta = O(\alpha/\log(R))$  and  $\alpha \leq 1/2$ , then w.h.p.:

$$(1 - \gamma)\Sigma_{\text{eff}} \preceq \tilde{\Sigma} \preceq (1 + \gamma)\Sigma$$

where  $\Sigma_{\text{eff}} = \frac{1}{n} \sum_{x \in X \setminus P_{\text{tail}}} xx^T$ .

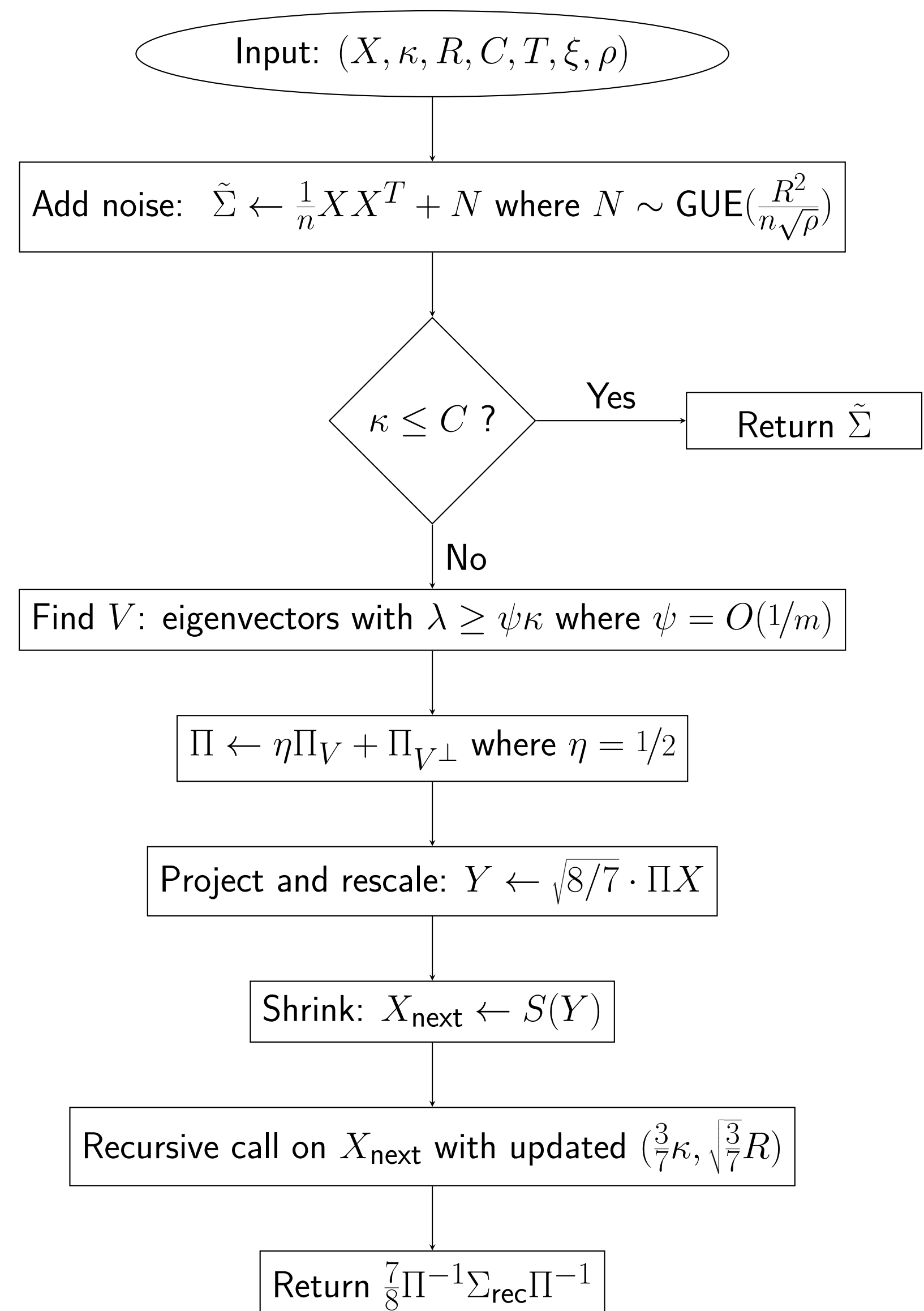


Figure 2. Highlighting the Effective Spectral Guarantee

## Comparison to Prior Work

- **Brown et al. [2]:** Require that every data point has bounded leverage. Their algorithm fails or degrades under even a small number of high-leverage outliers.
- **Previous Work:** Prior algorithms rely on strong distributional assumptions (e.g., Gaussians, sub-Gaussian tails) or on restrictive bounded-norm constraints.
- **Beyond Subsample and Aggregate:** Our algorithm improves over classical frameworks like Subsample-and-Aggregate, which tolerate far fewer outliers and require smaller contamination rates ( $\eta = \tilde{O}(\gamma^2/d)$  vs. ours:  $\eta = \tilde{O}(1/d)$ ).

## Main Algorithm



## Applications

We study the problem of estimating the second moment matrix of a distribution  $\mathcal{D}$ , where the input is corrupted by an  $\eta$ -fraction of arbitrary outliers. Consider  $\mathcal{D}$  to be a distribution that for any  $\alpha, \beta > 0$  is  $m(\alpha, \beta)$ -subsamplable for  $m = O(\frac{d \ln(d/\beta)}{\alpha^2})$ . We thus denote the second moment matrix of the input as  $\Sigma = (1 - \eta)\Sigma_{\mathcal{D}} + \eta\Sigma_{\text{out}}$ .

- **Our Guarantee:** Our method succeeds if a small random subsample is a good spectral approximation. It tolerates many high-leverage outliers so long as  $\eta = \tilde{O}(1/d)$  and  $\Sigma_{\text{out}} \preceq O(1/\eta)\Sigma_{\mathcal{D}}$ . Under these conditions, our DP algorithm returns  $\tilde{\Sigma} \succeq (1 - O(\gamma))\Sigma_{\mathcal{D}}$  w.p.  $\geq 1 - \xi$ .
- **Sample Complexity:** baseline:  $\tilde{O}(\frac{d^2}{\gamma^3\sqrt{\rho}})$  vs. ours:  $\tilde{O}(\frac{d}{\gamma^2} + \frac{d^{3/2}}{\gamma\sqrt{\rho}})$

Want to know more?

<https://arxiv.org/abs/2505.14251>



## References

- [1] Hassan Ashtiani and Christopher Liaw. Private and polynomial time algorithms for learning gaussians and beyond. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 1075–1076. PMLR, 02–05 Jul 2022. URL <https://proceedings.mlr.press/v178/ashtiani22a.html>.
- [2] Gavin Brown, Samuel Hopkins, and Adam Smith. Fast, sample-efficient, affine-invariant private mean and covariance estimation for subgaussian distributions. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 5578–5579. PMLR, 12–15 Jul 2023. URL <https://proceedings.mlr.press/v195/brown23a.html>.
- [3] Gautam Kamath, Jerry Li, Vikrant Singhal, and Jonathan Ullman. Privately learning high-dimensional distributions. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 1853–1902. PMLR, 25–28 Jun 2019. URL <https://proceedings.mlr.press/v99/kamath19a.html>.