# Supplementary Material

In Section A, we provide more data from our linear SVM experiments (as outlined in Section 4.3) for classification performance under limited supervision. In Section B we derive the theoretical relationship between $f_\psi$ and $\xi$. In Section C, we provide ablation experiments to validate our choice of $J = 32$ in the main paper.

## A. Linear SVM with Limited Labels

In this section we provide additional information regarding the Linear SVM experiment from Section 4.3 for the case of limited labeled data. Our setup matches the setup used in [42]: We use a DGCNN backbone pre-trained on ShapeNet. We generate a set of smaller datasets given as a percentage of the original ModelNet40 training split. To generate a dataset, we first randomly sample one labeled point cloud from each class (ensuring there is at least one label of each category). Then we sample the labeled data uniformly randomly across the entire training set until the desired percentage of labeled data is reached with respect to the original dataset size. Note that since ModelNet40 has 9843 training objects, training an SVM on $1\%$ of labels would mean roughly 1-4 examples per class.
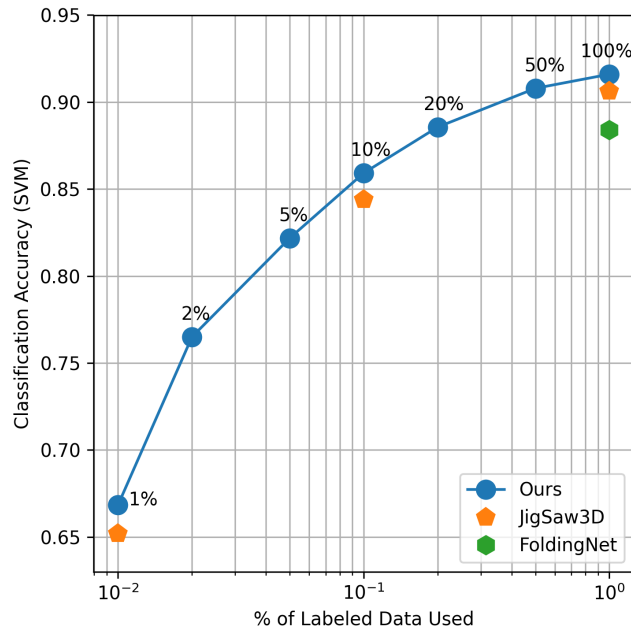


Figure 5. **Linear SVM Performance with Limited Labels:** With only $50\%$ of the training data, our method outperforms Jigsaw3D [42] trained on $100\%$ of the data (90.9 vs. 90.6). Additionally, our method outperforms FoldingNet [52] while using 5 times fewer labels (88.6 vs. 88.4).

Figure 5 shows the test accuracy (best of 3 runs) of a Linear SVM trained with varying amounts of labels, from $1 - 100\%$, using the frozen features of a pre-trained DGCNN. Surprisingly, our method trained on $50\%$ of the data outperforms Jigsaw3D [42] trained on $100\%$ of the data (90.9 vs. 90.6). In general, we gain a couple points in accuracy relative to Jigsaw3D, for example, at $1\%$ data (66.8 vs. 65.2) and $10\%$ data (85.9 vs. 84.4). For reference, we also place FoldingNet's SVM accuracy trained on $100\%$ of the data [52] on the same graph. Our method outperforms FoldingNet while using a fifth as many labels (88.6 vs. 88.4).

## B. Parametric Bottleneck as KL-Divergence Minimization

In this section, we provide more mathematical intuition for Equations 5, 6, and 7 in the main paper.

First, recall from Equation 1 that we redefined the logit matrix $\mathcal{S}$ of a segmentation network $f_\psi$ in terms of a set of joint log probabilities of points $\boldsymbol{p}_i \in \mathcal{P}$ and latent binary variables $c_{ij} \in \mathcal{C}$. We restate this equation here but with added notation showing the dependence on parameters $\boldsymbol{\psi}$,

$$\boldsymbol{s}_{ij} \stackrel{\text{def}}{=} \log p(\boldsymbol{p}_i, c_{ij} = 1; \boldsymbol{\psi}) \tag{12}$$

Under this interpretation, $f_{\psi}$ defines the log joint distribution over observed 3D points and latent binary correspondence variables (*i.e.* complete log likelihood), which we can marginalize over to compute the model evidence with respect to $\psi$,

$$p(\mathcal{P}; \psi) = \prod_{i=1}^{N} p(\boldsymbol{p}_i; \psi) = \prod_{i=1}^{N} \sum_{j=1}^{J} p(\boldsymbol{p}_i, c_{ij} = 1; \psi) = \prod_{i=1}^{N} \sum_{j=1}^{J} \exp(\boldsymbol{s}_{ij}) \tag{13}$$

Directly maximizing Equation 13 for $\psi$ would not produce a meaningful result. Instead, we maximize the data likelihood with respect to a parametric generative model (GMM) whose specific parameter set $\Theta$ can be analytically derived, as we will show in this section, as the result of an optimization problem.

Let $q_{\psi} \stackrel{\text{def}}{=} p(\boldsymbol{p}, \mathbf{c}; \psi) = \exp(f_{\psi})$ and $p_{\Theta} \stackrel{\text{def}}{=} p(\boldsymbol{p}, \mathbf{c}; \Theta)$. Our choice of $\xi$ functionally minimizes the KL-Divergence between $q_{\psi}$ and $p_{\Theta}$ with respect to $\Theta$. Restating in mathematical form,

$$\xi(\mathcal{P}, \mathcal{S}) = \operatorname*{argmin}_{\Theta} D_{KL}(q_{\psi} \,||\, p_{\Theta}) \tag{14}$$

Equation 14 provides additional theoretical insight into the "parametric bottleneck" concept: given a neural network parameterized by $\psi$ and a parametric model parameterized by $\Theta$, the purpose of $\xi$ is to squeeze $\psi$ into the specific set of parameters $\Theta$ of the parametric model such that the KL-Divergence between these two distributions is minimized. To prove that the specific form of $\xi$ solves the optimization in Equation 14, we first expand the KL-Divergence,

$$D_{KL}(q_{\psi} \,||\, p_{\Theta}) = -\int_{\boldsymbol{p}} \sum_{\mathbf{c}} p(\boldsymbol{p}, \mathbf{c}; \psi) \log \frac{p(\boldsymbol{p}, \mathbf{c}; \Theta)}{p(\boldsymbol{p}, \mathbf{c}; \psi)} d\boldsymbol{p} \tag{15}$$

$$= -H(q_{\psi}) - \int_{\boldsymbol{p}} \sum_{\mathbf{c}} p(\boldsymbol{p}_i, \mathbf{c}; \psi) \log p(\boldsymbol{p}_i, \mathbf{c}; \Theta) d\boldsymbol{p} \tag{16}$$

$$= -H(q_{\psi}) - \int_{\boldsymbol{p}} p(\boldsymbol{p}) \sum_{\mathbf{c}} p(\mathbf{c}|\boldsymbol{p}; \psi) \log p(\boldsymbol{p}_i, \mathbf{c}; \Theta) d\boldsymbol{p} \tag{17}$$

$$\approx -H(q_{\psi}) - \frac{1}{N} \sum_{\boldsymbol{p}_i \in \mathcal{P}} \sum_{\mathbf{c}} p(\mathbf{c}|\boldsymbol{p}_i; \psi) \log p(\boldsymbol{p}_i, \mathbf{c}; \Theta) \tag{18}$$

$H(\cdot)$ denotes the entropy functional and the approximation sign is due to the summation over $\boldsymbol{p}_i \in \mathcal{P}$. Since the entropy term doesn't depend on $\Theta$, it can be ignored when minimizing the KL-Divergence. The posterior $p(\mathbf{c}|\boldsymbol{p}; \psi)$ can be calculated via Bayes' rule (see Equations 2-4 in the main text). As in the main paper, we use the shorthand $\gamma_{ij} \stackrel{\text{def}}{=} p(c_{ij} = 1|\boldsymbol{p}_i; \psi)$.

$$\operatorname*{argmin}_{\Theta} D_{KL}(q_{\psi} \,||\, p_{\Theta}) = \operatorname*{argmax}_{\Theta} \sum_{\boldsymbol{p}_i \in \mathcal{P}} \sum_{\mathbf{c}} p(\mathbf{c}|\boldsymbol{p}_i; \psi) \log p(\boldsymbol{p}_i, \mathbf{c}; \Theta) \tag{19}$$

$$= \operatorname*{argmax}_{\Theta} \sum_{\boldsymbol{p}_i \in \mathcal{P}} \sum_{j=1}^{J} \gamma_{ij} \log \{p(\boldsymbol{p}_i|c_{ij} = 1; \Theta) p(c_{ij} = 1; \Theta)\} \tag{20}$$

$$= \operatorname*{argmax}_{\Theta} \sum_{\boldsymbol{p}_i \in \mathcal{P}} \sum_{j=1}^{J} \gamma_{ij} \{\log \mathcal{N}(\boldsymbol{p}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) + \log \pi_j\} \tag{21}$$

The optimization problem in Equation 21 shares the general underlying form of the M Step in a traditional Expectation Maximization (EM) algorithm [11] for optimizing the set of GMM parameters $\Theta$ given responsibilities $\gamma_{ij}$ computed in a previous E Step. Using the method of Lagrange multipliers to ensure $\sum_j^J \pi_j = 1$, the gradient of Equation 21 with respect to each $\boldsymbol{\mu}_j$, $\boldsymbol{\Sigma}_j$, and $\pi_j$ can be easily computed, set to zero, and solved, yielding the following analytical expressions,

$$\pi_j = \frac{1}{N} \sum_{i=1}^{N} \gamma_{ij} \tag{22}$$

$$\boldsymbol{\mu}_j = \frac{1}{N\pi_j} \sum_{i=1}^{N} \gamma_{ij}\boldsymbol{p}_i \tag{23}$$

$$\boldsymbol{\Sigma}_j = \frac{1}{N\pi_j} \sum_{i=1}^{N} \gamma_{ij}(\boldsymbol{p}_i - \boldsymbol{\mu}_j)(\boldsymbol{p}_i - \boldsymbol{\mu}_j)^T \tag{24}$$

Noting that $\gamma_{ij} = \frac{p(\boldsymbol{p}_i, c_{ij}=1; \boldsymbol{\psi})}{p(\boldsymbol{p}_i; \boldsymbol{\psi})} = \frac{\exp(\boldsymbol{s}_{ij})}{\sum_{j'}^{J} \exp(\boldsymbol{s}_{ij'})}$ via Equations 2-4 in the main text, we can substitute the latter expression for $\gamma_{ij}$ into Equations 22-24 to obtain the exact form of Equations 5, 6, and 7. Thus in computing these expressions, $\xi$ functionally maps the output of a point-wise classification network $f_\psi$ to the minimum KL-Divergence GMM parameter set $\Theta$, under the interpretation that $f_\psi$'s output $\mathcal{S}$ comprises joint likelihoods of observed and latent data.

## C. Ablation Experiments for $J$, the Number of Gaussian Components

On average, it seems that a larger $J$ relative to $N$ helps produce better representations. However, there are a few drawbacks: numerical stability, computational cost, and memory consumption during pre-training. For our ablation, we pre-trained on ShapeNet for various $J$, trained a linear SVM on the ModelNet40 training split using features learned during the pre-training stage, and report accuracy on the ModelNet40 test split. Times and memory consumption are reported from a single NVIDIA Titan RTX, with batch size equal to 32. Results are shown in Table 4. Larger $J$ incurs larger compute and memory costs for marginal downstream accuracy gains. We use $J = 32$ in the experiments in our paper as it performed the best while having reasonable memory and compute costs.

| $J$ | Memory Consumption | Time/Epoch (MM:SS) | Test Accuracy (SVM) |
|---|---|---|---|
| 2 | 2.1GB | 1:48 | 88.51 |
| 4 | 2.1GB | 2:09 | 89.00 |
| 8 | 2.1GB | 2:15 | 89.20 |
| 16 | 2.2GB | 2:21 | 89.37 |
| **32** | **2.2GB** | **3:30** | **90.30** |
| 64 | 2.3GB | 3:36 | 90.27 |
| 128 | 2.4GB | 5:01 | 89.89 |
| 256 | 3.4GB | 9:10 | 90.14 |
| 512 | 4.5GB | 16:32 | 89.04 |

Table 4. **Ablation Study for Various $J$**