

Exact Gaussian Processes on a Million Data Points

Ke Alexander Wang^{*1}, Geoff Pleiss^{*1}, Jacob R. Gardner², Stephen Tyree³, Kilian Q. Weinberger¹, Andrew Gordon Wilson⁴

¹Cornell University, ²Uber AI Labs, ³NVIDIA, ⁴New York University

Summary

- 1 We scale exact GPs to over 10^6 training points using multi-GPU parallelism and conjugate gradients-based inference. On 10^6 data points, training takes less than 2 hours.
- 2 We demonstrate that computing predictive distributions with exact GPs is fast and practical on consumer-grade GPUs.
- 3 We perform the first-ever comparison of exact GPs against GP approximations on datasets with 10^4 – 10^6 data points.

Background: matrix-multiplication inference

Given n training data points (X, \mathbf{y}) , GP training requires optimizing model hyperparameters θ (e.g. kernel lengthscale, observed noise):

$$\min_{\theta} \mathbf{y}^T \widehat{K}(\theta)_{XX}^{-1} \mathbf{y} + \log |\widehat{K}(\theta)_{XX}| \quad (1)$$

($\widehat{K}(\theta)_{XX}$ is the $n \times n$ covariance matrix with observational noise.) We use Black-Box Matrix-Matrix (BBMM) inference [1] to reduce the time per optimization iteration from $O(n^3)$ to $O(n^2)$ by relying only on matrix multiplications and conjugate gradients (CG).

Scaling GPs from $n = 10^4$ to $n = 10^6$

Reducing BBMM's memory from $O(n^2)$ to $O(n)$:

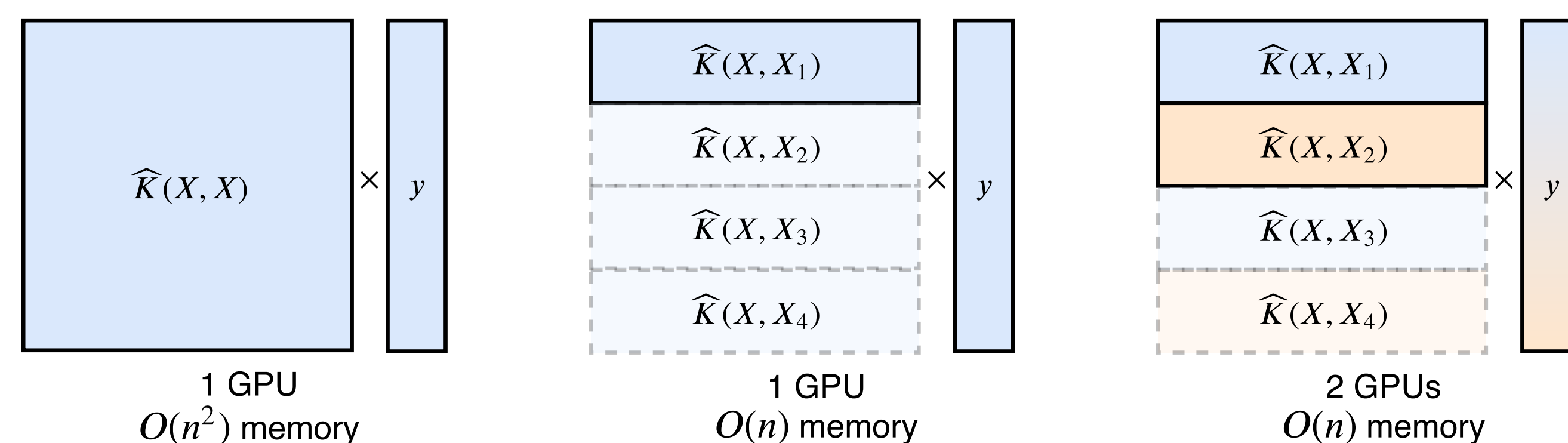


Figure 1: Partitioned kernel matrix-vector multiplication.

Additional techniques to speed up BBMM training:

- 1 Distribute the multiplication partitions across multiple GPUs.
- 2 Initializing hyperparameters θ from a GP trained on $< 10\%$ of data.
- 3 Increasing the rank of the pivoted Cholesky preconditioner (from rank-5 to rank-100) to compute $\widehat{K}(\theta)_{XX}^{-1} \mathbf{y}$
- 4 Using a looser CG convergence criterion at train time.

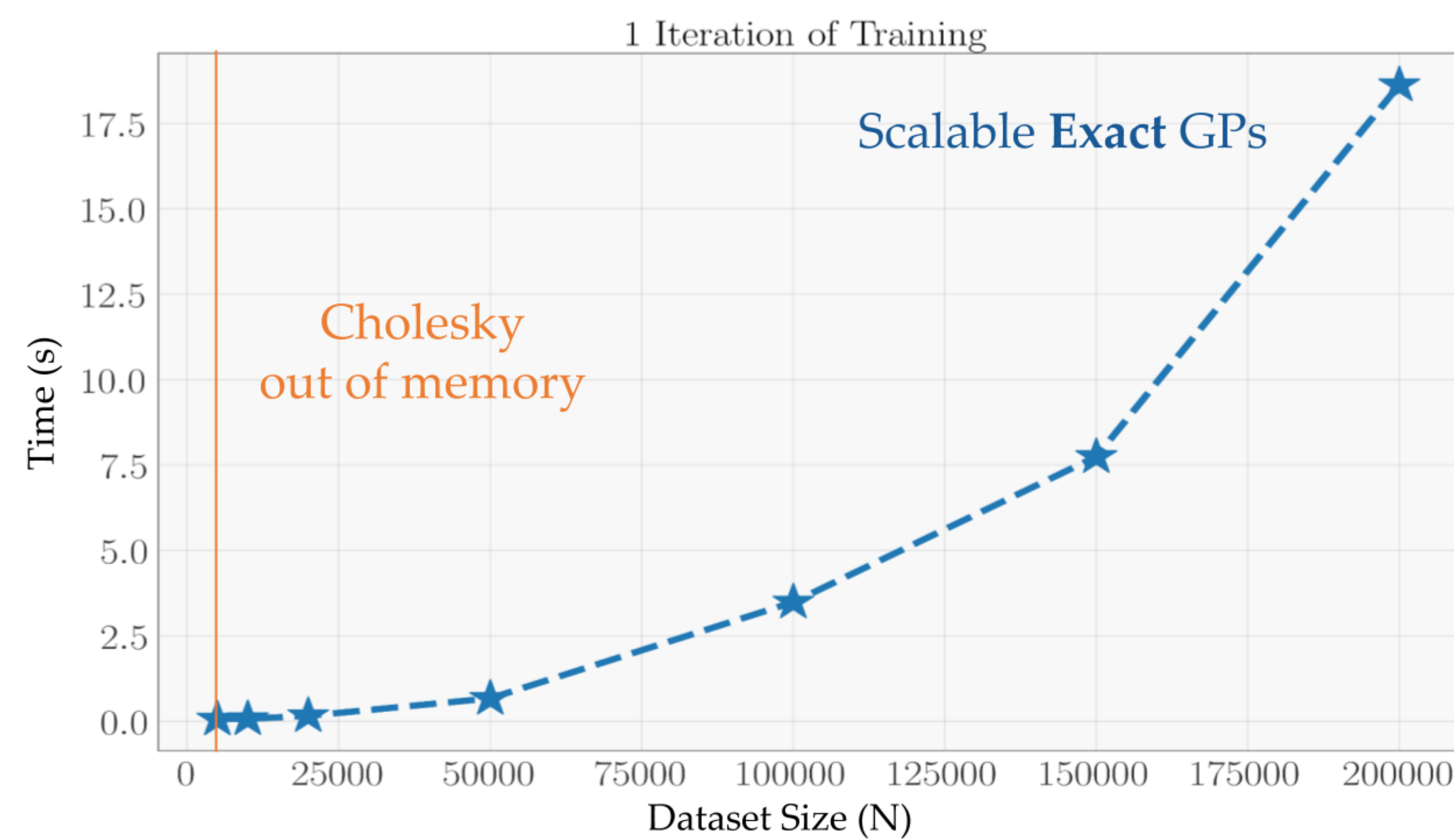


Figure 2: Conventional Cholesky-based inference versus multi-GPU CG-based inference.

Comparisons against approximate methods

We compare against approximate GP methods SGPR [2] and SVGP [3], two popular inducing point methods for large scale GP regression.

Dataset	n	d	RMSE			NLL		
			Exact GP (BBMM)	SGPR ($m=512$)	SVGP ($m=1,024$)	Exact GP (BBMM)	SGPR ($m=512$)	SVGP ($m=1,024$)
PoleTele	9,600	26	0.088 \pm 0.003	0.113 \pm 0.005	0.109 \pm 0.002	-0.660 \pm 0.081	-0.817 \pm 0.005	-0.644 \pm 0.008
Elevators	10,623	18	0.399 \pm 0.011	0.426 \pm 0.007	0.388 \pm 0.010	0.626 \pm 0.043	0.528 \pm 0.015	0.486 \pm 0.019
Bike	11,122	17	0.043 \pm 0.012	0.094 \pm 0.010	0.077 \pm 0.005	-1.323 \pm 0.170	-0.805 \pm 0.005	-0.984 \pm 0.021
Kin40K	25,600	8	0.080 \pm 0.001	0.225 \pm 0.026	0.240 \pm 0.007	-0.755 \pm 0.009	-0.073 \pm 0.055	0.091 \pm 0.033
Protein	29,267	9	0.511 \pm 0.009	0.619 \pm 0.003	0.613 \pm 0.011	0.960 \pm 0.033	0.915 \pm 0.004	0.952 \pm 0.018
KeggDirected	31,248	20	0.083 \pm 0.001	0.104 \pm 0.002	0.105 \pm 0.003	-0.838 \pm 0.031	-1.163 \pm 0.005	-0.853 \pm 0.033
CTslice	34,240	385	0.497 \pm 0.029	0.217 \pm 0.009	1.004 \pm 0.005	0.939 \pm 0.004	-0.037 \pm 0.060	1.423 \pm 0.005
KEGGU	40,708	27	0.120 \pm 0.001	0.130 \pm 0.001	0.126 \pm 0.002	-0.540 \pm 0.035	-1.049 \pm 0.010	-0.653 \pm 0.013
3DRoad	278,319	3	0.110 \pm 0.017	0.578 \pm 0.001	0.390 \pm 0.005	1.239 \pm 0.025	0.791 \pm 0.033	0.486 \pm 0.010
Song	329,820	90	0.774 \pm 0.001	0.816 \pm 0.038	0.998 \pm 0.000	1.162 \pm 0.002	1.243 \pm 0.083	1.417 \pm 0.000
Buzz	373,280	77	0.279 \pm 0.002	0.289 \pm 0.001	0.270 \pm 0.012	0.161 \pm 0.026	0.092 \pm 0.017	0.119 \pm 0.042
HouseElectric	1,311,539	9	0.054 \pm 0.000	—	0.127 \pm 0.046	-0.207 \pm 0.001	—	0.024 \pm 0.984

Table 1: Exact GP vs SGPR vs SVGP using a Matern 3/2 kernel with independent lengthscales.

Training times and prediction times

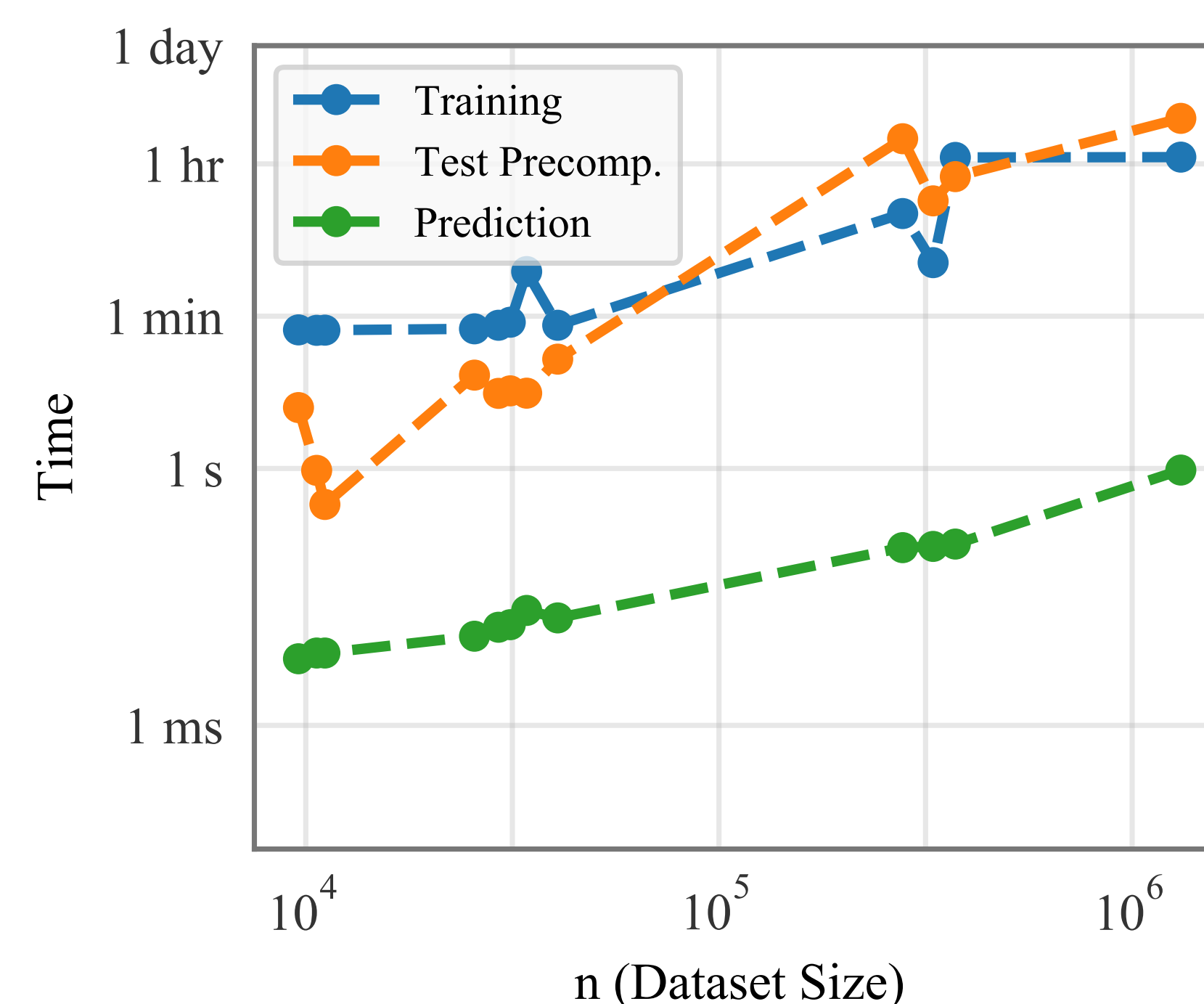


Figure 3: Training, test-time precomputation, and prediction times as a function of the number of inducing points. All predictions for exact GPs can be done in less than a second.

How much does pretraining help?

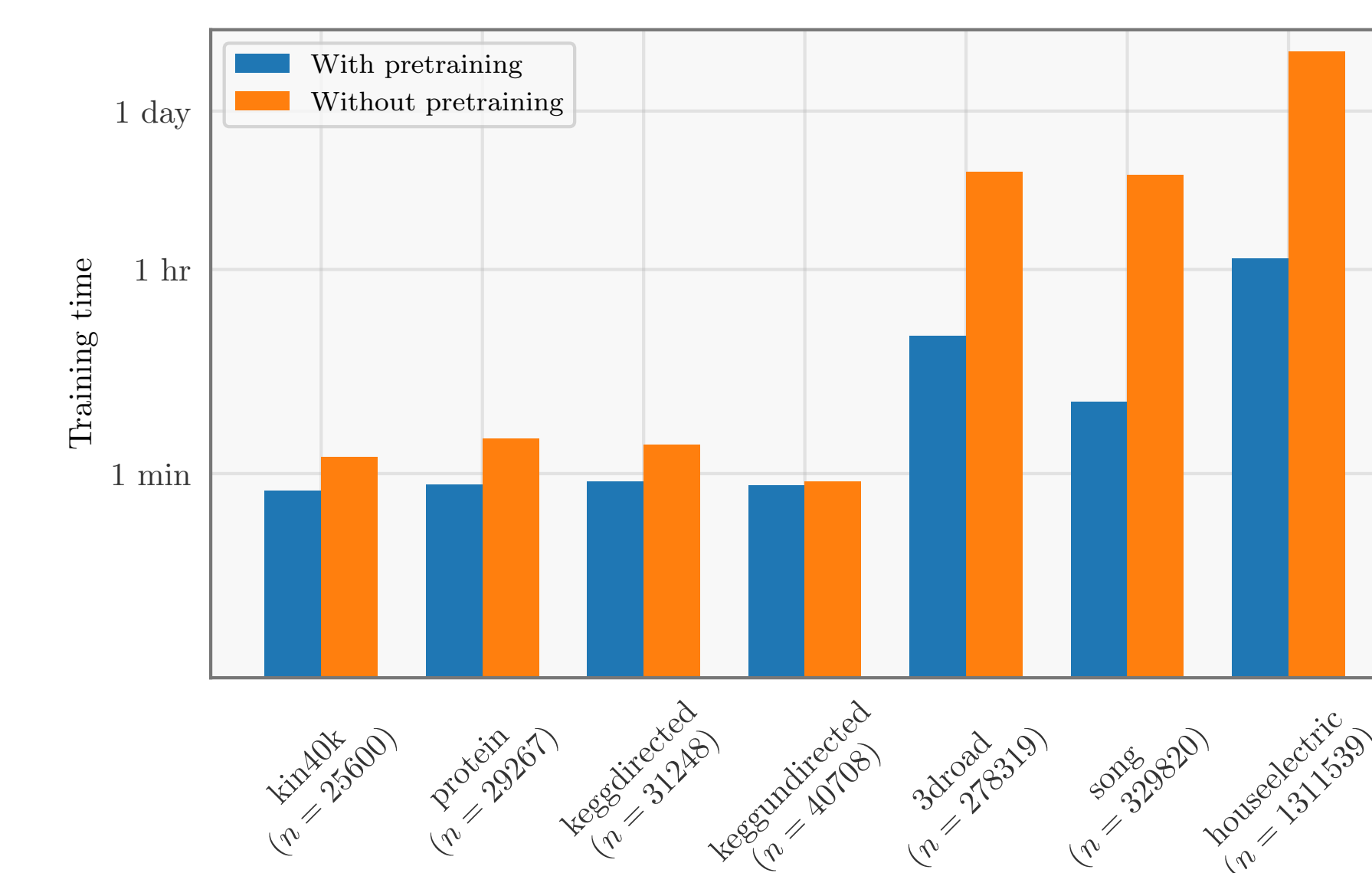


Figure 4: Fine-tuning on the full training set after pre-training on a smaller subset significantly reduces training time. Models with and without pretraining achieve similar errors with less than 5% difference (not shown here).

Do GPs need the entire dataset?

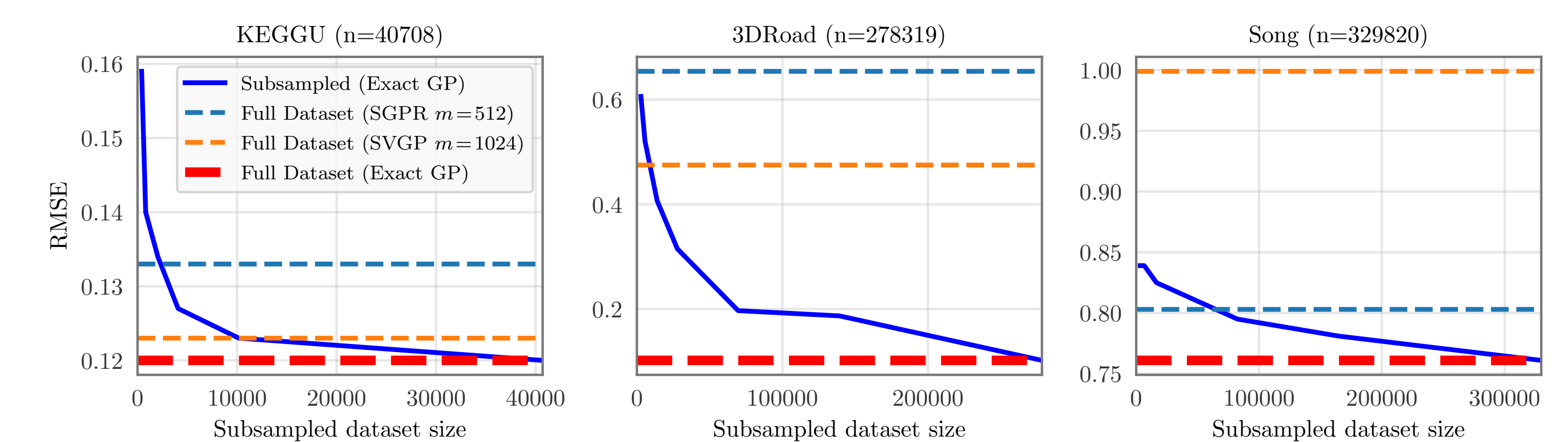


Figure 5: Test root-mean-square error (RMSE) vs. subsampled dataset size.

Would more inducing points help?

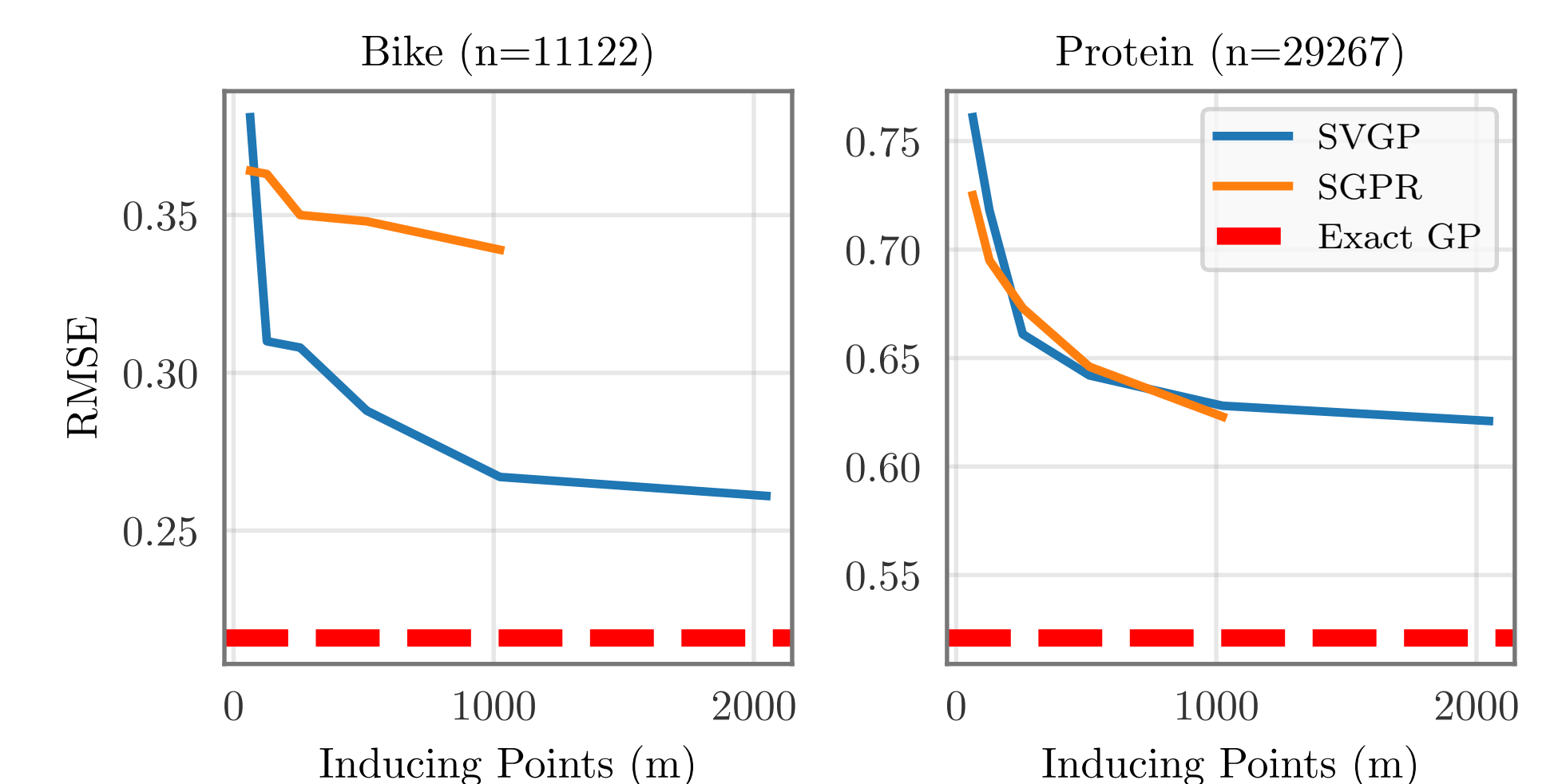


Figure 6: Test root-mean-square error (RMSE) of SVGP and SGPR methods as a function of the number of inducing points.

References

- [1] Jacob Gardner, Geoff Pleiss, Kilian Q. Weinberger, David Bindel, and Andrew G. Wilson. GPyTorch: Blackbox matrix-matrix Gaussian process inference with GPU acceleration. In *NeurIPS*, pages 7587–7597, 2018.
- [2] Michalis K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *AISTATS*, pages 567–574, 2009.
- [3] James Hensman, Nicolo Fusi, and Neil D. Lawrence. Gaussian processes for big data. In *UAI*, 2013.