## DIFFUSION MODELS ARE FEW-SHOT LEARNERS FOR DENSE VISION TASKS

Anonymous authors

Paper under double-blind review

Appendix

004

006

008 009

A BROADER IMPACT

014 Our method, which employs diffusion for general few-shot dense tasks, offers significant advantages 015 beyond technical improvements. It substantially reduces labor costs associated with pixel-by-pixel 016 annotation of visual dense tasks, making model deployment more cost-effective and accessible, 017 especially for resource-limited projects. Additionally, the few-shot nature of our approach reduces 018 energy consumption, lowering the environmental impact by decreasing the need for extensive data 019 and computational resources. This aligns with broader goals of energy conservation and emission reduction. By democratizing access to advanced machine learning technologies, our method enables smaller entities and individuals to innovate and implement AI solutions, promoting more responsible 021 and ethical AI development. 022

023 024

025

## B RESULTS ON FULLY TRAINING SET

We include the results of full training set in Table. 1. Although VPD's performance in the few-shot setting is not strong, with more training data, we can see that its performance improves significantly because it fine-tunes more parameters. In contrast, we only fine-tune the concept embeddings with a few hundred parameters. However, our method still outperforms VTM even after training on the full training set, demonstrating the higher potential of the diffusion prior. We also reported the 95% confidence interval, and it can be seen that our method, leveraging a very general prior, achieved more stable results compared to VTM.

033 034

## C RESULTS WITH DIFFERENT NUMBER OF TRAINING SAMPLES

In Fig 1, we illustrate the impact of using 10, 20, 50, and 100 training samples on our method and VPD across all 12 tasks. It can be observed that our method consistently adapts better to new tasks compared to VPD when fewer than 100 training examples are provided. Moreover, as the number of training samples increases, the performance of both methods improves accordingly.

- 040
- 041
- 042 043
- 044
- 045
- 046
- 047
- 048 049
- 049
- 050
- 051
- 053

Table 1: We present the results on 10 tasks from Taskonomy and 2 tasks from NYUv2. For Taskonomy tasks, 10-shot training examples are used for each of them, and for NYU tasks, we use 20 examples. To also evaluate the statistical robustness, we run each number for 100 times and report the 95% confidence interval. Besides segmentation task, lower number indicates better performance. Our method consistently outperforms VTM on all few-shot tasks, especially on out-of-domain tasks. And our method better unleashes the power of diffusion prior for few-shot dense prediction compared to VPD.

068 069		Few-shot			Fully Supervised		
070 071 072		VTM	VPD	Ours	VTM	VPD	Ours
073 074 075 076 077	EucDepth	$0.0812 \pm 0.0065$	$0.1056 \\ \pm 0.0102$	<b>0.0776</b> ±0.0072	0.0524	0.0456	0.0498
	Z-depth	$0.0347 \pm 0.0035$	$0.0404 \pm 0.0037$	<b>0.0308</b> ±0.0038	0.0257	0.0210	0.0236
078 079 080	2DEdge	$0.0818 \pm 0.0021$	$0.0965 \\ \pm 0.0023$	<b>0.0625</b> ±0.0022	0.0154	0.0131	0.0136
081 082	3DEdge	$0.0917 \pm 0.0028$	$0.1226 \\ \pm 0.0044$	<b>0.0812</b> ±0.0040	0.0638	0.0564	0.0599
083 084	2DKeypoint	$0.0671 \pm 0.0038$	$0.0697 \\ \pm 0.0035$	<b>0.0626</b> ±0.0040	0.0337	0.289	0.306
085 086 087	3DKeypoint	$0.0512 \pm 0.0018$	$0.0670 \\ \pm 0.0027$	<b>0.0389</b> ±0.0014	0.0360	0.0298	0.0324
088 089	Reshading	$0.1308 \\ \pm 0.0058$	$0.1609 \\ \pm 0.0044$	<b>0.1284</b> ±0.0049	0.834	0.756	0.772
090 091	Curvature	$0.0413 \pm 0.0010$	$0.0498 \\ \pm 0.0019$	<b>0.0376</b> ±0.0023	0.0345	0.0291	0.329
092 093 094	Normal	$11.7850 \pm 0.4580$	$14.4381 \pm 0.3097$	<b>10.1346</b> ±0.0361	6.2418	5.7963	5.9821
095 096	SemSeg	$0.3980 \\ \pm 0.0350$	$0.3484 \pm 0.0308$	<b>0.4178</b> ±0.0361	0.4618	0.4905	0.4784
097 098 099	NYUDepth	$\begin{array}{c} 0.73 \\ \pm 0.09 \end{array}$	$0.49 \\ \pm 0.11$	<b>0.43</b> ±0.08	0.35	0.25	0.29
100 101	NYUNormal	$\begin{array}{c} 26.1 \\ \pm 3.8 \end{array}$	$18.5 \pm 1.7$	<b>16.4</b> ±1.6	18.2	14.8	14.9
102	-						



Figure 1: We present the impact of using different numbers of training samples on our method and VPD across all 12 tasks.