# Hybrid Concept-based Models: Using Concepts to Improve Neural Networks' Accuracy

Anonymous Full Paper Submission 3

# **OD1** Abstract

Most datasets used for supervised machine learning 002 consist of a single label per data point. However, in 003 cases where more information than just the class la-004 bel is available, would it be possible to train models 005 006 more efficiently? We introduce two novel model architectures, which we call hybrid concept-based mod-007 els, that train using both class labels and additional 008 information in the dataset referred to as *concepts*. In 009 order to thoroughly assess their performance, we in-010 troduce ConceptShapes, an open and flexible class of 011 datasets with concept labels. We show that the hy-012 brid concept-based models can outperform standard 013 computer vision models and previously proposed 014 concept-based models with respect to accuracy. We 015 also introduce an algorithm for performing adversar-016 *ial concept attacks*, where an image is perturbed in 017 a way that does not change a concept-based model's 018 concept predictions, but changes the class prediction. 019 The existence of such adversarial examples raises 020 questions about the interpretable qualities promised 021 by concept-based models. 022

# 023 1 Introduction

Understanding model behavior is a crucial chal-024 lenge in deep learning and artificial intelligence [1-4]. 025 Deep learning models are inherently chaotic, and 026 give little to no insight into why a prediction was 027 made. In computer vision, early attempts of ex-028 plaining a model's prediction consisted of assigning 029 pixel-wise feature importance, referred to as *saliency* 030 maps [5-8]. Despite gaining popularity and being 031 visually appealing, a large number of experiments 032 show that saliency maps perform a poor job at ac-033 tually explaining model behavior [2, 3, 9-13]. 034

Recently, several *concept-based models* have been 035 proposed as inherently interpretable [14–18]. These 036 models are restricted to perform the downstream 037 prediction only based on whether it thinks some 038 predefined *concepts* are present in the input or not, 039 where concepts are defined as human meaningful 040 features. This way, the downstream predictions can 041 be interpreted by which concepts the model thought 042 were in the data. 043

However, recent experiments have highlighted issues with concept-based models' interpretability.
This is mainly due to the concept predictions encod-



Figure 1. Adversarial Concept Attack. Images are perturbed in a way that does not change a conceptbased model's concept predictions, but change the class prediction. This brings into question the interpretable qualities of these models.

ing more information than just the concepts, referred 047 to as *concept leakage* [19, 20]. We further add evidence to the lack of interpretability in concept-based 049 models by introducing *adversarial concept attacks* 050 (see Figure 1). 051

Due to the evidence demonstrating the limitations of of interpretability in concept-based models, we will os shift our focus away from interpretability and instead use the framework of concept-based models to improve the performance of the models. We present os two new model architectures aimed to achieve this. os

Our proposed model architectures use both con-058 cept predictions and information not interfering with 059 the concepts to make the downstream prediction. 060 This way, the models can use the concept predictions 061 if they are helpful for the downstream task, but can 062 also rely on a skip connection to encode information 063 about the data not present in the concepts. We 064 propose these models to better utilize the available 065 information in datasets with concepts. 066

A challenge in this research area is that the most 067 popular datasets used for benchmarking concept-068 based models have shortcomings that we argue 069 make them unsuitable to use as benchmarks, and 070 we therefore developed a new class of flexible con-071 cept datasets. The Caltech-USCD Birds-200-2011 072 dataset (CUB) [21] is the most widely used concept 073 dataset, where the downstream task is to classify im-074

ages among 200 classes of bird species, and is widely 075 used to benchmark concept-based models [14–18, 076 22–24]. Despite its popularity, there are various 077 problems with the concept labeling, which was done 078 by non-experts. Therefore, the dataset is processed 079 with a class-wise majority vote of the concept labels, 080 so every class has the exact same concept labels [15]. 081 Unfortunately, not only has this led to mistakes 082 where instances of a class have different concepts 083 [18], but ambiguity of concepts is very common. For 084 instance, there are images of birds where one can 085 not see its tail, belly or wings properly, but it may 086 still be labeled with concepts relating to those body 087 parts [18]. 088

Another popular concept dataset is Osteoarthritis 089 Initiative (OAI) [15, 17, 24, 25]. The main problem 090 with this dataset is the lack of availability. Since it 091 uses medical data, access to the dataset needs to be 092 requested, and the processed version is not directly 093 available. Moreover, the computational resources 094 used to process it was "several terabytes of RAM 095 and hundreds of cores" [26], which is not available 096 for many researchers. 097

Our contributions are as follows:

- Novel model architectures: We propose 099 novel model architectures which we call hy-100 brid concept-based models. Unlike previously 101 proposed concept-based models, ours are moti-102 vated by performance, not interpretability. We 103 conduct experiments that show that they can 104 outperform other computer vision and concept-105 based models. 106
- New concept datasets: In order to properly assess the performance of concept-based models, we propose a new set of openly available datasets with concepts called *ConceptShapes*.
- Adversarial concept attacks: We propose an algorithm for generating adversarial examples specifically for concept-based models. These examples further question concept-based models' interpretable qualities.

# 116 2 Related Work

#### 117 2.1 Concept-based Models

Concept-based models first predict some predefined 118 concepts in the dataset, then use those concept pre-119 dictions to predict the downstream task. This way, 120 the final prediction can be interpreted by which con-121 cepts the model thought were present in the input. 122 One of the first and most popular concept-based 123 models is the *concept bottleneck model* (CBM) [15], 124 which is a neural network with a *bottleneck layer* 125 that predicts the concepts. The model is trained 126 127 both using the concept labels and the target labels.



Figure 2. Architecture of a CBM-Res used for computer vision. The model uses both concept predictions and a skip connection that hops over the bottleneck layer to perform the output prediction. The architecture can be adapted to a CBM-Skip by performing concatenation instead of addition before the output layer.

Several alternatives to the CBM architecture have 128 been proposed. Concept-based model extraction 129 (CME) [14] may use a different hidden layer for 130 the various concepts. *Post-hoc concept bottleneck* 131 models (PCBM) [23] first learn the concept activa-132 tion vectors (CAVs) [27] of the concepts, and then 133 project embeddings down on a space constructed by 134 CAVs. Concept embedding models (CEM) [22] pro-135 duces vectors in a latent space of concepts that are 136 different for presence and absence of a concept and 137 predicts the probabilities of concepts being present. 138

Several experiments show that the concept pre-139 dictions encode more information than just the con-140 cepts, and therefore that they are unsuitable to use 141 as interpretation of the models' behavior [19, 20]. It 142 has also been shown that concept-based models are 143 susceptible to adversarial attacks that change the 144 concept predictions [24], but not the class predic-145 tions. 146

## 3 Methods

#### 147

148

#### 3.1 New Model Architectures

We propose two novel model architectures. The first 149 is based on a CBM [15], but uses an additional skip 150 connection that does not go through the concept 151 bottleneck layer (see Figure 2). The skip connection 152 can be implemented either as a residual connection 153 [28] or a concatenation [29], and we refer to the mod-154 els as CBM-Res and CBM-Skip, respectively. This 155 way, the model can use both the concept prediction 156 and information not interfering with the concepts 157 to make the final downstream prediction. 158

The other proposed architecture predicts the concepts sequentially throughout the neural network's 160 layers, instead of all at once (see Figure 3). All of 161 the concept predictions are concatenated together, 162 along with the final hidden layer, and given as input to the output layer. We refer to this model as the 164 Sequential Concept Model (SCM). 165

We use a loss function constructed by a weighted 166 sum of a concept loss and a task loss, similarly to 167



Figure 3. Architecture of a Sequential Bottleneck Model (SCM). The concepts are predicted sequentially throughout the layers, and concatenated together with the final hidden layer before the output layer, which produces the downstream predictions.

the *joint bottleneck* proposed for the CBM [15].

All of the proposed models are compatible with transfer learning, where the first part of the model can be a large pre-trained network. We present the models in the domain of classification, but they can easily be adapted to regression by replacing the

174 output layer with a single node.

#### 175 3.2 Introducing ConceptShapes

To accurately assess the performance of concept-176 based models, we have developed a class of flexible 177 synthetic concept datasets called *ConceptShapes*. 178 The input images consist of two shapes, where 179 the position and orientation are random, and the 180 downstream task is to classify which combination 181 of shapes that are present (see Figure 4). Some 182 examples of target classes are "triangle-rectangle", 183 "triangle-triangle" and "hexagon-pentagon". Depend-184 ing on how many shapes that are used, the dataset 185 186 contains 10, 15 or 21 classes.

The key feature of the datasets are that various 187 binary concepts are present, such as the color of 188 the shapes, outlines and background. Given a class, 189 some predefined concepts are drawn with a high 190 probability  $s \in [0.5, 1]$ , and the others with a low 191 probability 1 - s. The hyperparameter s can be 192 chosen by the user. When s = 0.5, the concepts 193 are drawn independently of the classes, and when 194 s = 1, the concepts are deterministic given the class. 195 The datasets can be created with either five or nine 196 concepts. 197

The datasets are flexible with regards to the rela-198 tionship between the concepts and the classes, and 199 the number of concepts, classes and data. This 200 way, the difficulty of correctly classifying the im-201 ages can be tuned by the amount of classes and the 202 amount of data, and the information in the concepts 203 can be tuned by the amount of concepts and the 204 value of s. Further details about the dataset can be 205 found in Appendix B. The code for generating the 206 ConceptShapes datasets can be found at https:// 207 208 anonymous.4open.science/r/ConceptShapes/.



Figure 4. Images from different classes of two ConceptShapes datasets. Left: Nine different images from a 10-class 5-concept dataset. Right: Nine different images from a 21-class 9-concept dataset.

Although synthetic datasets may have less com-209 plex patterns than datasets with real images, there 210 are also clear benefits of using them. They are pre-211 cisely labeled, there is no ambiguity in the concepts 212 and they have many flexible parameters. Therefore, 213 we believe that ConceptShapes can provide a use-214 ful addition to the existing benchmark datasets for 215 concept-based models. 216

#### 3.3 Adversarial Concept Attacks 217

We propose an algorithm for producing *adversarial* 218 *concept attacks*, which given a concept-based model 219 and input images, produces identically looking im-220 ages that give the same concept predictions, but 221 different output predictions. The algorithm is based 222 on projected gradient descent (PGD) [30], which it-223 eratively updates an input in the direction which 224 maximizes the classification error of the model, and 225 projects the alteration on an L-infinity ball around 226 the original input. In each iteration, we add a step 227 where we check if the alteration causes the model to 228 almost change the predictions of the concepts. If so, 229 we check which pixels that are altered in a direction 230 that changes the concept predictions, and multiply 231 those pixels' alterations by a number in [-1, 0]. The 232 complete algorithm is covered in Appendix C. 233

Our approach differs from the adversarial attacks 234 for concept-based models done by Sinha et al. [24], 235 which altered the concept predictions, and not the 236 class predictions. 237

# 4 Experiments

238

We show that the hybrid concept-based models 239 achieve the highest test set accuracy on multiple 240 datasets. In order to examine how the models perform with different amounts of data, we train and 242 test the models on various smaller subsets of the 243 datasets. We also investigate how well the concepts 244 are learned. 245

#### 247 4.1.1 Caltech-USCD Birds-200-2011 (CUB)

The CUB dataset [21] consists of N = 11,788 images 248 of birds, where the target is labeled among 200 bird 249 species. The original dataset contains 28 categorical 250 concepts, which makes 312 binary concepts when 251 one-hot-encoded. The processed version used for 252 benchmarking concept-based models [15] removed 253 sparse concepts and used a majority vote on the 254 concepts labels, so that every class has the exact 255 same concepts, ending up with 112 binary concepts. 256 The dataset is split in a 50%-50% training and test 257 split, and we use 20% of the training images for 258 259 validation. We train and evaluate on six different subset sizes. 260

#### 261 4.1.2 ConceptShapes

We experiment with many different ConceptShapes 262 configurations. First, we set the probability s to 263 be 0.98, in order to make the concepts useful, but 264 not deterministic given the class, and experiment 265 with different amounts of classes. Additionally, we 266 use different values of s to explore how the correla-267 tion between the concepts and the classes influence 268 the models' performance. We explore even more 269 configurations of ConceptShapes in Appendix A.2. 270 The datasets are generated with 1000 images in 271

each class, and we split them into 50%-30%-20%
train-validation-test sets. We train and evaluate on
subsets sizes with 50, 100, 150, 200 and 250 images
in each class, drawn from the 1000 images created.

#### 276 4.2 Setup

We compare our proposed models against a CBM, 277 which we refer to as *vanilla CBM*, and a convo-278 lutional neural network (CNN) [31] not using the 279 concepts at all, referred to as the standard model. 280 Additionally, we also include an *oracle model*, which 281 is a logistic regression model trained only on the 282 true concept labels, not using the input images. We 283 call it an *oracle* since it uses true concept labels at 284 test time, which are usually unknown, and do this 285 in order to measure how much information there is 286 in the concepts alone. 287

The models' accuracies are evaluated on a held-out test set, which is the same for every subset configuration. We perform a hyperparameter search for each model and each subset configuration using the validation sets. The details of the hyperparameter settings are covered in Appendix D.

The models trained on CUB use a pre-trained and frozen ResNet 18 [28] as the convolutional part of the model, while the models trained on ConceptShapes are trained from scratch. The details about the setup are explained in Appendix E. The code for running the experiments can be



Figure 5. Test set accuracies on the CUB dataset. The x-axis indicates the average amount of images included in the training and validation dataset for each class, where the rightmost point corresponds to the full dataset. The results are averaged over three runs and include 95% confidence intervals. The oracle model consistently got 100% accuracy and is omitted.



Figure 6. MPO scores on the CUB dataset. The y-axis indicates the proportion of images with m or more concept prediction mistakes. The results are averaged over three runs and include 95% confidence intervals. We used the full dataset.

found at https://anonymous.4open.science/r/ 300 Hybrid-Concept-based-Models-22C3. 301

We also record the *Misprediction overlap* (MPO) 302 [14] to measure the quality of the concept predictions. 303 The MPO measures the ratio of images that had 304 m or more concept mispredictions. We use m = 305  $1, 2, \ldots, k$ , where k is the amount of concepts in the 306 dataset. 307

We run a grid search to find the most successful adversarial concept attack ratio, meaning the ratio of images that changes the model's class prediction, but not the concept predictions. We use the best vanilla CBM found in the hyperparameter searches. and compare the results to PGD [30].



Figure 7. Test set accuracy on ConceptShapes with nine concepts and s = 0.98. The plots show that the hybrid concept-based models perform better than the benchmark models. The x-axis denotes how many training and validation images that were included in each class. The metrics are averaged over ten runs and include 95% confidence intervals.



Figure 8. MPO scores on ConceptShapes. All three plots shows that the concepts are properly learned by all models, where about 85% of the images has no concept mispredictions, and only about 2% has more than two. We used 250 training and validation images in each class. The metrics are averaged over ten runs and include 95% confidence intervals.

## 314 5 Results and Discussion

#### 315 5.1 Improved Accuracy

NLDL

#3

The test set accuracies can be found in Figure 5, 7 316 and 9. We observe that the hybrid concept-based 317 models generally have the best performance. When 318 s = 0.5 (Figure 9), the concepts provide no addi-319 tional information that helps predict the classes. 320 However, the hybrid concept-based models do not 321 perform worse than the CNN, suggesting they are 322 able to assign low weights to the bottleneck layer 323 when the concepts are irrelevant. When s increases, 324 the performance of all models does as well, and the 325 gap between the hybrid concept-based models and 326 the benchmark models becomes larger. 327

The oracle model, using only the true concept la-328 bels at test time, serves as a baseline for how much 329 information there is in the concepts alone. When 330 s = 0.5 (Figure 9), the oracle model has a test set 331 accuracy of about 10%, similarly to random guess-332 ing. The oracle model's accuracy increases when s333 increases, but decreases when there are more classes. 334 Since s controls the information in the concepts, and 335 an increase in the number of classes also increases 336 the difficulty of classifying, this behavior is expected. 337

With a low s or a large number of classes, the hybrid concept-based models perform better than the oracle model (see Figure 7, 9). 340 NLDL

#3

#### 5.2 CUB Concepts are not Learned 341

When inspecting the MPO plots for CUB in Fig-342 ure 6, we see that none of the concept-based models 343 learn the concepts properly. About 50% of the images are predicted with 15 or more mistakes in the 345 concept predictions for all models, which is about 346 as good as random guessing, since the labels are 347 one-hot-encoded and sparse. Because earlier work 348 has pointed out that the concepts in CUB are am-349 biguous and sometimes wrong [18], this might not 350 be surprising. 351

When inspecting the MPO plot for ConceptShapes 352 in Figure 8 and 10, we do however see that the concepts are properly learned, even when s = 0.5 and 354 the concepts are useless for predicting the class. This suggests that the concepts in ConceptShapes are not 356 ambiguous, and the datasets therefore serve as better benchmark datasets for concept-based models. 358



Figure 9. Test set accuracy on ConceptShapes with different values of s, using ten classes and nine concepts. Higher values of s means more correlation between the concepts and the classes.



Figure 10. MPO scores on ConceptShapes with different values of s using ten classes and nine concepts. The majority of images are have less than two concept mispredictions, even when the concepts are irrelevant for the classes (s = 0.5). The SCM performs better than the other models.

	Adversarial Concept	PGD
	Attack Success Rate	Success Rate
CUB		
112 concepts	57.4%	16.2%
ConceptShapes with		
10 classes and 5 concepts	35.5 %	31.4%
ConceptShapes with		
21 classes and 9 concepts	26.6%	22.5%

**Table 1.** Success rate of adversarial concept attacks on images in the test sets. An attack is considered a success when the class prediction is changed, but not the concept predictions.

#### **5.3** Adversarial Concept Attacks

The results of the adversarial concept attacks can be found in Table 1. We see that a substantial amount of images are perturbed with success, and the algorithm is more effective than PGD.

Since the concept predictions are used as the interpretation of the model behavior, but the same interpretation can lead to vastly different model behavior, we suggest that this experiment questions the interpretable qualities of concept-based models.

# **369 6 Conclusions**

We proposed new hybrid concept-based models motivated by improving performance, and demonstrated their effectiveness on CUB and several Concept-Shapes datasets. The proposed models train using 373 both the class label and additional concept labels. 374 In all of the datasets we experimented with, the 375 hybrid concept-based models performed better than 376 previously proposed concept-based models and the 377 standard computer vision models. 378

We also introduced ConceptShapes, a flexible 379 class of synthetic datasets for benchmarking conceptbased models. Finally, we demonstrated that 381 concept-based models are susceptible to adversarial 382 concept attacks, which we suggest are problematic 383 for their promised interpretable qualities. 384

In future work, we would like to apply hybrid 385 concept-based models in the domain of reinforcement learning, where concepts such as agent rotation, position and velocity can be automatically calculated, and do not need manual labeling. 389

## References

 F. Doshi-Velez and B. Kim. "Towards a rigorous science of interpretable machine learning". 392 In: arXiv preprint arXiv:1702.08608 (2017). 393

390

[2] Z. C. Lipton. "The mythos of model interpretability: In machine learning, the concept 395 396of interpretability is both important and slip-397pery." In: Queue 16.3 (2018), pp. 31–57.

- [3] C. Rudin. "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead". In: *Nature machine intelligence* 1.5 (2019), pp. 206–215.
- [4] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser,
  A. Bennetot, S. Tabik, A. Barbado, S. García,
  S. Gil-López, D. Molina, R. Benjamins, et
  al. "Explainable Artificial Intelligence (XAI):
  Concepts, taxonomies, opportunities and challenges toward responsible AI". In: *Information fusion* 58 (2020), pp. 82–115.
- K. Simonyan, A. Vedaldi, and A. Zisserman.
  "Deep inside convolutional networks: Visualising image classification models and saliency
  maps". In: arXiv preprint arXiv:1312.6034
  (2013).
- [6] R. R. Selvaraju, M. Cogswell, A. Das, R.
  Vedantam, D. Parikh, and D. Batra. "Gradcam: Visual explanations from deep networks
  via gradient-based localization". In: *Proceed- ings of the IEEE international conference on computer vision.* 2017, pp. 618–626.
- A. Shrikumar, P. Greenside, and A. Kundaje.
  "Learning important features through propagating activation differences". In: *International conference on machine learning*. PMLR. 2017, pp. 3145–3153.
- M. Sundararajan, A. Taly, and Q. Yan. "Axiomatic attribution for deep networks". In: *International conference on machine learning*.
  PMLR. 2017, pp. 3319–3328.
- [9] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim. "Sanity checks
  for saliency maps". In: Advances in neural information processing systems 31 (2018).
- [10] A.-K. Dombrowski, M. Alber, C. Anders, M.
  Ackermann, K.-R. Müller, and P. Kessel. "Explanations can be manipulated and geometry
  is to blame". In: Advances in neural information processing systems 32 (2019).
- 438 [11] A. Ghorbani, A. Abid, and J. Zou. "Interpretation of neural networks is fragile". In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 33. 01. 2019, pp. 3681–3688.
- [12] P.-J. Kindermans, S. Hooker, J. Adebayo, M.
  Alber, K. T. Schütt, S. Dähne, D. Erhan,
  and B. Kim. "The (un) reliability of saliency
  methods". In: *Explainable AI: Interpreting, explaining and visualizing deep learning* (2019),
  pp. 267–280.

- M. Ghassemi, L. Oakden-Rayner, and A. L. 448 Beam. "The false hope of current approaches 449 to explainable artificial intelligence in health 450 care". In: *The Lancet Digital Health* 3.11 451 (2021), e745–e750. 452
- [14] D. Kazhdan, B. Dimanov, M. Jamnik, P. Liò, 453 and A. Weller. "Now you see me (CME): 454 concept-based model extraction". In: arXiv 455 preprint arXiv:2010.13233 (2020). 456
- P. W. Koh, T. Nguyen, Y. S. Tang, S. Muss-457 mann, E. Pierson, B. Kim, and P. Liang. "Con-458 cept bottleneck models". In: *International con-*459 *ference on machine learning*. PMLR. 2020, 460 pp. 5338-5348.
- Y. Sawada and K. Nakamura. "Concept bottleneck model with additional unsupervised concepts". In: *IEEE Access* 10 (2022), pp. 41758– 464 41765.
- K. Chauhan, R. Tiwari, J. Freyberg, P. Shenoy, 466
  and K. Dvijotham. "Interactive concept bot-467
  tleneck models". In: Proceedings of the AAAI 468
  Conference on Artificial Intelligence. Vol. 37. 469
  5. 2023, pp. 5948–5955. 470
- E. Kim, D. Jung, S. Park, S. Kim, and S. Yoon. 471
   "Probabilistic Concept Bottleneck Models". In: 472
   arXiv preprint arXiv:2306.01574 (2023). 473
- [19] A. Mahinpei, J. Clark, I. Lage, F. Doshi-474
  Velez, and W. Pan. "Promises and pitfalls of 475
  black-box concept learning models". In: arXiv 476
  preprint arXiv:2106.13314 (2021). 477
- [20] A. Margeloiu, M. Ashman, U. Bhatt, Y. Chen, 478
  M. Jamnik, and A. Weller. "Do concept bottleneck models learn as intended?" In: arXiv 480
  preprint arXiv:2105.04289 (2021). 481
- [21] C. Wah, S. Branson, P. Welinder, P. Perona, 482
   and S. Belongie. *The Caltech-UCSD Birds-200-* 483
   *2011 dataset.* Tech. rep. California Institute of 484
   Technology, 2011. 485
- M. Espinosa Zarlenga, P. Barbiero, G. 486 Ciravegna, G. Marra, F. Giannini, M. Diligenti, Z. Shams, F. Precioso, S. Melacci, A. 488 Weller, et al. "Concept embedding models: 489 Beyond the accuracy-explainability trade-off". 490 In: Advances in Neural Information Processing 491 Systems 35 (2022), pp. 21400–21413. 492
- M. Yuksekgonul, M. Wang, and J. Zou. "Posthoc concept bottleneck models". In: arXiv 494 preprint arXiv:2205.15480 (2022).
- [24] S. Sinha, M. Huai, J. Sun, and A. Zhang. 496
  "Understanding and enhancing robustness of 497
  concept-based models". In: Proceedings of the 498
  AAAI Conference on Artificial Intelligence. 499
  Vol. 37. 12. 2023, pp. 15127–15135. 500

- [25] M. Nevitt, D. Felson, and G. Lester. "The
  osteoarthritis initiative". In: Protocol for the
  cohort study 1 (2006).
- E. Pierson, D. M. Cutler, J. Leskovec, S. Mullainathan, and Z. Obermeyer. "An algorithmic
  approach to reducing unexplained pain disparities in underserved populations". In: *Nature Medicine* 27.1 (2021), pp. 136–140.
- 509 [27] B. Kim, M. Wattenberg, J. Gilmer, C. Cai,
  510 J. Wexler, F. Viegas, et al. "Interpretability
  511 beyond feature attribution: Quantitative test512 ing with concept activation vectors (tcav)". In:
  513 International conference on machine learning.
  514 PMLR. 2018, pp. 2668–2677.
- [28] K. He, X. Zhang, S. Ren, and J. Sun. "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [29] G. Huang, Z. Liu, L. Van Der Maaten, and
  K. Q. Weinberger. "Densely connected convolutional networks". In: *Proceedings of the IEEE conference on computer vision and pat- tern recognition*. 2017, pp. 4700–4708.
- [30] A. Madry, A. Makelov, L. Schmidt, D. Tsipras,
  and A. Vladu. "Towards deep learning models
  resistant to adversarial attacks". In: arXiv
  preprint arXiv:1706.06083 (2017).
- [31] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D.
  Jackel. "Backpropagation applied to handwritten zip code recognition". In: *Neural computation* 1.4 (1989), pp. 541–551.
- T. Akiba, S. Sano, T. Yanase, T. Ohta, and M.
  Koyama. "Optuna: A next-generation hyperparameter optimization framework". In: Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining. 2019, pp. 2623–2631.
- 540 [33] D. P. Kingma and J. Ba. "Adam: A method for
  541 stochastic optimization". In: arXiv preprint
  542 arXiv:1412.6980 (2014).
- [34] A. Paszke, S. Gross, F. Massa, A. Lerer, J.
  Bradbury, G. Chanan, T. Killeen, Z. Lin, N.
  Gimelshein, L. Antiga, et al. "Pytorch: An imperative style, high-performance deep learning
  library". In: Advances in neural information
  processing systems 32 (2019).
- [35] N. Srivastava, G. Hinton, A. Krizhevsky, I.
  Sutskever, and R. Salakhutdinov. "Dropout: a simple way to prevent neural networks from overfitting". In: *The journal of machine learning research* 15.1 (2014), pp. 1929–1958.



Figure A.1. Test set accuracies on the CUB dataset with a hard bottleneck. The vanilla CBM suffers from substantially lower performance, while the hybrid concept-based models' performances are not changed much. The metrics are averaged over three runs and contains 95% confidence intervals.

[36] Russakovsky et al. "Imagenet large scale visual 554 recognition challenge". In: International jour-555 nal of computer vision 115 (2015), pp. 211-556 252.

# A More Experiments 558

## A.1 Hard Bottleneck 559

We also conducted experiments where we rounded 560 off the concept predictions to binary values in the 561 concept-based models, referred to as a hard bottle-562 *neck.* The results can be seen in Figure A.1 and 563 Figure A.2. We see that the vanilla CBM's perfor-564 mance is reduced, while the hybrid concept-based 565 models are not changed much. Comparing the MPO 566 plots (Figure A.2 and Figure 6) shows little change 567 in how well the concepts were learned. 568

#### A.2 ConceptShapes Experiments 569

We experiment with using five concepts instead of 570 nine, with the results plotted in Figure A.3. We see 571 that the oracle model has a lower accuracy, indicat- 572 ing there is less information in the concepts alone. 573 Thus, the gap between the standard model and the 574 hybrid concept-based models is a little narrower, 575 although the hybrid concept-based models still per-576 form the best in general. When inspecting the MPO 577 plot in Figure A.4, we see that the concepts are still 578 learned by all the models. 579

# B ConceptShapes Details 580

We now explain the ConceptShapes datasets in 581 greater detail. The crucial feature of the datasets 582



Figure A.2. MPO scores on the CUB dataset with a hard bottleneck. There is little to no change in how well the concepts are predicted when using a hard bottleneck. We used the full dataset.

are the concepts. All of them are binary and inde-583 pendent, meaning any combination of concepts are 584 possible. The five first concepts are based on the 585 two shapes in the image, while the last four optional 586 concepts are based on the background. We now 587 describe the concepts one-by-one, and visualizations 588 are available in Table B.1 and Table B.2. We start 589 with the five concepts that influence the shapes: 590

- 591 1. Big shapes. Every shape had two intervals of
  592 sizes to be randomly drawn from. One interval
  593 corresponded to the small figures, and the other
  594 to big ones.
- 2. Thick outlines. The outlines of the shapes
  were drawn from one of two intervals. One
  corresponded to a thin outline, and the other
  to a thick one.
- 599 3. Facecolor. There were two possible colors for600 the shapes, blue and yellow.
- 4. Outline color. The shapes had two possible
  outline colors, red and white.
- 5. Stripes. Some shapes were made with stripes,
  and some were not. The stripes were in the
  same color as the outline.

All of the concepts apply to the whole image. For
instance, if the image gets the thick outline concept,
both shapes in the image get a thick outline.

The datasets that use nine concepts have all five of the concepts above, in addition to four more. While all of the five-concept datasets have black backgrounds, the nine-concept datasets split the background in two and use the color and stripes as concepts.

615 6. Upper background color. The upper-half
616 of the background would either be magenta or
617 pale-green.



Table B.1. Overview of the five concepts regarding the shapes. Each row corresponds to one concept. The two leftmost columns of images have the concept, and the two rightmost columns do not have the concept. All of the images are from a 5-concept dataset, hence the black background. These five concepts are also present in the 9-concept datasets.

- 7. Lower background color. The lower-half 618 of the background would be either indigo or 619 dark-sea-green. 620
- 8. **Upper background stripes**. This represented whether there were black stripes present 622 in the upper background or not. 623
- 9. Lower background stripes. This represented 624 whether there were black stripes present in the 625 lower background or not. 626

To summarize, some of the image's visuals are 627 determined by the concepts, some by the classes and 628 some by randomness. The two shapes (from triangle, 629 square, pentagon, hexagon, circle and wedge) are 630 determined by the class. The shapes' size, color 631 and outline are determined by the concepts. If the 632 dataset uses nine concepts, the background color 633 and stripes are also determined by the concepts. 634 The shapes' position and rotation are determined 635 randomly, regardless of which class or concepts they 636 have. 637

# C Adversarial Concept At- 638 tacks 639

We describe a high level overview of the algorithm 640 for performing adversarial concept attacks, while 641 the full algorithm can be found in Table 1. We do 642 PGD [30] with an additional step. For each concept, 643



Figure A.3. Test set accuracy on ConceptShapes with five concepts and s = 0.98. The metrics are averaged over ten runs and include 95% confidence intervals.



Figure A.4. MPO scores on ConceptShapes with s = 0.98. We used 250 images in each class. The metrics are averaged over ten runs and include 95% confidence intervals.

if the concept prediction for the perturbed image is 644 close to changing compared to the original image, 645 we consider the concept as *sensitive*. We control 646 this with a sensitivity threshold  $\gamma \in [0, \infty)$ , so that 647 a concept is considered sensitive if its logits is in the 648 interval  $[-\gamma, \gamma]$ . We get an initial perturbation in 649 each iteration from following the gradient of the loss 650 with respect to the pixels of the images, similar to 651 PGD. The perturbation is multiplied with a mask 652 M, which is constructed so that it is 1 for pixels that 653 do not influence sensitive concepts to be even closer 654 655 to a change of prediction, and  $\beta \in [-1, 0]$  for pixels that do. We loop over the sensitive concepts and 656 iteratively update M, and use the notation  $\mathbb{I}_{\beta}(A =$ 657 B) to denote an elementwise indicator function, so 658 that it is 1 if elements in the same place in A and 659 B are equal, and  $\beta$  if not. The Hadamard product 660  $\odot$  represent elementwise multiplication. 661

The algorithm might terminate with failure due to several reasons. In any steps, if the concept predictions are changed, we terminate. If the mask M has only  $\beta$  as elements, the perturbation will not go in a direction that changes the class, and we therefore stop. Finally, if the maximum amount of steps are taken, we also terminate.

The algorithm can be improved in many ways, but our intention is to demonstrate that such adversarial examples are possible and easy to generate, not to make the best algorithm to do so. One possible improvement might for instance be to be able to 673 backtrack if the concept predictions are changed. 674

When tuning hyperparameters for the adversarial  $_{775}$  concept attacks, we chose to tune the step size  $\alpha$ ,  $_{676}$  which is multiplied with the perturbation in each  $_{777}$  step, and the sensitivity threshold  $\gamma$ , which determines how close a concept prediction needs to be to  $_{679}$  change before we try to cancel out its changes.  $_{680}$ 

For the ConceptShapes datasets, we used a grid 681 search with step size  $\alpha \in \{0.003, 0.001, 0.00075\}$  and 682 sensitivity threshold  $\gamma \in \{0.1, 0.05, 0.01\}$ . For CUB, 683 the values were  $\alpha \in \{0.0001, 0.000075, 0.00005\}$  and 684  $\gamma \in \{0.1, 0.075, 0.05, 0.02\}$ . These values were cho-685 sen after some initial experimentation. The best val-686 ues were  $\alpha = 0.001, \gamma = 0.1$  and  $\alpha = 0.000075, \gamma =$ 687 0.1, respectively. In order to reduce the running 688 time, we sampled 200 images from the training set 689 that was correctly predicted. We used 800 max steps 690 for ConceptsShapes and 300 for CUB. 691

We ran the grid search with  $\beta = -0.3$ , which is 692 the weight multiplied with pixels that would make 693 concept predictions change, and  $\epsilon = 1$ , which deter-694 mines where the adversarial images get projected 695 back on. After the grid search, we performed a line 696 search on  $\beta \in \{0.1, 0, -0.1, -0.3, -0.5, -0.7, -1\}$ . 697 The results were very similar, but slightly better for 698  $\beta = -0.1$ . The success rates were calculated using 699 all of the images in the test set where the model 700 originally predicted the correct class. 701

# Algorithm 1: Adversarial Concept Attack Algorithm

NLDL

#3

**Result:** Perturbed image  $\widetilde{\mathbf{x}}$  of  $\mathbf{x}$ , such that CMB h misclassifies  $\tilde{\mathbf{x}}$ , but the concept predictions are the same for  $\widetilde{\mathbf{x}}$  and  $\mathbf{x}$ , or 0 for a failed run. Input: Input image  $\mathbf{x} \in \mathbb{R}^d$ . Class label  $y \in [1, \ldots, p]$ . CBM  $h: \mathbb{R}^d \to \mathbb{R}^p$  with input-to-concept function  $g: \mathbb{R}^d \to \mathbb{R}^k$ , such that  $\operatorname{argmax}(h(\mathbf{x})) = y.$ Sensitivity threshold  $\gamma \in [0, \infty)$ . Step size  $\alpha \in (0, 1)$ . Deviation threshold  $\epsilon \in \mathbb{R}^d$ . Max iterations  $t_{\max} \in \mathbb{N}_1$ . Gradient weight  $\beta \in [-1, 0]$ . Valid pixel range  $[x_{\min}, x_{\max}]$ .  $\widetilde{\mathbf{x}}_0 \leftarrow \mathbf{x} \mathrel{/\!/} \mathsf{Adversarial}$  example  $\mathbf{\hat{c}} = g(\mathbf{x})$  // Original concept logits  $\mathbf{\hat{c}}_b = \mathbb{I}(\sigma(\mathbf{\hat{c}}) > 0.5)$  // Original binary predictions for t = 0 to  $t_{max}$  do  $\widetilde{\mathbf{c}} \leftarrow g(\widetilde{\mathbf{x}}_t)$  $\widetilde{\mathbf{c}}_b = \mathbb{I}(\sigma(\widetilde{\mathbf{c}}) > 0.5)$  // New concept predictions if  $\widetilde{\mathbf{c}}_b \neq \mathbf{\hat{c}}_b$  then return 0 // Changed concept predictions if  $\operatorname{argmax}(h(\widetilde{\mathbf{x}}_t)) \neq y$  then return  $\widetilde{\mathbf{x}}_t$  // Success  $\mathbf{\hat{p}}_t = \operatorname{sign}(\nabla_{\mathbf{\widetilde{x}}_t} L(h(\mathbf{\widetilde{x}}_t), \mathbf{y})) // \text{Initial}$ perturbation Initialize  $M_t \in \mathbb{R}^d$  with all elements as ones for j = 0 to k do if  $\tilde{c}_i$  in  $[-\gamma, \gamma]$  then  $\mathbf{q}_j = \operatorname{sign}(\nabla_{\widetilde{\mathbf{x}}_t} g(\widetilde{\mathbf{x}}_t)_j)$  $M_{t,j} \leftarrow \mathbb{I}_{\beta}(\hat{\mathbf{p}}_t \odot \mathbf{q}_j \neq \operatorname{sign}(g(\mathbf{x})_j))$  $M_t \leftarrow \min(M_t, M_{t,i})$ if All entries in  $M_t$  equal  $\beta$  then return 0 // All  $\beta$  mask  $\mathbf{p}_t = \mathbf{\hat{p}}_t \odot M_t$  // Final perturbation  $\widetilde{\mathbf{x}}' = \Pi_{[\mathbf{x}-\epsilon,\mathbf{x}+\epsilon]}(\widetilde{\mathbf{x}}_t + \alpha \mathbf{p}_t) // \text{Projection}$  $\widetilde{\mathbf{x}}_{t+1} = \operatorname{clamp}(\widetilde{\mathbf{x}}', x_{\min}, x_{\max})$ return 0 // Max iterations exceeded



Table B.2. Overview of the four concepts that relate to the background. These four concepts are only present in the 9-concept datasets. Each row corresponds to one concept. The two leftmost columns of images have the concept, while the two rightmost columns have not.

# D Hyperparameter Details 702

For the ConceptShapes datasets, the learn- 703 rates were sampled from values in ing 704  $\{0.05, 0.01, 0.005, 0.001\}$ , and dropout proba-705 bilities from  $\{0, 0.2, 0.4\}$ . The standard model also 706 searched for an exponential decay parameter to 707 the linear learning scheduler in  $\{0.1, 0.5, 0.7, 1\}$ , 708 applied every five epochs. The concept-based 709 models set the decay parameter to 0.7 and 710 searched for a weight balancing the concept loss 711 function and the class loss function. They were 712  $\{(100, 0.8), (100, 0.9), (5, 1), (10, 1)\},$  where the 713 first element in the tuples represents the weight 714 multiplied with the concept loss, and the second 715 is an exponential decay parameter, applied to the 716 weight every epoch. All of the models had an equal 717 amount of hyperparameter trials. 718

For the CUB dataset, we sat the dropout probability to 0.15 and searched for learning rates 720 in  $\{0.001, 0.0005, 0.0001\}$ . The standard model 721 searched for the exponential learning rate decay 722 parameter in  $\{0.1, 0.7\}$ , applied every ten epochs. 723 The concept-based models set this 1 and searched 724 for concept weight in  $\{10, 15\}$ . 725

The oracle models had a learning rate of 0.01 for 726 ConceptShapes and 0.001 for CUB. They quickly 727 converged and did not require further hyperparameter tuning. The grid search was implemented with 729 the python library Optuna [32]. 730

Model	Total Parameters	Trainable Parameters	Frozen Parameters	
ConceptShapes models with 10 classes and 5 concepts				
Standard model	139 578	139 578	0	
Vanilla CBM	137 583	137 583	0	
CBM-Res	138 527	138 527	0	
CBM-Skip	138 399	138 399	0	
SCM	152 849	152 849	0	
ConceptShapes models with 21 classes and 9 concepts				
Standard model	139 941	139 941	0	
Vanilla CBM	138 053	138 053	0	
CBM-Res	139 758	139 758	0	
CBM-Skip	139 614	139 614	0	
SCM	167 579	167 579	0	
CUB models				
Standard model	11 425 032	248 520	11 176 512	
Vanilla CBM	11 371 880	195 368	11 176 512	
CBM-Res	11 416 952	240 440	11 176 512	
CBM-Skip	11 428 088	251 576	11 176 512	
SCM	11 786 488	609 976	11 176 512	

Table E.1. Amount of parameters in the different models. There are many variations of the Concept-Shapes datasets, here we show the one resulting in the fewest parameters (top) and the most parameters (middle).

# 731 E Training Details

We used the Adam optimizer [33] with a linear learn-732 ing rate scheduler and the Pytorch deep learning 733 library [34]. The standard model was constructed to 734 be as similar as possible to the concept-based mod-735 els. Instead of a bottleneck layer, it had an ordinary 736 linear layer. We added dropout [35] after the convo-737 lutional part of the models. The amount of model 738 parameters can be found in Table E.1. We used a 739 cross entropy loss function as the class loss and a 740 binary cross entropy loss function for the concepts. 741 The models trained on CUB used a frozen 742 Resnet 18 [28], pre-trained on Imagenet [36], as 743 the base model, while the ConceptShapes models 744 were trained from scratch. We used three layers of 745  $3\times 3$  convolutional blocks, with padding and  $2\times 2$ 746 maxpooling. We used 256 nodes in the first linear 747 layer for the models trained on CUB, and 64 nodes 748 for the models trained on ConceptShapes. 749

The pixel values were scaled down to [0, 1] and 750 normalized. The models trained on CUB used Ima-751 genet [36] normalization parameters, and the models 752 on the ConceptShapes datasets used means of 0.5 753 and standard deviations of 2 for all channels. We 754 performed random cropping on the training images, 755 and center cropping when evaluating and testing. 756 We did not perform any data augmentation that 757 changed colors of the images, in order to not inter-758 fere with the concepts relating to colors. 759