

Figure 9: In experiments, we used a common feature-extractor (F), decoder (D), and predictor (P) network backbone while replacing different encoder heads (E).

440 7 Implementation details

441 Here, we include implementation details omitted from the main paper for brevity. Code for recreating
 442 the experiments in the paper is available at https://anonymous.4open.science/r/complexity_concepts-16A2. Upon acceptance, a deanonymized repository will be released.
 443

444 7.1 Pretraining

445 In pre-training (before the finetuning with a small number of examples on a cruder task), we used
 446 the following setup. In general, the overall network architecture comprised a feature extractor, an
 447 encoder head, a decoder, and a predictor, as depicted in Figure 9.

448 In all experiments, the predictor was parametrized as a two-layer feedforward neural network with
 449 hidden dimension 128 and a ReLU activation. The last layer’s dimension depended upon the exact
 450 prediction task (e.g., 10 neurons for FashionMNIST, 100 for CIFAR100, and 1010 for iNat) and used
 451 a softmax activation.

452 The feature extractors and decoders varied by domain. For FashionMNIST, the feature extractor
 453 used 3 2D convolution layers, followed by one fully connected layer. The decoder used two linear
 454 layers, followed by 3 inverse convolution layers. We again emphasize that the code for these models
 455 is available in our codebase, linked to at the beginning of this section.

456 For CIFAR100 and iNat, we pre-processed the images to extract the 512-dimensional activations
 457 from the penultimate layer of a ResNet18 pretrained on ImageNet [9]. These features were used as
 458 inputs to the feature extractor (x in Figure 9). For both CIFAR100 and iNat, the feature extractor used
 459 two linear layers, with a ReLU activation after the first layer, which had hidden dimension 128. The
 460 decoder was used to reconstruct the 512-dimension outputs of the ResNet18, using 3 fully-connected
 461 layers of dimension 128, 256, and 512, with ReLU activations between layers.

462 The different encoder heads were β -VAE, VQ-VIB $_C$, and VQ-VIB $_N$ models. β -VAE models used
 463 two linear layers, branching off the output of the feature extractor, to generate μ and σ from which
 464 to sample a continuous latent variable. VQ-VIB $_C$ directly passed the output of the feature extractor
 465 into the vector quantization layer, from which the discrete latent representations were sampled, as
 466 described in the main paper. In VQ-VIB $_N$, the output of the feature extractor was passed through two
 467 linear layers to generate a μ and a σ (exactly as in the β -VAE case) before the sampled continuous
 468 representation was discretized via vector quantization. Across experiments, the only differences
 469 among encoder heads that could arise were due to different latent dimensions (although we fixed it to
 470 32 for all experiments) or, for VQ-VIB $_C$ and VQ-VIB $_N$, the number of elements in the learnable
 471 codebook or n , the number of quantized vectors to combine into a latent representation.

472 In the main paper, we discussed the VQ-VIB $_C$ training loss (Equation 1), maximizing utility, mini-
 473 mizing reconstruction loss, and minimizing the entropy of the categorical distribution over codebook
 474 elements. A strict generalization of Equation 1, in which a variational bound on the complexity of
 475 representations is also penalized, is included in Equation 2:

$$\begin{aligned}
& \max \quad \lambda_U \mathbb{E}[U(x, y)] - \lambda_I \mathbb{E}[\|x - \hat{x}\|^2] \\
& - \lambda_H \mathbb{E} \left[\sum_{i \in [1, n]} H(\mathbb{P}(\zeta | h_i(x))) \right] \\
& - \lambda_C \mathbb{E} [D_{\text{KL}}[\mathbb{P}(\zeta | h(x)) \| \mathcal{U}(C)]] \\
& - \|\text{sg}[h_i(x)] - \zeta_i(x)\|^2 - \alpha \|h_i(x) - \text{sg}[\zeta_i(x)]\|^2
\end{aligned} \tag{2}$$

476 Equation 2 differs from Equation 1 via the third line, penalizing the KL divergence between the
477 conditional categorical distribution over codebook elements and a uniform distribution over the C
478 elements. This provides a variational bound on $I(X, Z)$, dubbed the complexity of representations
479 in prior literature [35, 27]. In our main experiments, we set $\lambda_C = 0$ and vary λ_H ; ablation studies
480 in which we varied λ_C instead of λ_H are included in Appendix 10 and confirm that controlling the
481 entropy of representations supported better finetuning accuracy.

482 In training β -VAEs, we trained to maximize the function described in Equation 3, where $\mu(x)$ and
483 $\sigma(x)$ represent the μ and σ parameters output by the encoder.

$$\begin{aligned}
& \max \quad \lambda_U \mathbb{E}[U(x, y)] - \lambda_I \mathbb{E}[\|x - \hat{x}\|^2] \\
& - \lambda_C \mathbb{E} [D_{\text{KL}}[\mathcal{N}(\mu(x), \sigma(x)) \| \mathcal{N}(0, 1)]]
\end{aligned} \tag{3}$$

484 Equation 3 trains agents to maximize classification accuracy, minimize MSE, and minimize the
485 complexity of representations. The scalar weight λ_C can be viewed as a Lagrange multiplier,
486 constraining how much information can be encoded in representations. This equation is closely
487 related to Equation 2 but, given the continuous nature of encodings in β -VAE, we could not penalize
488 the entropy of a categorical distribution.

489 In training VQ-VIB $_{\mathcal{N}}$, we used the training objective proposed by Tucker et al. [27], which closely
490 resembles the training loss for VQ-VIB $_C$ and is shown in Equation 4 (and closely matches Equation 2).
491 We use the same notation as for the β -VAE and VQ-VIB $_C$ models.

$$\begin{aligned}
& \max \quad \lambda_U \mathbb{E}[U(x, y)] - \lambda_I \mathbb{E}[\|x - \hat{x}\|^2] \\
& - \lambda_H \mathbb{E} \left[\hat{H}(\mathbb{P}(\zeta | \mu(x))) \right] \\
& - \lambda_C \mathbb{E} [D_{\text{KL}}[\mathcal{N}(\mu(x), \sigma(x)) \| \mathcal{N}(0, 1)]] \\
& - \|\text{sg}[h_i(x)] - \zeta_i(x)\|^2 - \alpha \|h_i(x) - \text{sg}[\zeta_i(x)]\|^2
\end{aligned} \tag{4}$$

492 The two key differences between Equation 4 and Equation 2 (used for training VQ-VIB $_C$) are bounds
493 on entropy and complexity (on the second and third lines of Equation 4). Just as for β -VAEs, VQ-
494 VIB $_{\mathcal{N}}$ models uses a KL divergence loss to regulate the complexity of representations. Increasing λ_C ,
495 as we did while annealing complexity in experiments, decreases the amount of encoded information.
496 However, as shown in our results, simply increasing λ_C does not ensure that VQ-VIB $_{\mathcal{N}}$ models will
497 use fewer discrete representations. (For visualizations of this effect, see Appendix 9.) Tucker et al.
498 [27] advocate for using a small positive λ_H to penalize the estimated entropy over codebook elements.
499 We explore varying λ_H in Appendix 10 and find some benefits relative to our main experiments, in
500 which we set $\lambda_H = 0$. However, given that VQ-VIB $_{\mathcal{N}}$ only supports an *approximation* of the entropy
501 term, we find that controlling the entropy for VQ-VIB $_{\mathcal{N}}$ is not as effective as controlling the entropy
502 for VQ-VIB $_C$.

503 7.2 Finetuning

504 In finetuning, we loaded pretrained frozen encoders and trained new predictor models to map from
505 encodings to downstream predictions.

506 For a finetuning task with M distinct classes, and a ‘‘duplication factor,’’ k , for how many examples
507 of each class to train with, we randomly selected $k * M$ datapoints to train with. For example, when

508 finetuning on the binary CIFAR100 task of living vs. non-living things, $M = 2$, so we loaded 2
509 total datapoints for $k = 1$, 4 datapoints for $k = 2$, etc.. For each input in the finetuning dataset,
510 we generated an encoding by passing through the encoder once. This generated a new dataset of
511 encodings and labels, which we used to train the predictor. (Note that this approach is distinct
512 from passing the input through the encoder many times; given stochastic encoders, which we used,
513 the same input could result in many different encodings. Here, we assumed a limited budget of
514 *encodings*.)

515 Predictor neural networks were instantiated as feedforward networks with 4 fully connected layers,
516 with hidden dimension 256 and ReLU activations, and trained to map from shifted encodings (see
517 next paragraph) to classifications for 100 epochs using an Adam optimizer with default parameters,
518 with the learning rate decreasing by a factor of 10 based on plateauing training loss, with a patience
519 of 5 epochs, and early stopping if the learning rate fell below 10^{-8} .

520 One particular design choice that we made in finetuning predictors merits elaboration: shifting
521 encodings. Rather than directly training predictors to map from encodings to predictions, we applied
522 a simple linear transformation to the encodings before feeding them to the predictor. Specifically, we
523 multiplied all tensor elements by 5 and increased them by 1. This simple linear transformation does
524 not affect any relations between encodings except scale, and indeed we found that predictors could
525 be successfully trained with this rescaling. Nevertheless, this linear transformation was important to
526 provide a check against merely relying upon initialization conditions for good finetuning performance.
527 In particular, we found that if we did not apply this linear transformation (i.e., pass the raw encodings
528 to the predictor), predictors sometimes performed better than they should given the training data. For
529 example, as a sanity check, we trained a predictor on a binary classification task, but only provided
530 two positive datapoints and no negative data. In general, given this data, one would expect a trained
531 predictor to only predict positive labels. However, we observed that in several cases, the predictor
532 would achieve nearly perfect accuracy, including predicting negative labels for negative inputs.
533 This surprising result disappeared when we simply shifted encodings, indicating that the particular
534 initialization conditions of the predictor seemed to align well with pre-trained encoders. We wanted
535 to measure the effect of data on finetuning performance, rather than just initialization conditions, but
536 we note that this odd phenomenon of well-aligned initializations merits further investigation.

537 We ran 10 finetuning trials per model, which was important given the small amount of randomly-
538 sampled finetuning data.

539 7.3 Hyperparameters

540 In the following subsections, we present the hyperparameters used for training different encoders in
541 the different domains. In general, we used the following principles when choosing hyperparameters:

- 542 • For VQ-based methods, use a large enough codebook to have at least one element per class.
543 Larger C are also acceptable, as tuning weights should decrease the effective codebook size.
- 544 • When annealing, use a small enough weight increment to generate smooth changes during
545 training. Larger increments, however, speed up training.
- 546 • When annealing for larger n one can increase the annealing rate. Models with greater n
547 tended to use more complex representations, so annealing could be extremely slow for small
548 increments.

549 7.3.1 FashionMNIST

550 For FashionMNIST, we trained all models with batch size 64 for 200 epochs, using the hyperparam-
551 eters specified in Table 1. The only differences across methods were which hyperparameters we
552 annealed to penalize complexity. Other differences simply reflected differences in architecture (e.g.,
553 using a codebook for vector-quantization methods). Pre-training a single model for 200 epochs took
554 approximately 5 minutes on a desktop computer with one NVIDIA 2080 GeForce RTX.

555 7.3.2 CIFAR100

556 For CIFAR100, we trained all models with batch size 256 for 400 epochs, using the hyperparameters
557 specified in Table 2. As explained for FashionMNIST, the only substantial differences across

Table 1: Hyperparameters for FashionMNIST training.

ENCODER	LATENT DIM	C	n	λ_U	λ_I	λ_{C_0}	λ_C INCR	λ_{H_0}	λ_H INCR
β -VAE	32	NA	NA	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{N}	32	1000	1	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{N}	32	1000	2	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{N}	32	1000	4	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{C}	32	1000	1	10	10	0.0	0.0	0.001	0.2
VQ-VIB \mathcal{C}	32	1000	2	10	10	0.0	0.0	0.001	0.4
VQ-VIB \mathcal{C}	32	1000	4	10	10	0.0	0.0	0.001	0.8

Table 2: Hyperparameters for CIFAR100 training.

ENCODER	LATENT DIM	C	n	λ_U	λ_I	λ_{C_0}	λ_C INCR	λ_{H_0}	λ_H INCR
β -VAE	32	NA	NA	10	10	0.01	0.1	0.0	0.0
VQ-VIB \mathcal{N}	32	1000	1	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{N}	32	1000	2	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{N}	32	1000	4	10	10	0.01	0.5	0.0	0.0
VQ-VIB \mathcal{C}	32	1000	1	10	10	0.0	0.0	0.001	0.04
VQ-VIB \mathcal{C}	32	1000	2	10	10	0.0	0.0	0.001	0.08
VQ-VIB \mathcal{C}	32	1000	4	10	10	0.0	0.0	0.001	0.12

558 architectures were architecture-specific terms that needed to be specified, and which terms were
 559 annealed to penalize complexity. Pre-training a single model for 400 epochs took approximately 10
 560 minutes on a desktop computer with one NVIDIA 3080.

561 7.3.3 iNaturalist

562 For iNat, we trained all models with batch size 256, using the hyperparameters specified in Table 3.
 563 We trained β -VAE and VQ-VIB \mathcal{C} models for 300 epochs, while we trained VQ-VIB \mathcal{N} models for
 564 600. We used more annealing epochs for VQ-VIB \mathcal{N} simply because it seemed to need more epochs
 565 to eventually anneal to random chance. Likely, a larger annealing rate could also accomplish the
 566 desired effect, but initial experiments with faster annealing tended to induce rapid codebook collapse
 567 that did not generate the smooth spectrum of MSE values we desired. Pre-training a single model
 568 for 300 epochs took approximately 10 minutes on a desktop computer with one NVIDIA 3080 (and
 569 twice as long for 600 epochs).

570 8 Further Finetuning Results

571 In the main paper, we included only a small number of graphs highlighting our key results. Here,
 572 we include further results that corroborate the main trends stated in the paper. Primarily, these plots
 573 include further experiments for varying the amount of finetuning data, as well as varying n , the
 574 number of codebook elements to combine into a latent representation. Results are divided by domain:
 575 FashionMNIST, CIFAR100, and iNat.

576 8.1 FashionMNIST

577 Here, we include the finetuning results for the FashionMNIST domain for varying amounts of
 578 finetuning data, ranging over $k \in [1, 2, 5, 10, 50]$, and n , the number of quantized vectors to combine
 579 into a single representation. Results for each k are included in Figure 10.

580 As expected, increasing the amount of finetuning data improved performance for all models, and the
 581 gap between all model types (VQ-VIB \mathcal{C} , VQ-VIB \mathcal{N} , and β -VAE) shrank. It is noteworthy, however,
 582 that a VQ-VIB \mathcal{C} model, tuned to the right complexity level and trained with just one example per
 583 ternary class (Figure 10 a), achieved better accuracy than a β -VAE model trained with 50 examples
 584 per class (Figure 10 e). Further, for any fixed k , VQ-VIB \mathcal{C} consistently outperformed VQ-VIB \mathcal{N} ,
 585 suggesting that many recent works that use VQ-VIB \mathcal{N} could be improved by replacing the model

Table 3: Hyperparameters for iNaturalist training.

ENCODER	LATENT DIM	C	n	λ_U	λ_I	λ_{C_0}	λ_C INCR	λ_{H_0}	λ_H INCR
β -VAE	32	NA	NA	10	10	0.0001	0.03	0.0	0.0
VQ-VIB $_{\mathcal{N}}$	32	2000	1	10	10	0.001	0.1	0.0	0.0
VQ-VIB $_{\mathcal{N}}$	32	2000	2	10	10	0.001	0.2	0.0	0.0
VQ-VIB $_{\mathcal{N}}$	32	2000	4	10	10	0.001	0.4	0.0	0.0
VQ-VIB $_C$	32	2000	1	10	10	0.0	0.0	0.001	0.05
VQ-VIB $_C$	32	2000	2	10	10	0.0	0.0	0.001	0.10
VQ-VIB $_C$	32	2000	4	10	10	0.0	0.0	0.001	0.20

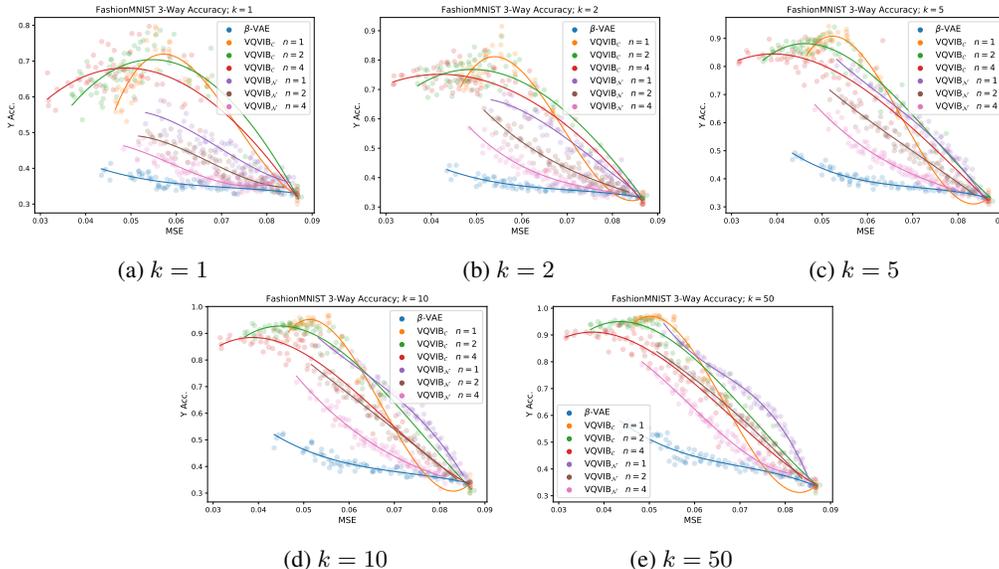


Figure 10: FashionMNIST finetuning results for varying k . As k increased, all models benefited. The data-efficiency of advantage of VQ-VIB $_C$ was most pronounced when using the least amount of data.

586 type [27, 28, 13, 8]. Lastly, for both VQ-VIB $_{\mathcal{N}}$ and VQ-VIB $_C$, increasing n tended to support
 587 lower MSE but worse finetuning accuracy. This supports an intuition that combining more discrete
 588 representations starts to more densely fill the representation space, trending towards continuous
 589 representations.

590 We note briefly that VQ-VIB $_{\mathcal{N}}$, both in this domain and others (explored in the next sections),
 591 typically failed to learn as complex representations as either VQ-VIB $_C$ or β -VAEs. This is apparent
 592 given the limited range of MSE values for the VQ-VIB $_{\mathcal{N}}$ curves. We consistently struggled to
 593 make VQ-VIB $_{\mathcal{N}}$ learn as rich representations as for the other model types, which led to worse
 594 reconstructions and higher MSE values.

595 8.2 CIFAR100

596 We found similar trends in the CIFAR100 to those in the FashionMNIST domain and plotted results
 597 in Figures 11 and 12 (for 2-way and 20-way finetuning tasks, respectively). In all experiments,
 598 VQ-VIB $_C$ outperformed both β -VAE and VQ-VIB $_{\mathcal{N}}$. In the 2-way finetuning example, we again
 599 found a peaked curve for VQ-VIB $_C$ finetuning accuracy as a function of MSE, indicating that tuning
 600 to the right complexity level induced the best accuracy. In the more complex 20-way classification
 601 task, however, we did not observe this peak.

602 This last result is unsurprising: the 20-way hierarchy in CIFAR100 is less semantically meaningful
 603 and likely less obvious in photos than the 2-way task of distinguishing living and non-living things.
 604 For example, two of the 20 categories are simply different sorts of vehicles. It would be extremely
 605 surprising for VQ-VIB $_C$ to learn such arbitrary groups automatically while compressing represen-

606 tations. Without learning the right groupings, VQ-VIB_C cannot benefit from learning less complex
 607 representations.

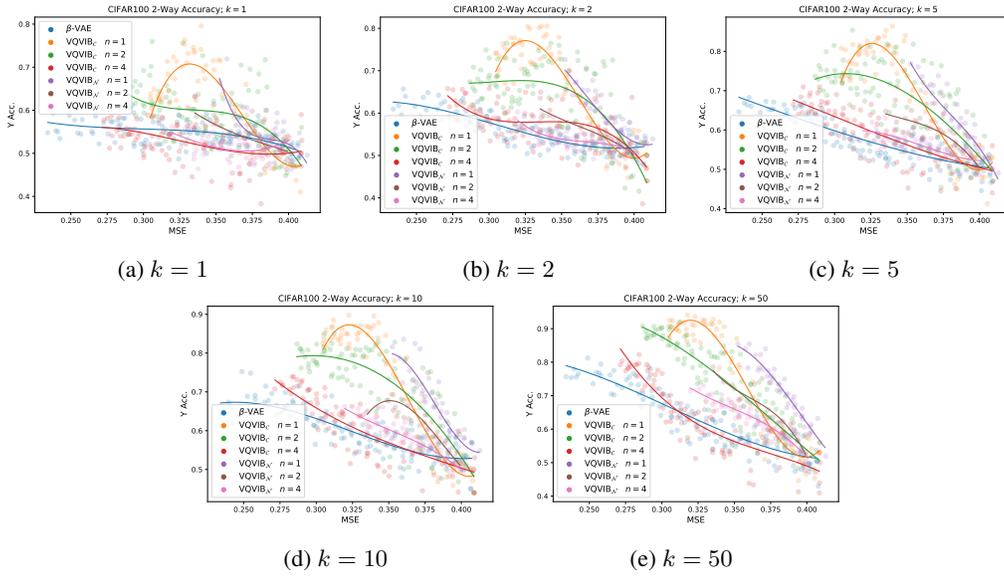


Figure 11: CIFAR100 2-way finetuning results for varying k .

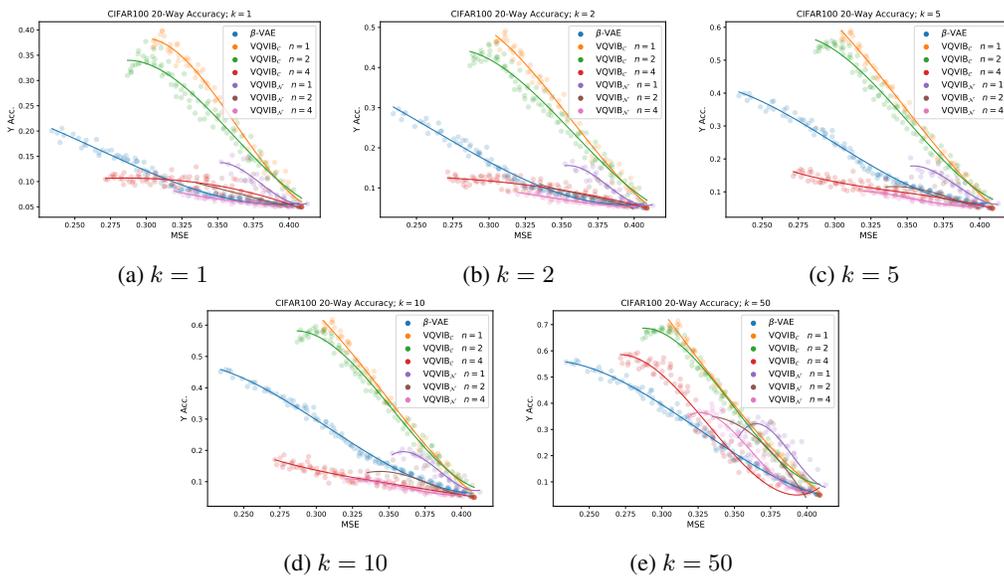


Figure 12: CIFAR100 20-way finetuning results for varying k .

608 **8.3 iNaturalist**

609 Lastly, we found similar trends in finetuning in the iNat domain, finetuned on a 3-way (Figure 13),
 610 34-way (Figure 14), and 1010-way (Figure 15) finetuning task.

611 On the 3-way finetuning task (between animals, plants, and fungi), we observed similar peaking
 612 behavior as in earlier experiments, indicating yet again the importance of tuning to the right complex-
 613 ity. In addition, as in prior results, we found a similar trend that greater n tended to allow greater
 614 complexity (lower MSE) but induced worse finetuning performance. For example, in Figure 13 b, the
 615 orange line, corresponding to $n = 1$ stays above and to the right of the green ($n = 2$) and red ($n = 4$)
 616 lines. Intuitively, this seems to indicate that the more combinatorial representations, with greater
 617 n , were somewhat of a midpoint between the continuous β -VAE representations and the discrete
 618 representations used by $VQ-VIB_C$ for $n = 1$.

619 Results from the 34-way finetuning followed similar patterns as before as well. Just as in CIFAR100
 620 wherein we tested both a 2-way and 20-way finetuning task, this 34-way finetuning task for iNat
 621 showed that $VQ-VIB_C$ continued to outperform $VQ-VIB_N$ and β -VAE for more complex finetuning
 622 tasks, although the performance gap shrank as k increased.

623 Most interestingly, perhaps, we conducted yet another iNat finetuning experiment, this time using
 624 the 1010 low-level labels that had originally been used during pre-training. As before, we used very
 625 small amounts of data in finetuning (e.g., for $k = 1$, only 1 example from each class, so 1010 labeled
 626 examples total). Results from those experiments are shown in Figure 15.

627 For small k , we again see that $VQ-VIB_C$ outperforms other model types. For larger k , however,
 628 we see one of the limitations of $VQ-VIB_C$. Because the discrete encoders learned less complex
 629 representations than β -VAEs (as shown by the fact that they never reach lower MSE values), with
 630 enough finetuning data, β -VAEs are able to capture distinctions between classes that $VQ-VIB_C$ models
 631 cannot. Thus, in the particular case of large amounts of finetuning data and complex finetuning tasks,
 632 more complex, continuous encoders continue to outperform our method.

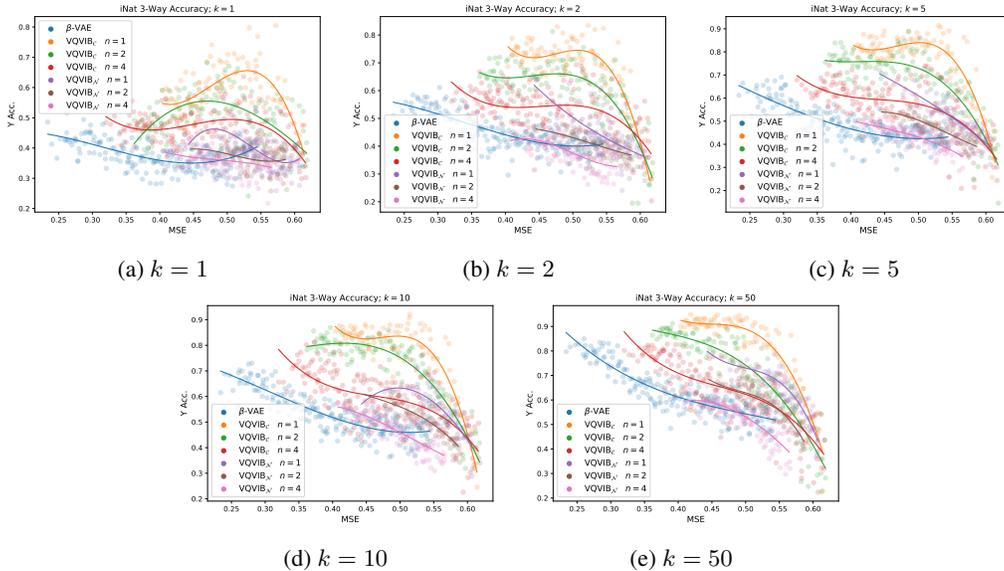


Figure 13: iNat 3-way finetuning results for varying k .

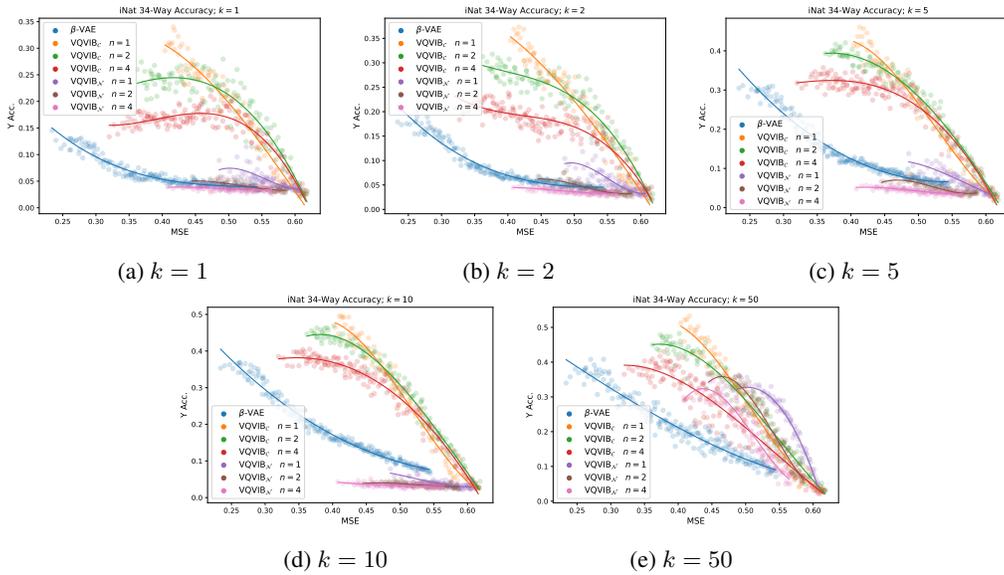


Figure 14: iNat 34-way finetuning results for varying k .

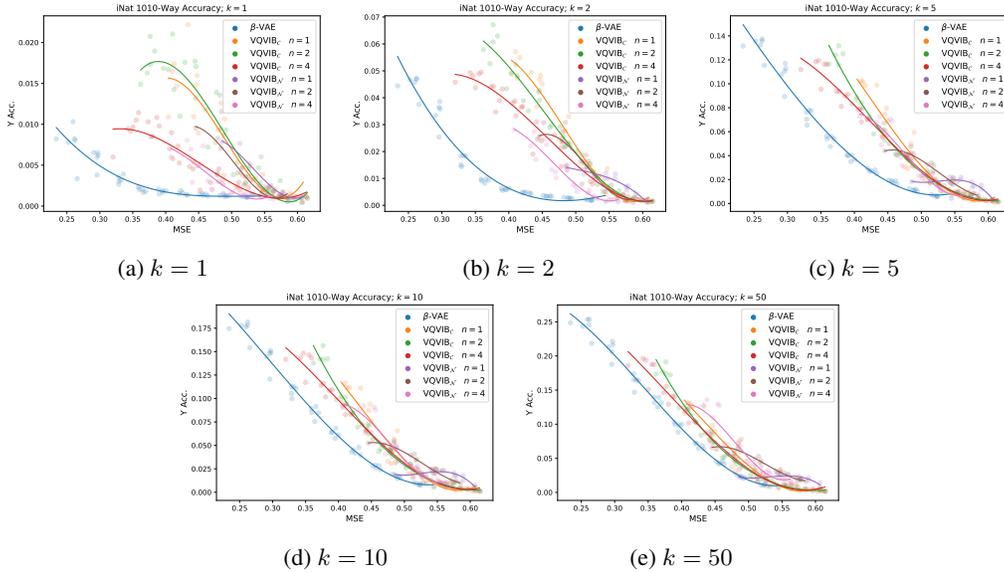


Figure 15: iNat 1010-way finetuning results for varying k .

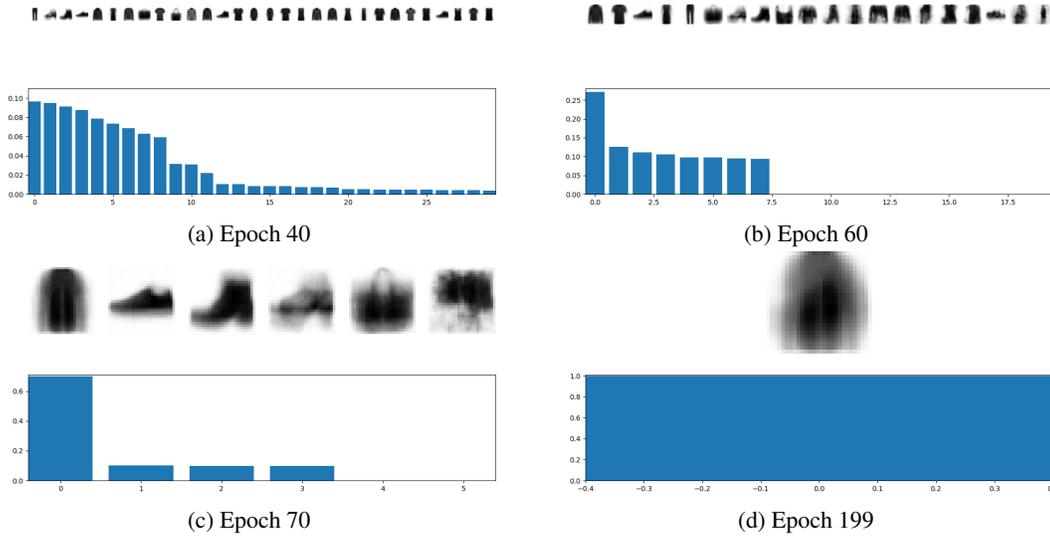


Figure 16: The evolution of the distribution over prototypes during annealing for VQ-VIB_C in the FashionMNIST domain. In early epochs, VQ-VIB_C uses many prototypes, with a long-tailed distribution. Over the course of annealing the entropy, the probability distribution becomes more concentrated (b and c) before collapsing to a single prototype (d).

633 9 Prototype Utilization: Further Visualizations

634 Here, we include some further visualizations that we omitted from the main paper due to space
 635 constraints. These visualization primarily illustrate the importance of entropy-regulated representation
 636 learning (for VQ-VIB_C) vs. complexity-regulated (for VQ-VIB_N).

637 Figures 16 and 17 shows the prototypes for VQ-VIB_C and VQ-VIB_N, respectively, in the FashionM-
 638 NIST domain over the course of training. Each subfigure consists of a top row of decoded prototypes,
 639 with associated probabilities (frequency of use measured when passing through images from the test
 640 set) below. The 30 most frequent prototypes are visualized, or fewer prototypes if fewer were used.

641 There is an important trend in Figures 16 and 17: the entropy-based annealing for VQ-VIB_C caused
 642 models to use fewer prototypes, while the complexity-based annealing for VQ-VIB_N did not. At
 643 epoch 40, just as both methods begin annealing, VQ-VIB_C and VQ-VIB_N use a large number of
 644 prototypes, as seen by the long-tailed distributions. Over the course of annealing, however, VQ-VIB_C
 645 uses fewer prototypes, and merges images of different classes into the same prototype. Thus, the
 646 degenerate encoder at the end of annealing (epoch 199) uses just a single prototype to represent all
 647 possible inputs (Figure 16). At the same time, VQ-VIB_N, during annealing, does not use fewer
 648 prototypes. Rather, the complexity-penalization term seems to induce the model to make the mapping
 649 from input to prototype more stochastic (Figure 17). Thus, the degenerate VQ-VIB_N encoder uses
 650 many prototypes, each of which is blurry because it could correspond to any input.

651 Visualizations of decoded prototypes for the CIFAR100 domain is more challenging. In richer image
 652 domains, prototype-based methods often use training examples as prototypes [3, 20, 4], which can
 653 make it more difficult to understand when a single prototype represents more than one concept.
 654 Nevertheless, by visualizing the distribution over prototypes (without decoding them), we see the
 655 same pattern that VQ-VIB_C tends to learn to use fewer prototypes over the course of annealing than
 656 VQ-VIB_N. Snapshots of the categorical distributions for CIFAR100 are included in Figure 18.

657 10 Ablation Study: Entropy vs. Complexity

658 Here, we present results motivating penalizing entropy, as opposed to complexity, in VQ-VIB_C.
 659 Appendix 9 showed how annealing entropy in VQ-VIB_C caused models to use fewer prototypes,
 660 whereas penalizing complexity in VQ-VIB_N did not induce similar reductions in effective codebook

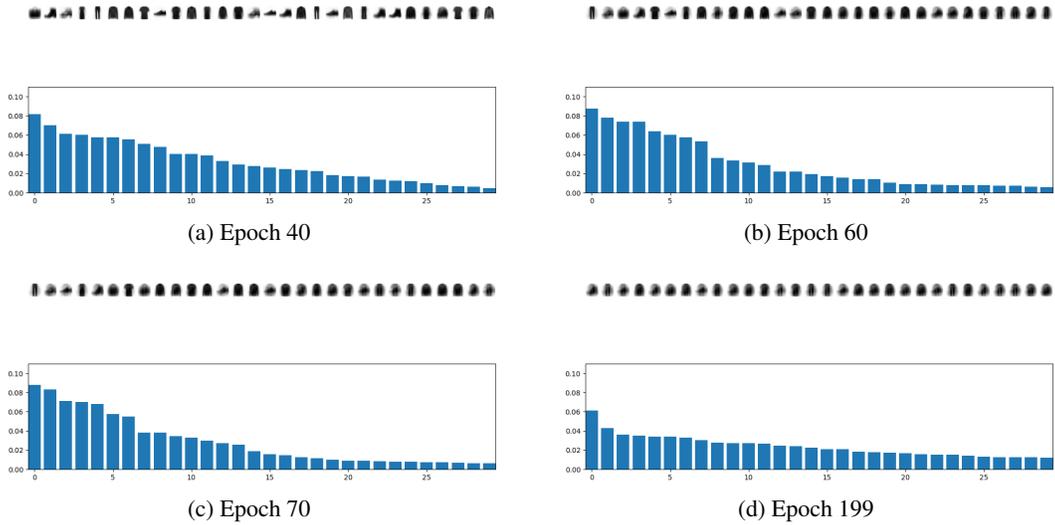


Figure 17: The evolution of the distribution over prototypes during annealing for VQ-VIB_N in the FashionMNIST domain. Unlike annealing entropy for VQ-VIB_C, annealing the complexity in VQ-VIB did not result in fewer prototypes being used. Instead, each prototype became blurrier, indicating that each prototype became more likely regardless of class.

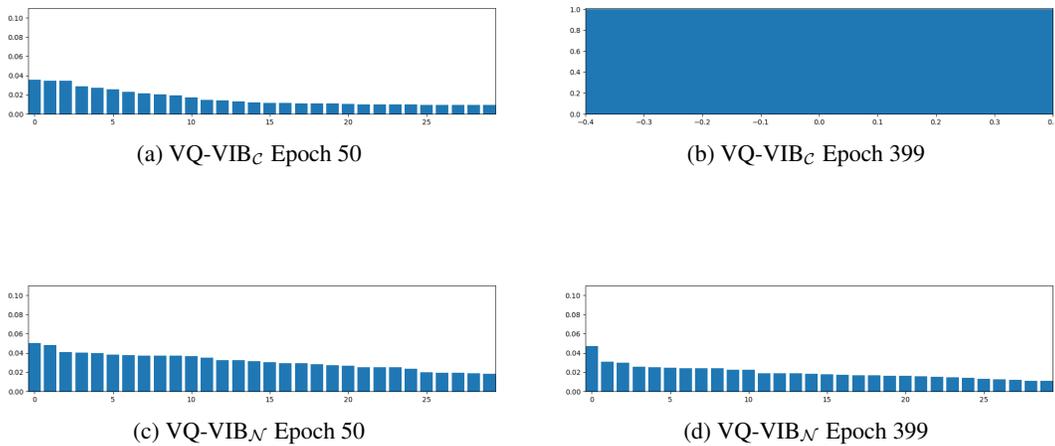


Figure 18: Categorical distribution over prototypes while annealing in the CIFAR100 domain at the start of annealing (Epoch 50) and at the end (Epoch 399) for VQ-VIB_C (top row) and VQ-VIB_N (bottom row). Annealing the entropy in VQ-VIB_C caused the model to use fewer prototypes (note the degenerate categorical distribution over only one prototype at Epoch 399), whereas annealing complexity for VQ-VIB_N did not cause similar concentration.

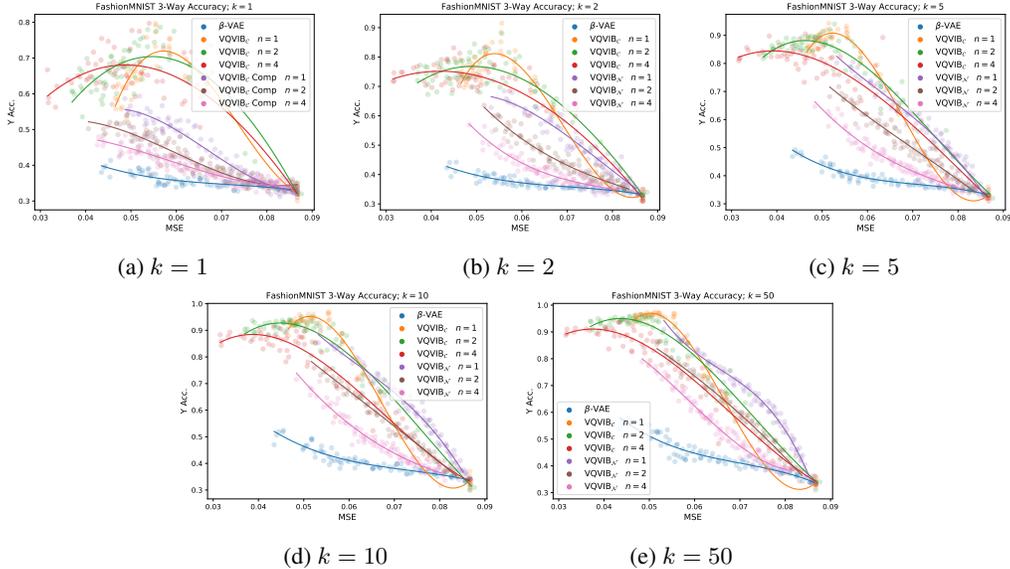


Figure 19: FashionMNIST finetuning results for varying k , comparing annealing by entropy (VQ-VIB $_C$) and annealing by complexity (VQ-VIB $_C$ Comp.). Annealing by complexity resulted in worse finetuning performance.

661 size. Further experiments corroborate our findings that penalizing entropy was the key to this
 662 difference in behavior.

663 We trained VQ-VIB $_C$ agents on the FashionMNIST task, using the same pre-training and finetuning
 664 procedures as in the main paper, with the only difference being that we annealed the complexity of
 665 representations instead of the entropy. Results from finetuning such models are included in Figure 19.

666 Figure 19 shows that annealing by entropy, as opposed to complexity, was the key factor in improving
 667 VQ-VIB $_C$ finetuning performance. The difference in performance when penalizing entropy vs. com-
 668 plexity closely matches the difference in performance between VQ-VIB $_C$ and VQ-VIB $_N$ examined
 669 in the main paper. Thus, the entropy-regularization term seems to explain much of the difference
 670 between VQ-VIB $_C$ and VQ-VIB $_N$.

671 In subsequent experiments in the CIFAR100 and iNat domains, therefore, we tested whether penaliz-
 672 ing the estimated entropy of VQ-VIB $_N$ models matched VQ-VIB $_C$ results from the main paper. We
 673 note that Tucker et al. [27] advocate for a small positive λ_H to penalize entropy, but the authors also
 674 acknowledge that exactly computing this entropy is impossible given the VQ-VIB $_N$ architecture.

675 Finetuning results for CIFAR100 (on the 2-way and 20-way finetuning tasks) and iNat (on the 3-way
 676 and 34-way finetuning tasks), for VQ-VIB $_C$ trained by varying λ_C and VQ-VIB $_N$ trained by varying
 677 λ_H , are included in Figures 20, 21, 22, and 23. Several important trends emerge from viewing these
 678 plots, especially compared to results from our main paper for VQ-VIB $_C$ controlled via λ_H .

679 First, finetuning performance is noisier using these models compared to results from the main text.
 680 This likely arises, for VQ-VIB $_N$ models, because increasing λ_H failed to consistently reduce the
 681 number of discrete representations used. Thus, for a given MSE value, different models used different
 682 numbers of representations, and therefore exhibited different finetuning performance.

683 Second, varying λ_H , instead of λ_C , seemed to somewhat improve VQ-VIB $_N$ performance, but
 684 not as much as when varying λ_H for VQ-VIB $_C$, as presented in our main paper. For example,
 685 consider Figure 21 a. The best-performing model, VQ-VIB $_N$, $n = 1$, peaks at finetuning accuracy of
 686 approximately 0.16, outperforming VQ-VIB $_C$ models when varying λ_C . However, in Figure 12 a, we
 687 found that VQ-VIB $_C$ models in the exact same setting achieved a mean accuracy of approximately
 688 0.38: more than double the VQ-VIB $_N$ performance. Thus, varying λ_H seemed to improve VQ-VIB $_N$
 689 performance somewhat, but VQ-VIB $_C$ better supports penalizing entropy, and therefore achieves
 690 higher performance.

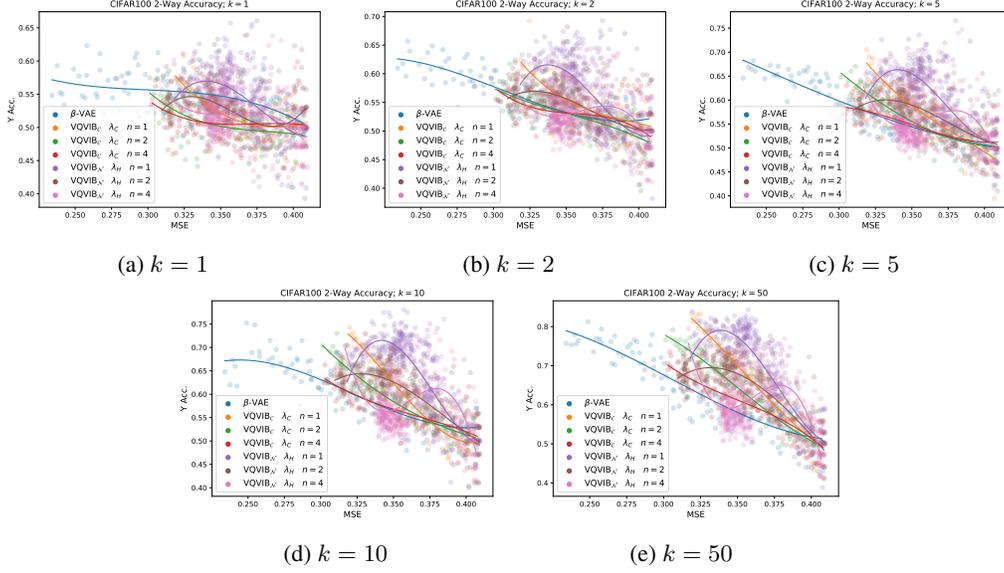


Figure 20: Ablation study results for the CIFAR100 2-way finetuning task. Tuning λ_C for VQ-VIB $_C$ or λ_H for VQ-VIB $_N$ led to worse results than in our main paper.

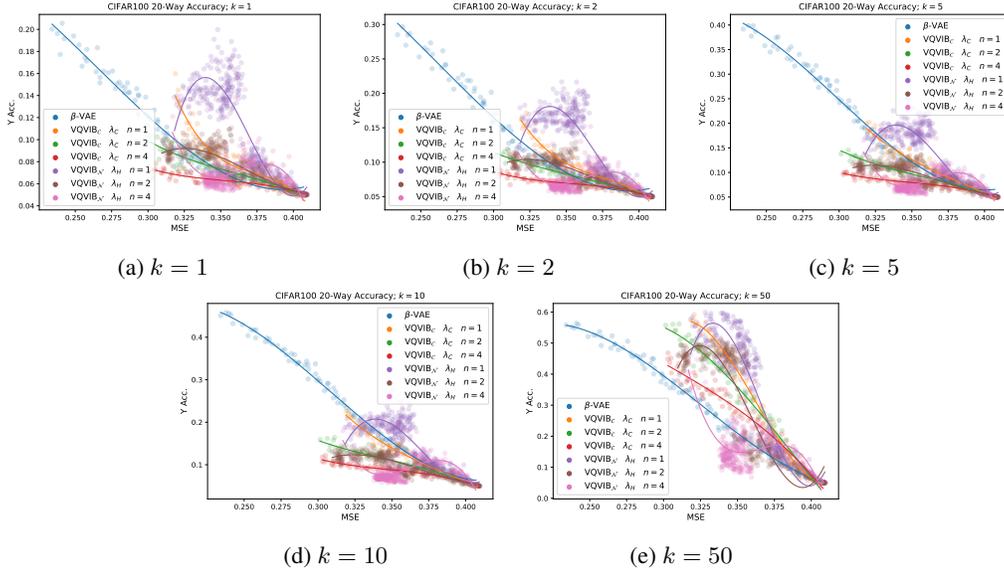


Figure 21: Ablation study results for the CIFAR100 20-way finetuning task. Tuning λ_C for VQ-VIB $_C$ or λ_H for VQ-VIB $_N$ led to worse results than in our main paper.

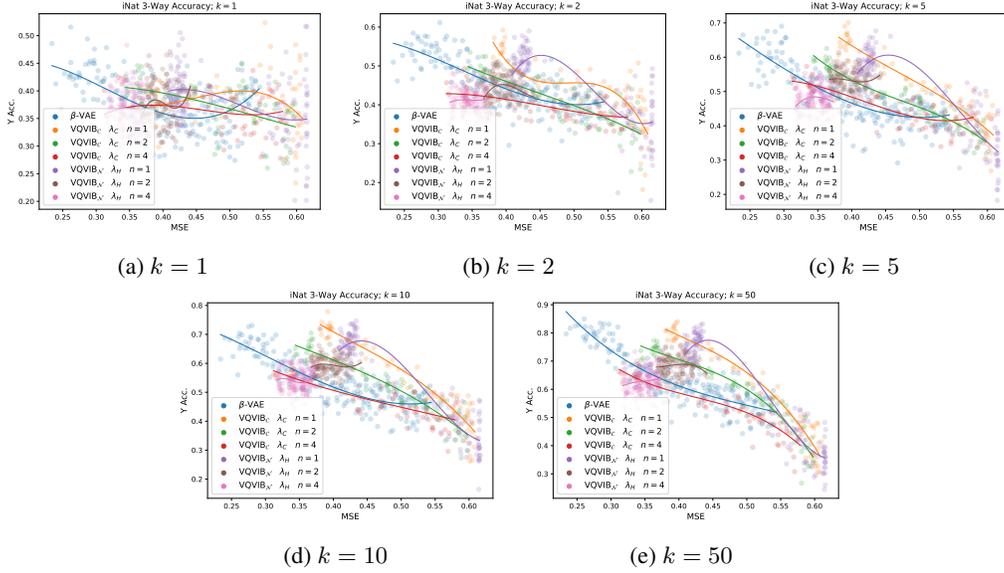


Figure 22: Ablation study results for the iNat 3-way finetuning task. Tuning λ_C for VQ-VIB $_C$ or λ_H for VQ-VIB $_N$ led to worse results than in our main paper.

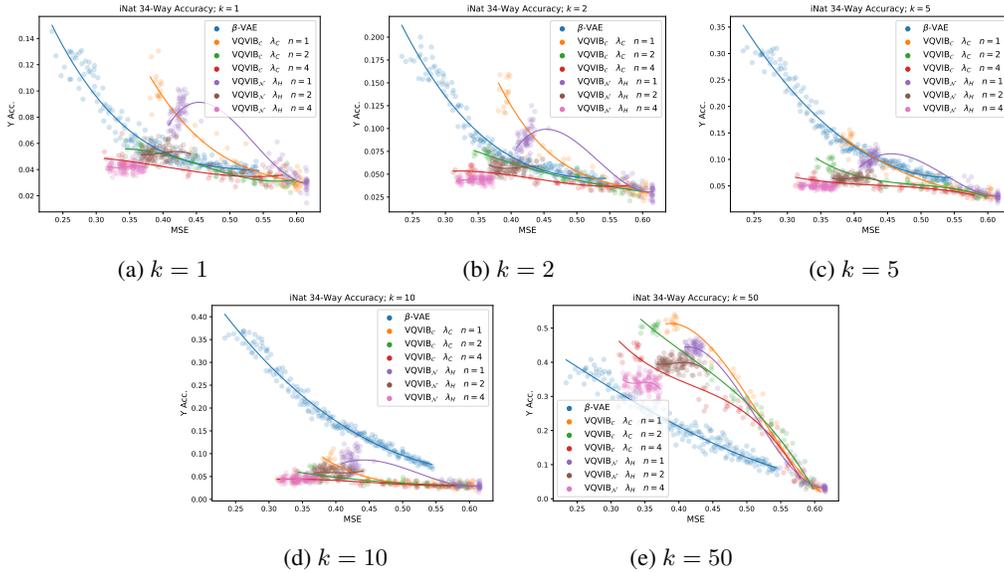


Figure 23: Ablation study results for the iNat 34-way finetuning task. Tuning λ_C for VQ-VIB $_C$ or λ_H for VQ-VIB $_N$ led to worse results than in our main paper.

691 Third, varying λ_C , instead of λ_H , for VQ-VIB_C worsened finetuning performance. Once again,
692 by comparing finetuning performance for VQ-VIB_C models in Figure 21 a (achieving a maximum
693 accuracy around 0.14), to results from our main paper, we note the importance of penalizing the
694 entropy of representations.

695 Thus, in general these ablation studies support many of the design decisions made in the main paper.

- 696 1. Varying λ_H , instead of λ_C for VQ-VIB_C improves finetuning performance by decreasing
697 the number of discrete representations used.
- 698 2. VQ-VIB_N benefits somewhat from penalizing entropy, but because it is architecturally
699 unable to support exact calculations of entropy, we were unable to match VQ-VIB_N perfor-
700 mance.

701 It is certainly possible that some optimal combination of λ_H and λ_C might further improve VQ-VIB_N
702 or VQ-VIB_C performance; initial explorations of such combinations with fixed λ_H values while
703 annealing λ_C did not yield obvious results. Most importantly, our current findings are enough to
704 indicate that controlling the entropy of discrete representations appears important for data-efficient
705 finetuning.

706 **11 User Study**

707 In the subsequent pages, we have included the exact pdf document shared with participants of the
708 user study, edited only to preserve anonymity during peer-review. Following standard user study
709 procedure, we initial briefed users by telling them how long the study was and that they were free
710 to leave at anytime. Demographic information was collected in person. After the study, users were
711 debriefed and given the email of the study designer to contact if they had any questions.

Block 1

Brief

You are free to leave this study at any time.

It will take less than 5mins.

Thank you for your participation!

Block 2

Introduction

A factory has created a robot that is good at sorting clothes into these 10 categories

- T-shirt
- Trouser
- Pullover
- Dress
- Coat
- Sandal
- Shirt
- Sneaker
- Bag
- Ankleboot

See below for examples of these 10 categories



Block 3

Introduction

However, in your home you don't care about these 10 categories.

You only sort your clothes into these 3 categories:

- Shirts
- Shoes
- Trousers/Bags

So, we want to teach the robot to be more general and only care about these three high level categories.

Block 10

Introduction

In doing so, we are going to re-train ***two different robots*** over a period of time.

At a certain point during this re-training, they will be able to sort these clothes into these three new categories best.

Your task is to pick the point in which they are going to categorize these three high-level categories best.

Block 4

Introduction

You will have to answer one question about each robot.

First, you will be shown one robot which only communicates with its error score. *The lower the score, the better the robot is at classifying the 10 categories.*

Secondly, you will be shown another robot. This robot has learned to communicate what it has learned through visualizing the most important images it uses for categorizing the 3 high-level concepts. You have to decide which visualization corresponds to the robot being able to categorize the three new high-level categories best.

Generally, the more images the second robot is using, the less general it will be, and the worse it will do at categorizing the 3 new high-level categories.

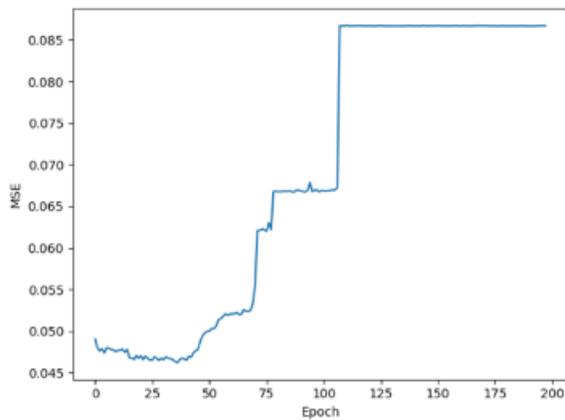
Block 5

Robot 1

Here is an example of the first robot communicating its error scores over training.

You will have to pick a score from 3 sampled options you think will perform best on your desired task of learning the 3 clothing concepts. ***A lower score is typically seen as better (which corresponds to a lower point on the blue curve below).***

Your task is to pick the point which you think will perform best at representing these three categories.



Block 6

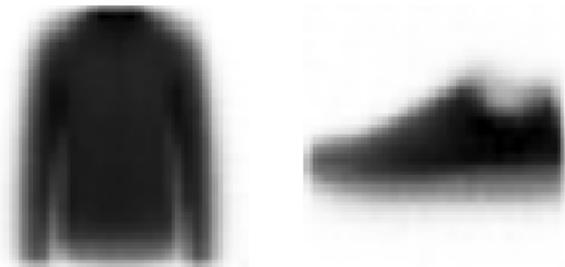
Robot 2

Then you will be shown visualizations like this, which is second robot's method of communicating its learned concepts.

If the visualization shows that (1) the robot is not using many images, and (2) they roughly represent the three high-level categories (shirts, shoes, trousers/bags), then the robot should perform well in the task of sorting clothes into these three categories.

Your task is to pick the visualization which you think will perform best at sorting the clothes into these 3 high-level categories.

For example, here is one robot's visualisation which is perhaps too general, as there is not even three categories being used.



And here is one which could be too specific, because there is a lot of detail.



Important: We want the robot to learn the three high-level concepts, but no more or no less!

Block 11

Click next to begin the study

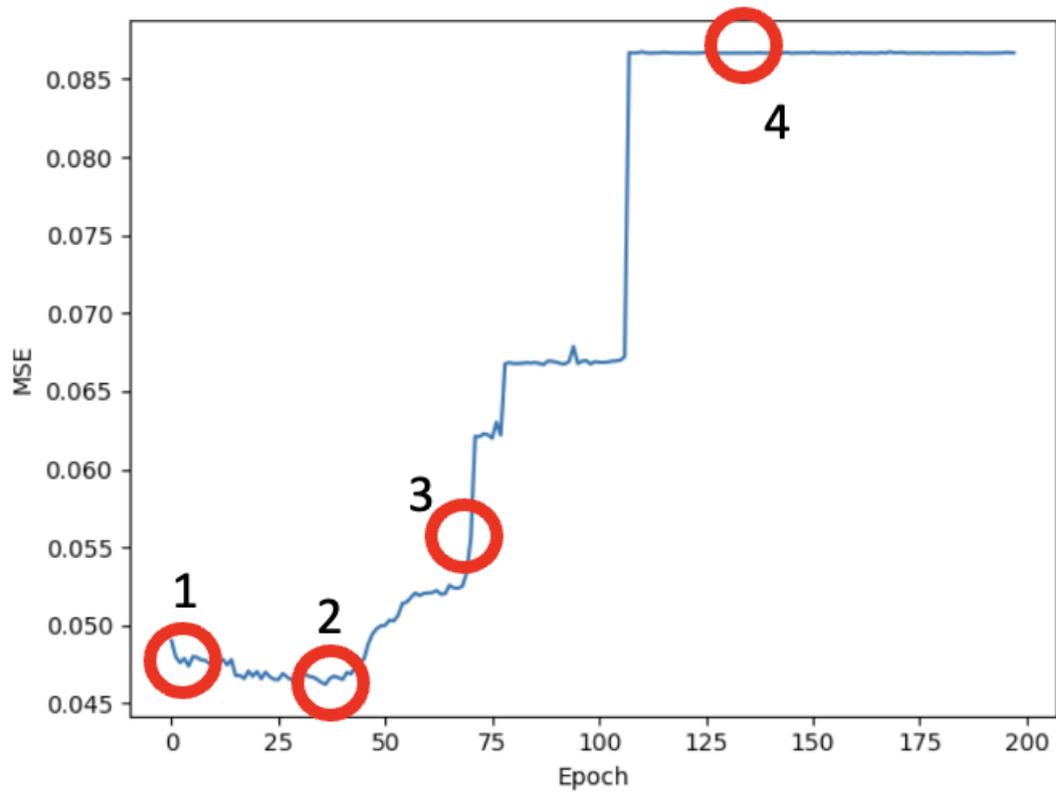
Block 7

Question 1:

Each circled region represents the robot trained with a different error score. Remember this error is on the 10-way clothing classification task.

Choose the point that you think the robot is going to perform the best at classifying the high-level categories.

- Shirts
- Shoes
- Trousers/Bags



Choose the robot now:

- 1
- 2
- 3
- 4

Block 8

Question 2:

Each visualization represents the robot trained with different high-level categories.

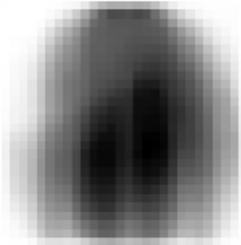
Choose the visualization that you think will help the robot perform best at sorting your clothes into these three high-level categories:

- Shirts
- Shoes
- Trousers/Bags

The first option is



The second option is



The Third Option is



The Fourth option is



Click to write the question text

- First option
- Second option
- Third option
- Fourth option

Block 9

Debrief Page

Thank you for your participation, this study was designed to evaluate people's ability to use an explainable artificial intelligence system to help fine-tuning towards down-stream tasks.

if you have any questions please contact 

Click next to finish

