# Project 1: Concept Bottleneck Models

*Abstract*—This paper is about exploring the concept bottleneck models in which concepts are used to predict the final output. Unlike the end-to-end models, correcting the concept prediction and propagating this modification to the target prediction are possible at the test time which leads to improving the model's accuracy. Furthermore, the different ways of learning the concept bottleneck models are proposed and the task error of these methods is replicated on the Caltech-UCSD Birds-200-2011 (CUB) dataset for comparison with the standard model. Two of these models are selected to train again from scratch to have concept bottleneck models. The effect of the test-time intervention on the performance of the concept bottleneck is studied.



Fig. 1: The concept bottleneck model for bird classification

## I. INTRODUCTION

Deep learning models are considered end-to-end models that are trained to predict the output directly from the input and there is no interpretation or intervention. For example, assume that a radiologist is assessing the severity of knee osteoarthritis with the help of a machine-learning model. The logic behind the model's prediction is important for the radiologist. If she adds extra information to the model based on her expertise, the model's prediction would be changed? Is it possible in an end-to-end machine learning model? To address this issue, the concept bottleneck models are proposed. Any end-to-end neural network can converted into this model by providing annotations for concepts in the training phase. It is shown that these bottleneck models achieve higher accuracies compared to standard models as we are allowed to modify the concept prediction to lead the model to predict correctly.

Three different learning methods are suggested to explore the interpretability and performance of concept bottleneck models for specific tasks. Assume $y$ is a target that should be predicted from input $x$ and $c$ is the concept, thus $f(.)$ predicts $c \rightarrow y$ and $g(.)$ predicts $x \rightarrow c$. $\lambda$ is the task-concept tradeoff hyperparameter. The standard model learns the target directly from the input without considering concepts $\hat{y} = f(g(x))$. The suggested learning methods are as follows:

- The Independent bottleneck: learn $\hat{y} = f(c)$ and $\hat{c} = g(x)$ separately.
- The Sequential bottleneck: first learn $\hat{c} = g(x)$ then use the prediction concept to learn $\hat{y} = f(g(x))$.
- The Joint bottleneck: learn $\hat{y} = f(g(x))$ and $\hat{c} = g(x)$ together in a way that minimizes the weighted sum $f$ and $g$ as $\hat{y}, \hat{c} = f(g(x)) + \lambda g(x)$ for some $\lambda > 0$.

Two different datasets are considered for assessing task accuracy and the high concept accuracy of the bottleneck model, an x-ray grading (OAI) [1] and a bird identification (CUB) [2] dataset. The OAI requires an application for data access, thus the rest of the paper is about replicating the result from the CUB dataset, which is public. The Caltech-UCSD Birds-200-2011 (CUB) dataset contains n = 11, 788 bird photographs that are classified into 200 bird species with k = 112 binary bird attributes representing wing color, beak shape, etc (Figure 1). The neural network architecture
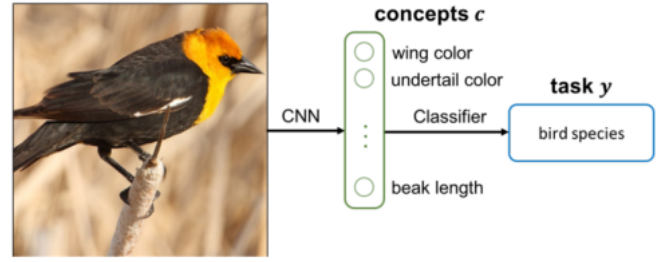
used for bird classification is Inception V3 [3]. This model without considering the fully-connected layers is pretrained on ImageNet and then the end-to-end model is fine-tuned [4] on the CUB dataset.

As the task of the CUB dataset is classification, $f$ and $g$ are turned into a probabilistic prediction to compute the real-valued scores. These models are trained to have a bottleneck layer and the predicted concepts can be read out from the bottleneck layer. This helps to manipulate the predicted concepts and explore the effect of this editing on the final prediction. The drawback of concept bottleneck models is requiring annotated concepts at training time. Also, some CUB concept annotations are noisy since they are provided by a single crowdworker (not a birding expert).

In contrast to the Independent bottleneck model, sequential bottleneck enables the $c \rightarrow y$ model to adjust based on its ability to forecast $x \rightarrow c$. In the joint bottleneck model, the concept model undergoes refinement in order to enhance predictive accuracy.

The rest of this paper is organized as follows. To validate the result obtained from the concept bottleneck model on the CUB dataset, the model is trained at least under one of the mentioned learning methods, and then the task error of Independent, Sequential, and Joint concept bottleneck models are compared with each other and with the standard model in section II. In section III, some examples of test-time intervention are provided to show the effect of editing the concept predictions in the target predictions. Conclusions are discussed in Section IV.

## II. VALIDATION

As aforementioned, there are three different kinds of methods to train the end-to-end model to have a concept bottleneck model. For this purpose, the CUB dataset and the processed CUB dataset are used to train the pretrained Inception V3 model for the bird classification task. The training process which is reported in [5] is done with a batch size of 64, and SGD with momentum of 0.9 as the optimizer. To assess the validity of this training process, it is tried to replicate this process to train the Joint model. This replication is done by

changing the batch size to 32 due to a GPU (graphics processing unit) memory issue. This model is successfully trained and automatically stopped in 690 epochs with a training accuracy of 92.5468 and training loss of 0.2408 as shown in Figure 2.

The training process of the Independent model, the part that predicts concepts from the input is also replicated. The training process lasts more than 24 hours and it is stopped in epoch 270 as reaches the training accuracy of 98.3087 and the training loss of 0.3173 as shown in Figure 3.

In addition, the performance of the obtained models (Independent, Sequential, and Joint models) is studied and compared with the standard model. In a way that each of these models is applied to test data to compute the error of concept prediction and target prediction. These models are obtained by training the end-to-end model under different training methods for bird classification. The predicted class labels obtained from each of these models are compared with ground truths to compute the error of concept prediction and target prediction as shown in Table I.

TABLE I: Task errors with $\pm 2SD$ over concept bottleneck and compared with the standard end-to-end model

| Model | Error of y | Error of c | Accuracy |
|---|---|---|---|
| Independent | $0.2400 \pm 0.0063$ | $0.0328 \pm 0.0009$ | 0.97 |
| Sequential | $0.2425 \pm 0.0028$ | $0.0328 \pm 0.0009$ | 0.97 |
| Joint(0.01) | $0.1987 \pm 0.0021$ | $0.0297 \pm 0.0002$ | 0.97 |
| Standard | $0.1746 \pm 0.0042$ | $0.5001 \pm 0.0029$ | 0.50 |

As shown in Table I, the error of the Independent and Sequential concept bottleneck models in predicting the output is worse than the standard model while the error of the Joint concept bottleneck model is close to the Standard model. However, the accuracy of concept bottleneck models is better than the standard model. The prediction of concepts in concept bottleneck models is much better than the Standard model. The accuracy of the Independent and Sequential concept bottleneck models is better than the Joint concept bottleneck model and all of the concept bottleneck models have better accuracy



Fig. 2: The details of Joint model training process



Fig. 3: The detail training process of Independent $x \rightarrow c$ model

compared to the Standard model. The values of this table are not exactly the same as those in Table 1 in the paper [5] but they follow the same pattern.

The aim of this paper is to show the effect of the intervention on the concept bottleneck model at test time on the final prediction. For this purpose, the paper shows the case study in two images of the CUB dataset which is shown in Figure 3 of the paper [5]. Replicating this part was challenging as the intervention process on concepts at test time is not elaborated in the paper or the GitHub repository. Based on Figure 3 of the paper [5], two images of the dataset are selected, it is tried to select those two same images (Glaucous_Winged_Gull_0112_44731.jpg and Worm_Eating_Warbler_0014_176042.jpg). To find those two images in the dataset, the test pickle file is read out and then tried to find images that have the same class labels (Glaucous_winged_Gull and Worm_eating_Warbler with class labels of 59 and 180) as those images in the paper.

As the paper has not specified which model was used to obtain those results, the performance of all models is surveyed on the new dataset. I figured out that seed1 of the Independent model predicts the class labels of 61 (Herring_Gull) and 98 (Ovenbird) as reported in the paper (Figure 4). Then, the test-time intervention with $n = 28$ is applied to those images, and the class prediction is changed to Glaucous_winged_Gull (59) and Ovenbird (98) (Figure 5). As shown, after concept intervention the model cannot predict the class label of the second image correctly. I checked all other models and different values for $n$ to reach the same values reported in the paper but I did not succeed.
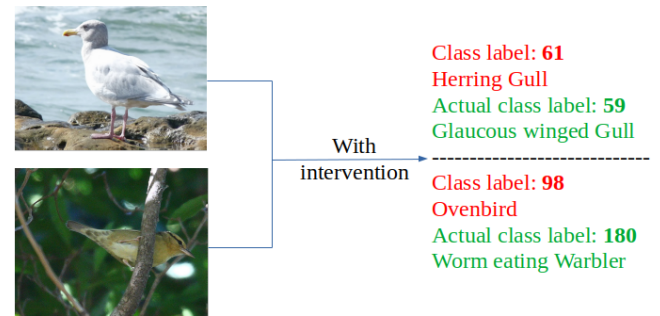


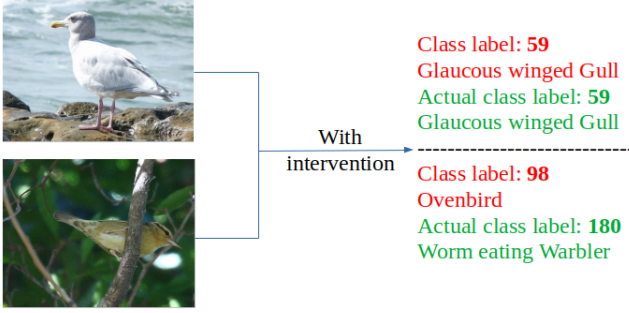Fig. 4: Class prediction of two examples of images without intervention

Fig. 5: Class prediction of two examples of images with intervention

TABLE II: Test-time intervention on different concept bottleneck models with $n = 1$

|  |  | Class prediction |  |  |
|---|---|---|---|---|
| Number of Concepts Intervened |  |  | $n = 1$ |  |
| Bottleneck models |  | Seed1 | Seed2 | Seed3 |
| Joint |  | [61, 98] | [61, 177] | [61, 98] |
| Independent |  | [61, 98] | [59, 98] | [65, 98] |
| Sequential |  | [61, 98] | [61, 98] | [61, 178] |

TABLE III: Test-time intervention on different concept bottleneck models with $n = 28$

|  |  | Class prediction |  |  |
|---|---|---|---|---|
| Number of Concepts Intervened |  |  | $n = 28$ |  |
| Bottleneck models |  | Seed1 | Seed2 | Seed3 |
| Joint |  | [61, 98] | [61, 180] | [61, 98] |
| Independent |  | [63, 98] | [59, 98] | [65, 98] |
| Sequential |  | [61, 98] | [59, 98] | [59, 36] |

TABLE IV: Test-time intervention on different concept bottleneck models

|  | Class prediction |  |
|---|---|---|
| Bottleneck models | $concepts = 0$ | $concepts = 1$ |
| Joint | [61, 180] | [61, 180] |
| Independent | [65, 98] | [65, 98] |
| Sequential | [59, 36] | [59, 36] |

## III. EXTENSION

In this section, it is tried to extend some methods proposed in the paper and study the effect of hyperparameters in the performance of concept bottleneck models. To better see the effect of editing the concepts at the test time on the class prediction, different bottleneck models and their seeds are considered as shown in Tables II and III. As you can see in some cases by changing at least one value of concepts the class prediction can be changed. Tables II and III show the class prediction of different models under the intervention of $n = 1$ and $n = 28$ ($n$ is the number that is intervened). The dataset is as same as the previous part and the ground truth class labels are [59, 180]. It seems that the number of concepts that intervened has an effect on the model performance. Based on the graphs in Figure 4 of the paper, increasing the number of concepts intervened leads to decreasing the task error.

Based on the paper, as the CUB dataset is a classification, not a regression, the intervention on the CUB dataset is complicated. Assume $\hat{l} = \hat{g}(x)$, to modify a concept $\hat{c}j$, the true value $cj$ cannot be replaced simply. Instead, we must adjust $\hat{l}j$ so that the probability of $\hat{c}j$ being 1 is close to the true value of $cj$. We intervene on $\hat{c}j$ by setting $\hat{l}j$ to either the 5th percentile (if $cj = 0$) or the 95th percentile (if $cj = 1$) of $\hat{l}j$ over the training distribution.

In addition, to better see the effect of concepts on class prediction, the whole concept values are changed to 0 and 1 and the results are shown in Table IV. As it is obvious, the second class label is finally predicted correctly when the whole values of concepts intervened are changed to 1 or 0 for the Joint model.

## IV. CONCLUSION

We are trying to duplicate the results that are reported in the main paper. At first, we train two models, Joint (from input to output) and Independent (from input to concept). Then, the results of task errors overall seeds of bottleneck models are computed and compared with the Standard model. The effect of test time intervention on two specific images is also studied but replicating the result is not successful. The specific model used to achieve those results has not been mentioned in the paper, after investigating all the models, the model that has similar behavior to those results is selected. The test time intervention is also applied to this model but it leads to predicting one of the class labels correctly instead of both.

In the extension section, some other test time intervention is studied to show the role of concepts in predicting class labels in bottleneck models. The effect of the values of concept and the number of concepts intervened is studied on the performance of models. Concept bottleneck models can compete on task accuracy while supporting intervention and interpretation, however, a drawback of these models is that they require annotated concepts at training time.

## REFERENCES

[1] Nevitt, M., Felson, D. T., and Lester, G. The Osteoarthritis Initiative. Cohort study protocol, 2006.
[2] Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The Caltech-UCSD Birds-200-2011 dataset. Technical report, California Institute of Technology, 2011.
[3] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the Inception architecture for com- puter vision. In Computer Vision and Pattern Recognition (CVPR), pp. 2818–2826, 2016.
[4] Cui, Y., Song, Y., Sun, C., Howard, A., and Belongie, S. Large-scale fine-grained categorization and domain-specific transfer learning. In Computer Vision and Pattern Recognition (CVPR), pp. 4109–4118, 2018.
[5] Koh, P. W., Nguyen, T., Tang, Y. S., Mussmann, S., Pierson, E., Kim, B. and Liang, P. Concept bottleneck models. In International conference on machine learning, 2020.