

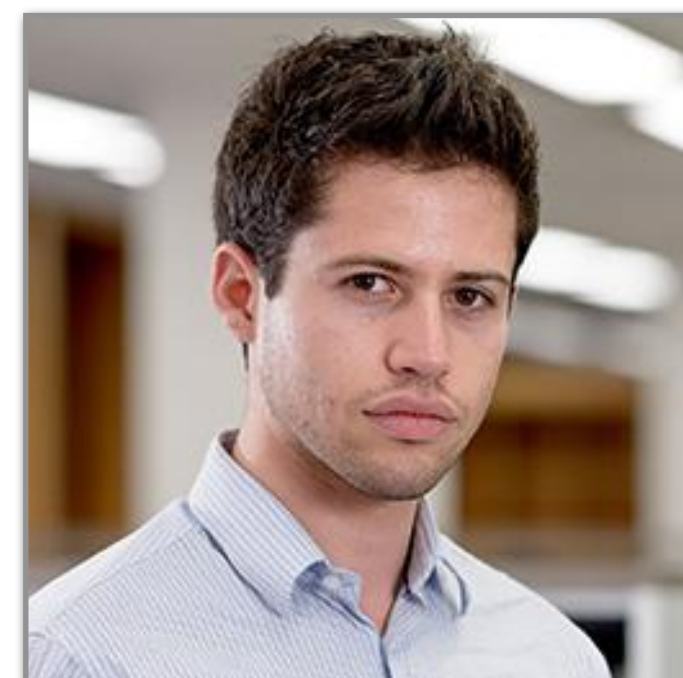
ON SIGNAL-TO-NOISE RATIO ISSUES IN VARIATIONAL INFERENCE FOR DEEP GAUSSIAN PROCESSES



TIM G. J. RUDNER*



OSCAR KEY*



YARIN GAL



TOM RAINFORTH

— ICML 2021 —

Correspondence to

`tim.rudner@cs.ox.ac.uk`

Conventional Deep GPs (Damianou et al., [2013]; Salimbeni et al. [2017])

$$y_n = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)} (x_n) \right) \dots \right) + \varepsilon_n \quad \varepsilon_n \sim \mathcal{N} (0, \sigma^2 I)$$

Conventional Deep GPs (Damianou et al., [2013]; Salimbeni et al. [2017])

$$y_n = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)} (x_n) \right) \dots \right) + \varepsilon_n \quad \varepsilon_n \sim \mathcal{N} (0, \sigma^2 I)$$

Latent-Variable Deep GPs (Salimbeni et al. [2018])

$$y_n = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)} ([x_n, z_n]) \right) \dots \right) + \varepsilon_n \quad \varepsilon_n \sim \mathcal{N} (0, \sigma^2 I)$$
$$z_n \sim \mathcal{N} (0, I_{\tilde{D}})$$

LATENT-VARIABLE DEEP GAUSSIAN PROCESSES

Conventional Deep GPs (Damianou et al., [2013]; Salimbeni et al. [2017])

$$y_n = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)} (x_n) \dots \right) \right) + \varepsilon_n \quad \varepsilon_n \sim \mathcal{N} (0, \sigma^2 I)$$

Latent-Variable Deep GPs (Salimbeni et al. [2018])

$$y_n = f^{(L)} \left(f^{(L-1)} \left(\dots f^{(1)} ([x_n, z_n]) \dots \right) \right) + \varepsilon_n \quad \varepsilon_n \sim \mathcal{N} (0, \sigma^2 I)$$

$$z_n \sim \mathcal{N} (0, I_{\tilde{D}})$$

Latent Variable



2-Layer Latent-Variable Deep GP Model

$$y_n = f^{(2)} \left(f^{(1)} \left([x_n, z_n] \right) \right) + \varepsilon_n$$

$$\varepsilon_n \sim \mathcal{N} \left(0, \sigma^2 I \right)$$

$$z_n \sim \mathcal{N} \left(0, I_{\tilde{D}} \right)$$

2-Layer Latent-Variable Deep GP Model

$$y_n = f^{(2)} \left(f^{(1)} \left([x_n, z_n] \right) \right) + \varepsilon_n$$

$$\varepsilon_n \sim \mathcal{N} \left(0, \sigma^2 I \right)$$

$$z_n \sim \mathcal{N} \left(0, I_{\tilde{D}} \right)$$

Variational Objective (Salimbeni et al. [2018])

$$\mathcal{L}_K \stackrel{\text{def}}{=} \mathbb{E}_{f_{1:K}^{(1)}, z_{n,1:K}} \left[\sum_n \log \frac{1}{K} \sum_{k=1}^K \frac{\mathcal{F} \left(x_n, y_n, f_k^{(1)}, z_{n,k} \right) p \left(z_{n,k} \right)}{q_\phi \left(z_{n,k} \right)} \right] - \sum_{\ell=1}^2 D_{\text{KL}} \left(q \left(f^{(\ell)} \right) \parallel p \left(f^{(\ell)} \right) \right)$$

2-Layer Latent-Variable Deep GP Model

$$y_n = f^{(2)} \left(f^{(1)} \left([x_n, z_n] \right) \right) + \varepsilon_n$$

$$\varepsilon_n \sim \mathcal{N} \left(0, \sigma^2 I \right)$$

$$z_n \sim \mathcal{N} \left(0, I_{\tilde{D}} \right)$$

Variational Objective (Salimbeni et al. [2018])

$$\mathcal{L}_K \stackrel{\text{def}}{=} \mathbb{E}_{f_{1:K}^{(1)}, z_{n,1:K}} \left[\sum_n \log \frac{1}{K} \sum_{k=1}^K \frac{\mathcal{F} \left(x_n, y_n, f_k^{(1)}, z_{n,k} \right) p \left(z_{n,k} \right)}{q_\phi \left(z_{n,k} \right)} \right] - \sum_{\ell=1}^2 D_{\text{KL}} \left(q \left(f^{(\ell)} \right) \parallel p \left(f^{(\ell)} \right) \right)$$

Importance Weight = $w_{n,k}$

MAIN RESULT: SNR DECREASES AS K INCREASES

Theorem 1.

Assume mild regularity conditions and let

$$\Delta_{n,M,K}^{\text{DGP}}(\phi) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \nabla_{\phi} \log \frac{1}{K} \sum_{k=1}^K w_{n,m,k}.$$

Then, the signal-to-noise ratio of gradient estimator $\Delta_{n,M,K}^{\text{DGP}}(\phi)$ is

$$\text{SNR}_{n,M,K}^{\text{DGP}}(\phi) = \frac{|\mathbb{E} [\Delta_{n,M,K}^{\text{DGP}}(\phi)]|}{\sqrt{\text{Var} [\Delta_{n,M,K}^{\text{DGP}}(\phi)]}} = \mathcal{O} \left(\sqrt{M/K} \right)$$

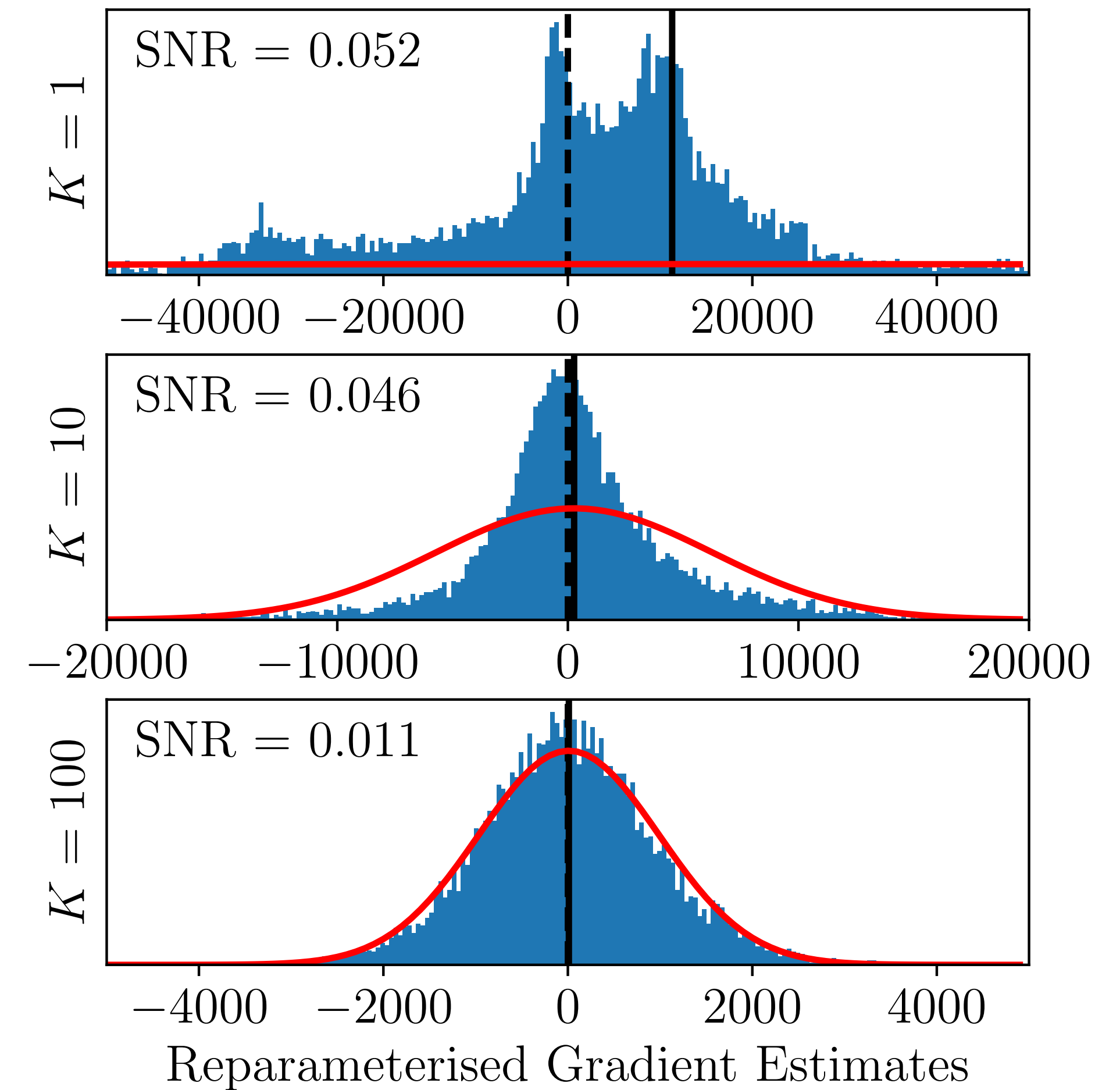
EMPIRICAL CONFIRMATION OF SNR DETERIORATION

► Gradient Estimate

$$\Delta_{n,M,K}^{\text{DGP}}(\phi) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \nabla_{\phi} \log \frac{1}{K} \sum_{k=1}^K w_{n,m,k}$$

► Signal-to-Noise Ratio

$$\text{SNR}_{n,M,K}^{\text{DGP}}(\phi) = \frac{|\mathbb{E} [\Delta_{n,M,K}^{\text{DGP}}(\phi)]|}{\sqrt{\text{Var} [\Delta_{n,M,K}^{\text{DGP}}(\phi)]}} = \mathcal{O} \left(\sqrt{M/K} \right)$$



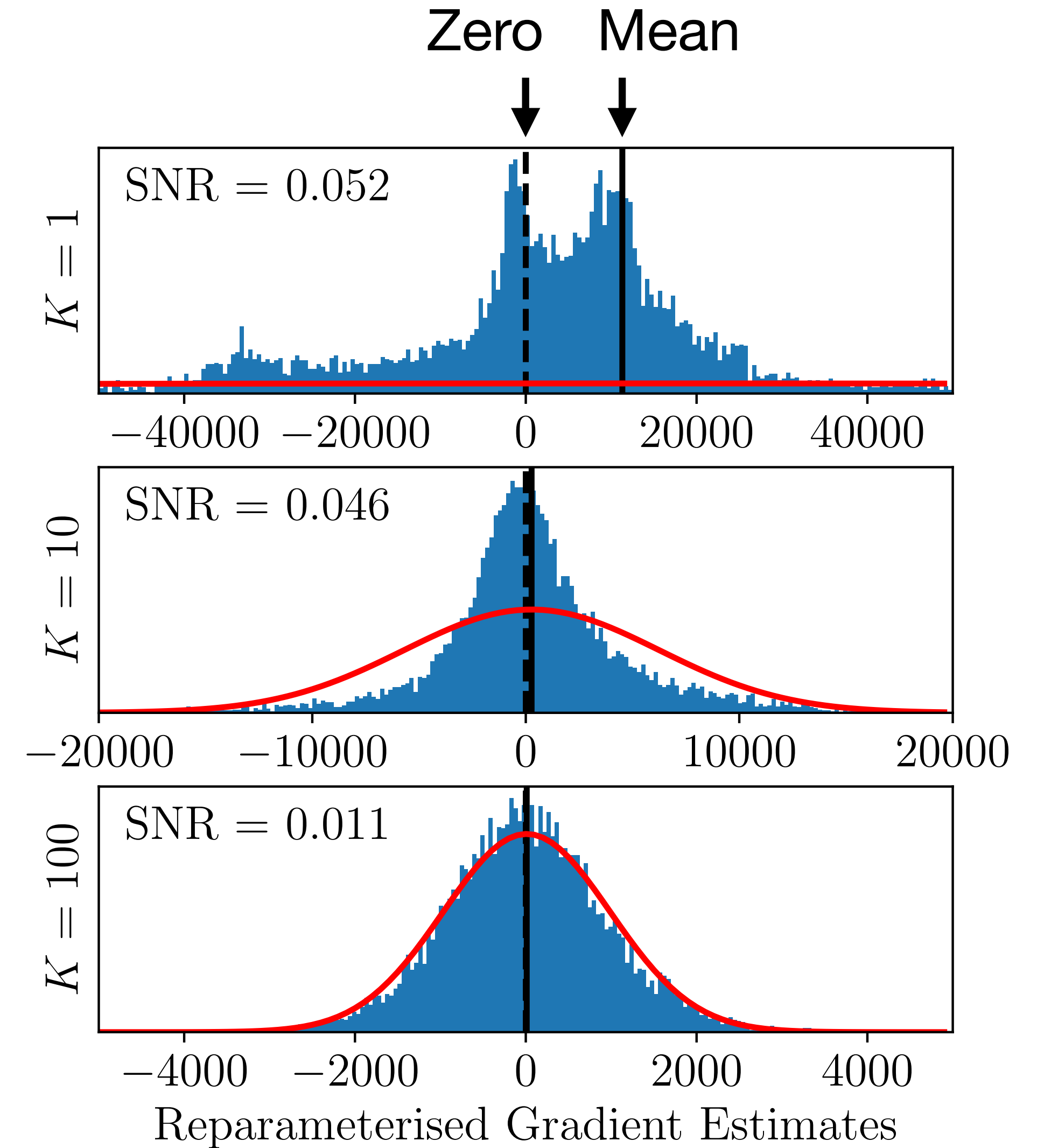
EMPIRICAL CONFIRMATION OF SNR DETERIORATION

► Gradient Estimate

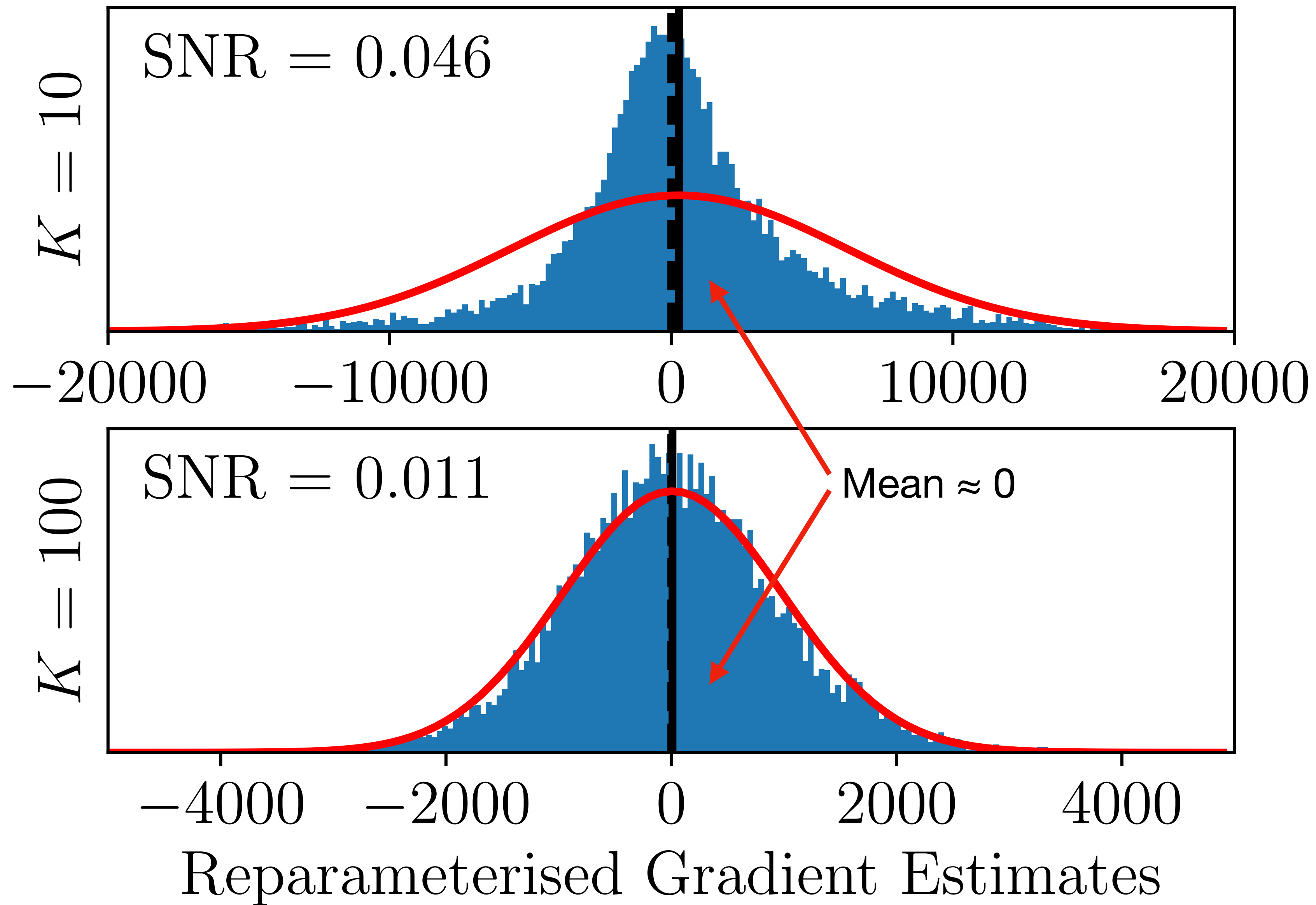
$$\Delta_{n,M,K}^{\text{DGP}}(\phi) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \nabla_{\phi} \log \frac{1}{K} \sum_{k=1}^K w_{n,m,k}$$

► Signal-to-Noise Ratio

$$\text{SNR}_{n,M,K}^{\text{DGP}}(\phi) = \frac{|\mathbb{E} [\Delta_{n,M,K}^{\text{DGP}}(\phi)]|}{\sqrt{\text{Var} [\Delta_{n,M,K}^{\text{DGP}}(\phi)]}} = \mathcal{O} \left(\sqrt{M/K} \right)$$

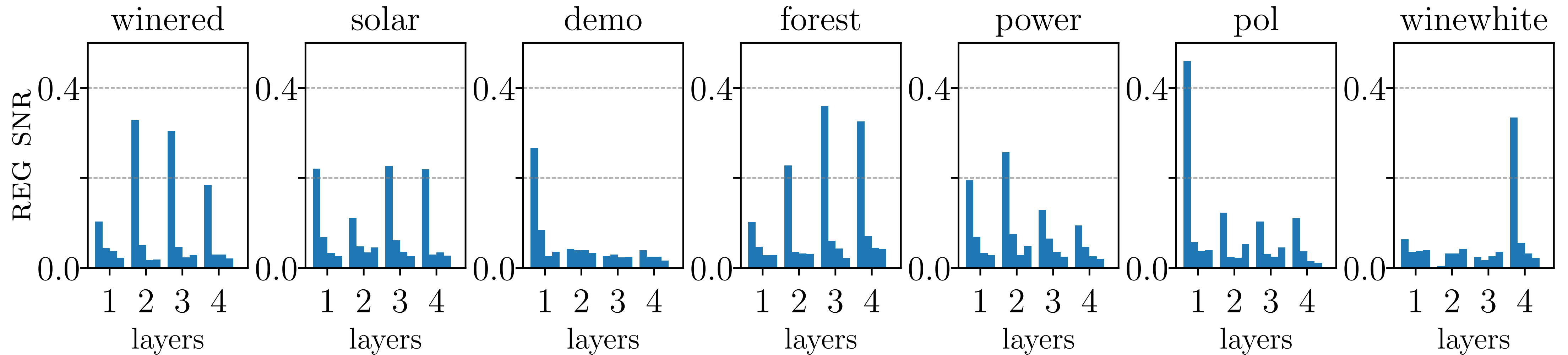


EMPIRICAL CONFIRMATION OF SNR DETERIORATION



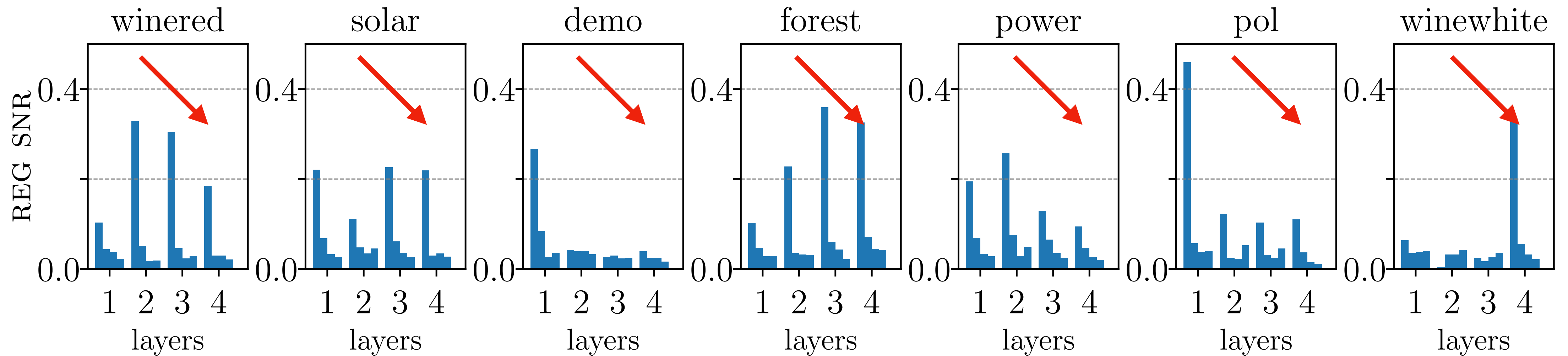
SIGNAL-TO-NOISE RATIO ISSUES IN DEEP GPs

Empirical Confirmation of SNR Deterioration



SIGNAL-TO-NOISE RATIO ISSUES IN DEEP GPs

Empirical Confirmation of SNR Deterioration



Deep GP Doubly Reparameterized Gradient Estimator

Deep GP Doubly Reparameterized Gradient Estimator

- ▶ Adapted to deep GPs from Tucker et al. [2018]:

$$\tilde{\Delta}_{n,M,K}^{\text{DGP}}(\phi) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \sum_{k=1}^K \left(\frac{w_{n,m,k}}{\sum_{j=1}^K w_{n,m,j}} \right)^2 \frac{\partial \log w_{n,m,k}}{\partial z_{n,m,k}} \frac{\partial z_{n,m,k}}{\partial \phi}$$

Deep GP Doubly Reparameterized Gradient Estimator

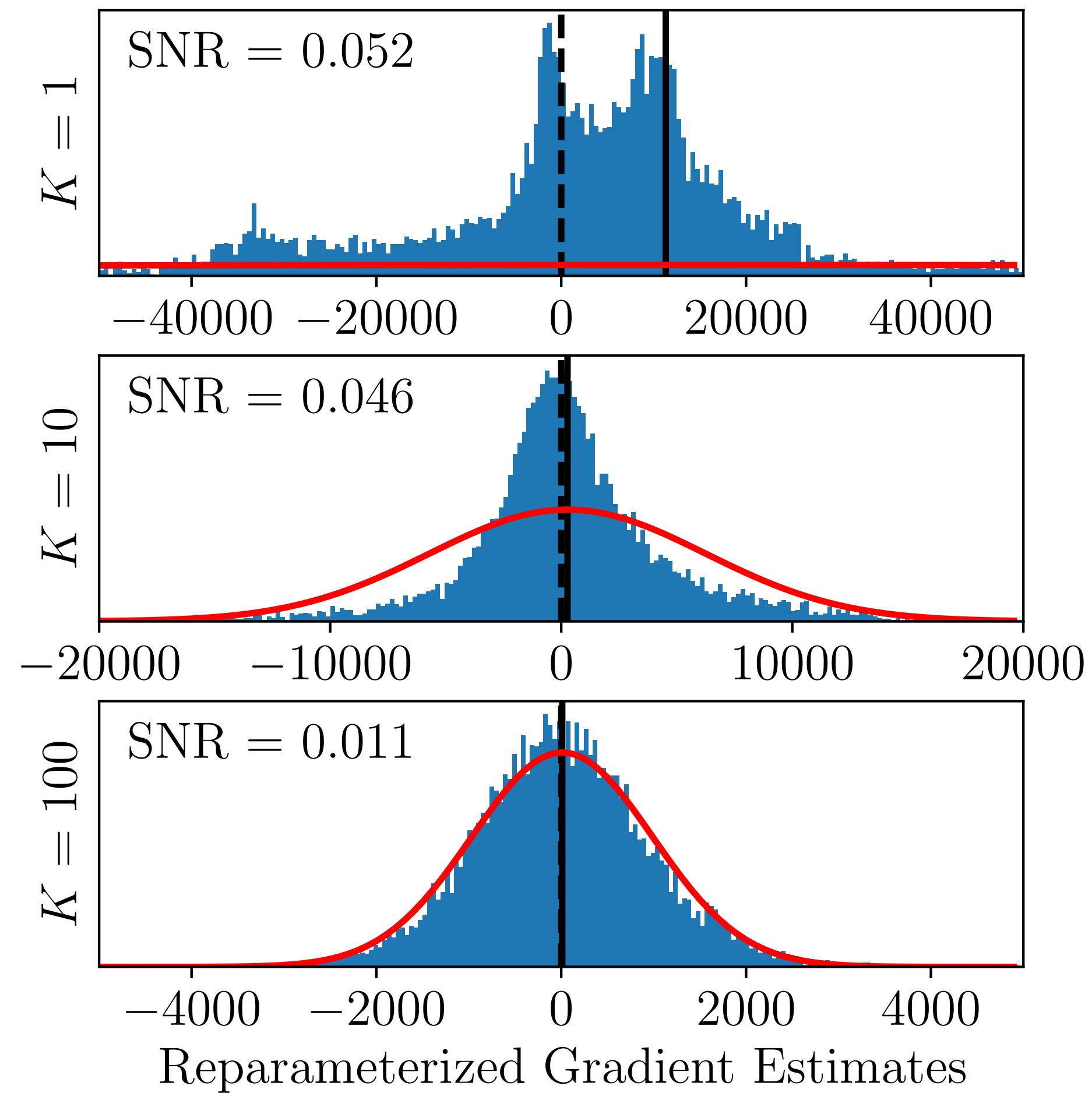
- ▶ Adapted to deep GPs from Tucker et al. [2018]:

$$\tilde{\Delta}_{n,M,K}^{\text{DGP}}(\phi) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \sum_{k=1}^K \left(\frac{w_{n,m,k}}{\sum_{j=1}^K w_{n,m,j}} \right)^2 \frac{\partial \log w_{n,m,k}}{\partial z_{n,m,k}} \frac{\partial z_{n,m,k}}{\partial \phi}$$

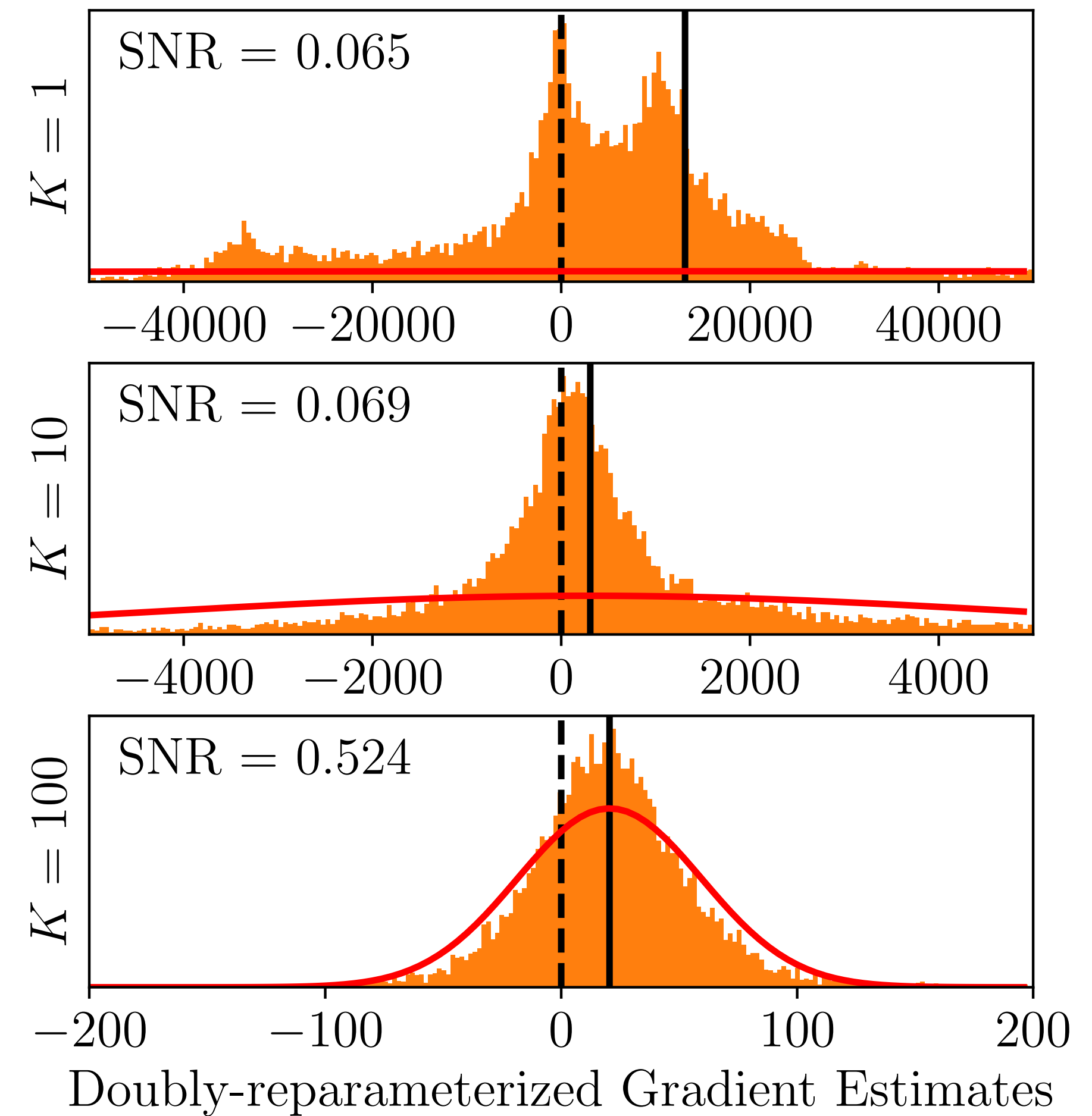
- ▶ More importance samples $K \rightarrow$ higher SNR

FIXING THE PATHOLOGY: EMPIRICAL CONFIRMATION

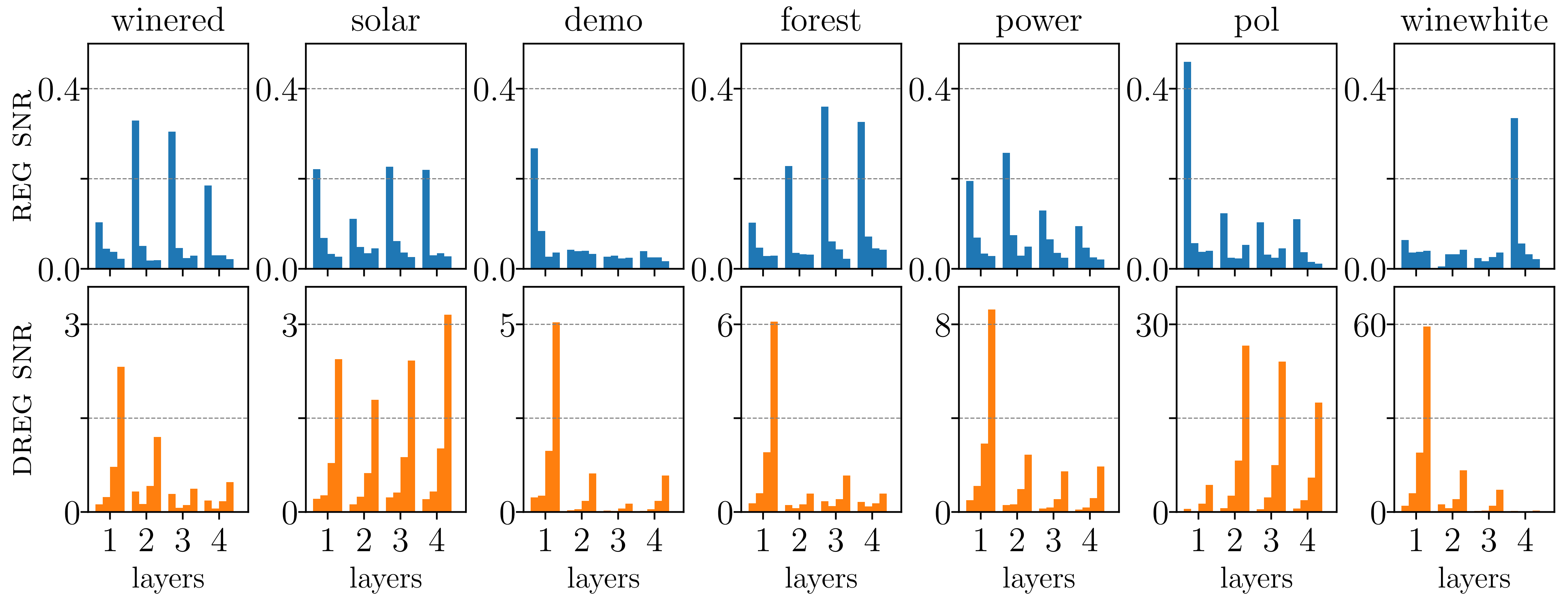
ReG



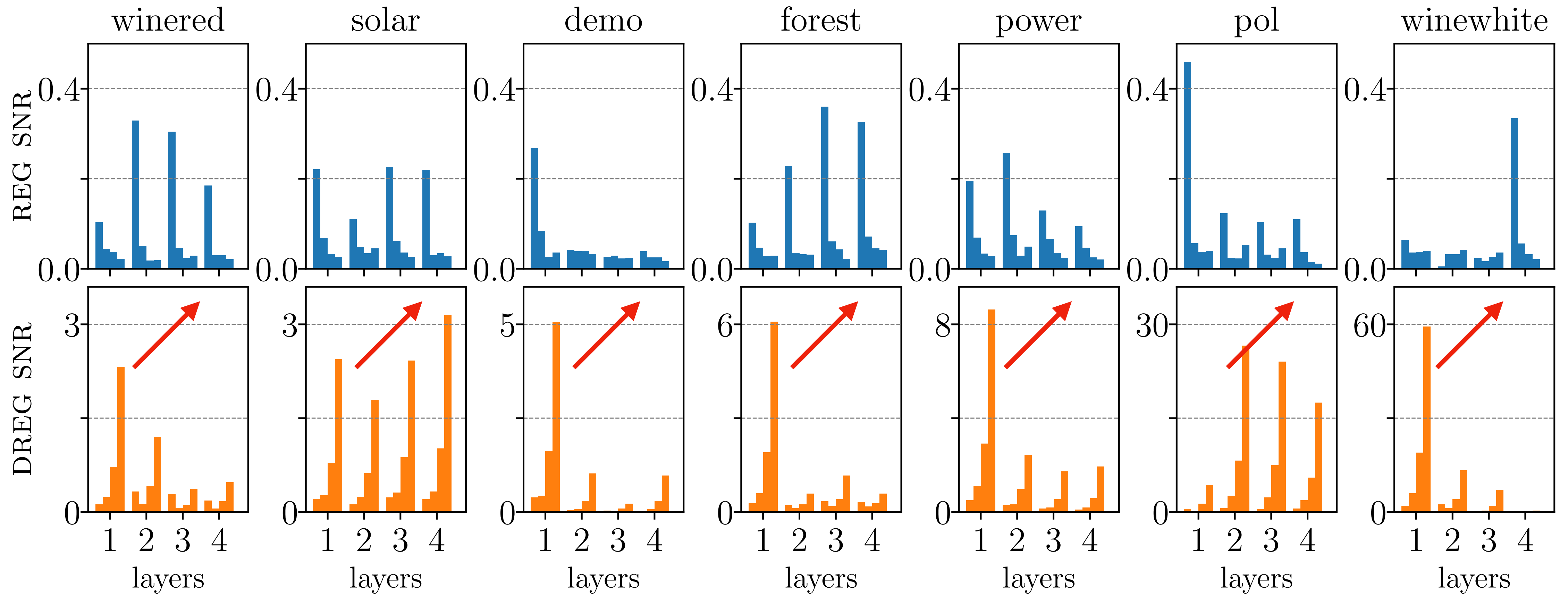
DReG



FIXING THE PATHOLOGY: EMPIRICAL CONFIRMATION



FIXING THE PATHOLOGY: EMPIRICAL CONFIRMATION



FIXING THE PATHOLOGY: EMPIRICAL CONFIRMATION

Improvement in Predictive Performance

Dataset	Train ELBO ($K = 50$)				Test log-likelihood				Wilcoxon Test
	REG		DREG		REG		DREG		
	Mean	SE	Mean	SE	Mean	SE	Mean	SE	p -value
forest	-97.56	(11.04)	-92.53	(10.42)	0.59	(0.08)	0.63	(0.08)	0.1%
solar	1657.41	(27.56)	1707.75	(42.20)	2.33	(0.17)	2.57	(0.11)	2.8%
pol	34610.49	(66.18)	34665.08	(70.34)	2.99	(0.01)	2.99	(0.01)	24.7%
power	1510.50	(10.62)	1515.60	(10.16)	0.21	(0.01)	0.21	(0.01)	67.3%
winewhite	-4701.26	(4.92)	-4703.14	(4.98)	-1.11	(0.01)	-1.11	(0.01)	50.0%
winered	447.91	(249.81)	314.75	(216.32)	0.57	(0.27)	0.61	(0.20)	41.1%
Across Datasets:									1.2%

FIXING THE PATHOLOGY: EMPIRICAL CONFIRMATION

Improvement in Predictive Performance

Dataset	Train ELBO ($K = 50$)				Test log-likelihood				Wilcoxon Test
	REG		DREG		REG		DREG		
	Mean	SE	Mean	SE	Mean	SE	Mean	SE	p -value
forest	-97.56	(11.04)	-92.53	(10.42)	0.59	(0.08)	0.63	(0.08)	0.1%
solar	1657.41	(27.56)	1707.75	(42.20)	2.33	(0.17)	2.57	(0.11)	2.8%
pol	34610.49	(66.18)	34665.08	(70.34)	2.99	(0.01)	2.99	(0.01)	24.7%
power	1510.50	(10.62)	1515.60	(10.16)	0.21	(0.01)	0.21	(0.01)	67.3%
winewhite	-4701.26	(4.92)	-4703.14	(4.98)	-1.11	(0.01)	-1.11	(0.01)	50.0%
winered	447.91	(249.81)	314.75	(216.32)	0.57	(0.27)	0.61	(0.20)	41.1%
Across Datasets:									1.2%

ON SNR ISSUES IN VARIATIONAL INFERENCE FOR DEEP GPs

Summary

- ▶ **Deteriorating gradient SNR** in variational inference for deep GPs.
- ▶ Fixing the SNR deterioration leads to **improved performance**.
- ▶ Improvement comes at **no additional computational cost**.

THANK YOU!



Tim G. J. Rudner
@timrudner



Oscar Key



Yarin Gal
@yaringal



Tom Rainforth
@tom_rainfroth

CORRESPONDENCE: `tim.rudner@cs.ox.ac.uk`

CODE: `https://github.com/timrudner/snr_issues_in_deep_gps`