
Post-hoc Stochastic Concept Bottleneck Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Concept Bottleneck Models (CBMs) are interpretable models that predict the
2 target variable through high-level human-understandable concepts, allowing users
3 to intervene on mispredicted concepts to adjust the final output. While recent
4 work has shown that modeling dependencies between concepts can improve CBM
5 performance, especially under interventions, such approaches typically require
6 retraining the entire model, which may be infeasible when access to the original
7 data or compute is limited. In this paper, we introduce Post-hoc Stochastic Concept
8 Bottleneck Models (PSCBMs), a lightweight method that augments any pre-trained
9 CBM with a multivariate normal distribution over concepts by adding only a small
10 covariance-prediction module, without retraining the backbone model. We propose
11 two training strategies and show on real-world data that PSCBMs consistently
12 match or improve both concept and target accuracy over standard CBMs at test
13 time. Furthermore, we show that due to the modeling of concept dependencies,
14 PSCBMs perform much better than CBMs under interventions, while remaining
15 far more efficient than retraining a similar stochastic model from scratch.

16 1 Introduction

17 The adoption of machine learning models in high-stakes domains is often limited by their lack of
18 interpretability due to their black-box nature [Lipton, 2016, Doshi-Velez and Kim, 2017]. Concept
19 Bottleneck Models (CBMs), first introduced by Koh et al. [2020], address this by inserting a layer
20 of neurons aligned with human-understandable concepts before the target prediction. A CBM is
21 composed of a concept encoder that predicts the concepts, and a prediction head that takes these
22 concepts as inputs; the final prediction is thus explained through the intermediate concept values.
23 While the original formulation assumes independence among concepts, accounting for correlations
24 has been shown to improve the model’s performance [Havasi et al., 2022, Singhi et al., 2024,
25 Vandenhirtz et al., 2024]. However, existing approaches that capture concept dependencies typically
26 require training the entire model with dedicated objectives. In contrast, we demonstrate that such
27 dependencies can be incorporated post-hoc into a pre-trained CBM. A more detailed discussion of
28 related work is provided in Appendix A.

29 Building on this idea, we introduce Post-Hoc Stochastic Concept Bottleneck Models (PSCBMs),
30 which extend pre-trained CBMs with a lightweight module that predicts concept covariance. Inspired
31 by Stochastic Concept Bottleneck Models (SCBMs) [Vandenhirtz et al., 2024], we model concept
32 dependencies with a multivariate normal distribution. This approach not only improves predictive
33 performance, but also enables more efficient concept interventions—the process by which a user
34 modifies predicted concepts to directly influence downstream predictions.

35 This work contributes to model interpretability through human-understandable concepts in three ways:
36 (i) We introduce PSCBM, a lightweight post-hoc module that incorporates concept correlations into
37 existing CBMs without requiring full retraining, thereby reducing data and compute requirements.
38 (ii) We propose a simple intervention-based training procedure that further improves intervention

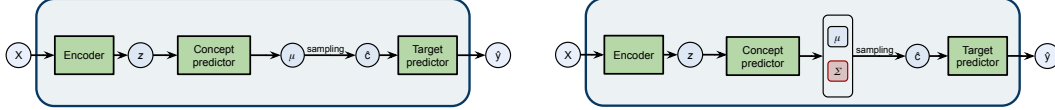


Figure 1: Comparison of CBM (left), SCBM (right), and our proposed PSCBM (right, red block). All methods input x to an encoder to produce the feature vector z from which concepts are obtained and fed to the target predictor. CBM directly predicts concept values \hat{c} , while in SCBM, the predictor outputs an expected value μ and a covariance matrix Σ that define a multivariate normal distribution to sample from. PSCBM incorporates the Σ predictor (red box) to a pre-trained CBM.

39 efficiency in (P)SCBM-like models. (iii) We show on real-world data that PSCBM improves both
 40 predictive accuracy and intervention effectiveness while remaining substantially more efficient than
 41 full model retraining.

42 2 Method

43 **Modeling concept dependencies with a normal distribution** In this paper we propose PSCBMs¹, a
 44 post-hoc extension of CBMs that enables to model dependencies between concepts without retraining
 45 the model from scratch. A CBM Koh et al. [2020] consists of three main components: a feature
 46 encoder, a concept predictor, and a target predictor. The feature encoder h (e.g., a CNN or ResNet)
 47 extracts features $z = h(x)$ from an input x . Then the concept predictor g outputs concept probabilities
 48 $p = g(z)$. When focusing on binary concepts for downstream prediction (i.e. *hard* CBMs [Havasi
 49 et al., 2022]), rather than their probabilities (*soft* CBMs), the next step is sampling M concept values
 50 from a Bernoulli distribution $\hat{c}^{(1)}, \dots, \hat{c}^{(M)} \sim \text{Bernoulli}(p)$. Then the target predictor f maps the
 51 sampled concepts to the final prediction $\hat{y} = \frac{1}{M} \sum_{m=1}^M f(\hat{c}^{(m)})$. SCBMs [Vandenhirtz et al., 2024]
 52 extend this framework by explicitly capturing correlations between concepts: the concept predictor is
 53 extended to feature a mean predictor g_μ and a covariance predictor g_Σ , which define a multivariate
 54 normal distribution that models concept dependencies. Concepts are sampled as $\hat{c}^{(m)} = \sigma(\eta^{(m)})$
 55 for $m = 1, \dots, M$, where $\eta^{(m)} \sim \mathcal{N}(\mu, \Sigma)$ and σ is the sigmoid function. Our proposed PSCBMs
 56 directly build upon this idea: given a pre-trained CBM, we reuse its concept predictor g as g_μ and
 57 add a lightweight covariance predictor g_Σ , training only the latter while keeping the rest of the model
 58 frozen. In this way, any existing CBM can be turned into a stochastic, dependency-aware model
 59 post-hoc, as illustrated in Figure 1.

60 **Interventions in Stochastic Concept Bottleneck Models** A distinguishing property of CBMs is
 61 that they allow users to modify predicted concept values at test time and thereby directly update the
 62 downstream prediction. This process, known as intervention, can be factorized into two parts: the
 63 policy π , which selects the concepts to be intervened on (e.g., randomly or based on uncertainty),
 64 and the strategy τ , which determines how their values are updated. We describe several policies
 65 and strategies in Appendix B. Intervening in SCBMs or PSCBMs additionally accounts for concept
 66 dependencies and proceeds in four steps:

- 67 1. **Select concepts for intervention** Use a policy π to choose a subset \mathcal{S} of concepts.
- 68 2. **Set intervened logits** Assign values to $\eta'_{\mathcal{S}}$, the logits of intervened concepts, according to a
 69 chosen intervention strategy τ .
- 70 3. **Update remaining logits** Compute the conditional normal distribution parameterized by
 71 $\bar{\mu}, \bar{\Sigma}$ for the non-intervened concept set $\setminus \mathcal{S}$, using the equations for a conditional normal
 72 distribution: $\eta_{\setminus \mathcal{S}} \mid x, \eta'_{\mathcal{S}} \sim \mathcal{N}(\bar{\mu}, \bar{\Sigma})$,
 73 where $\bar{\mu} = \mu_{\setminus \mathcal{S}} + \Sigma_{\setminus \mathcal{S}, \mathcal{S}} \Sigma_{\mathcal{S}, \mathcal{S}}^{-1} (\eta'_{\mathcal{S}} - \mu_{\mathcal{S}})$, $\bar{\Sigma} = \Sigma_{\setminus \mathcal{S}, \setminus \mathcal{S}} - \Sigma_{\setminus \mathcal{S}, \mathcal{S}} \Sigma_{\mathcal{S}, \mathcal{S}}^{-1} \Sigma_{\mathcal{S}, \setminus \mathcal{S}}$.
- 74 4. **Sample probabilities** Sample binary concept values $c_{\setminus \mathcal{S}}$ from the logits $\eta_{\setminus \mathcal{S}}$.

75 **Learning the covariance matrix post-hoc** The covariance predictor g_Σ is trained by minimizing
 76 the loss function of a regular SCBM. It teaches the model to predict both the concepts and target
 77 variable correctly and encourages sparsity in the covariance matrix:

¹The code is available at <https://anonymous.4open.science/r/PSCBM-B32D>.

$$\mathcal{L} = \underbrace{-\log \sum_{m=1}^M \exp \sum_{i=1}^C -\text{BCE} \left(c_i, \sigma(\eta_i^{(m)}) \right)}_{\text{Concept Loss}} + \lambda_1 \underbrace{\text{CE} \left(y, \frac{1}{M} \sum_{m=1}^M g_{\psi}(\mathbf{c}^{(m)}) \right)}_{\text{Target Loss}} + \lambda_2 \underbrace{\sum_{i \neq j} \Sigma_{i,j}^{-1}}_{\text{Regularization}}.$$

78 Here, M is the number of Monte Carlo concept samples, C is the number of concepts, BCE and CE
 79 denote the (Binary) Cross-Entropy, \mathbf{x} is the input, y the true label, c_i the true concept values, $\eta^{(m)}$
 80 the sampled logit vectors, $\sigma(\cdot)$ the sigmoid function, $\mathbf{c}^{(m)}$ the binary concept values sampled from
 81 the Bernoulli distribution defined by $\sigma(\eta^{(m)})$, $\Sigma(\mathbf{x})$ the predicted covariance matrix and $g(\cdot)$ the
 82 target predictor. Finally, λ_1 and λ_2 are weighting parameters regulating the loss strength.

83 We propose two training paradigms for PSCBM, which can be extended to SCBM: (i) The loss
 84 function is applied to the model’s predictions without concept interventions. (ii) We encourage more
 85 responsive predictors to concept interventions. In every training iteration, for every sample in the
 86 training dataset, we randomly select a set \mathcal{S} of concepts to intervene on using some strategy τ and
 87 calculate the loss after the intervention. The cardinality of \mathcal{S} is fixed throughout training. This is done
 88 N times for each data point and the average of the loss for all these interventions is calculated. This
 89 method is similar to *RandInt* introduced by Espinosa Zarlenga et al. [2022] but it differs in two ways:
 90 (i) Unlike *RandInt*, which decides independently for each concept whether to replace it with its true
 91 value, our method enforces that a fixed number of concepts is intervened on in every case. (ii) At each
 92 training iteration, multiple random interventions per sample are made and the loss is averaged. This
 93 reduces gradient variance and yields more stable training compared to applying a single intervention.
 94 A detailed pseudocode for the calculation of this loss can be found in Appendix C.

95 3 Results

96 In our experiments, we address two main questions: **1.** How does post-hoc learning of the covariance
 97 affect the test accuracy without interventions? **2.** How does it affect predictions after interventions?

98 **Experimental setup** We compare the performance of PSCBM to SCBM and regular CBM. PSCBM
 99 indicates training without concept interventions, while PSCBMi includes interventions during training.
 100 SCBM proposes two ways to learn the covariance matrix, amortized per instance, $\Sigma(\mathbf{x})$, or learnt
 101 globally, Σ . For brevity, we present the results on the global covariance, and include their amortized
 102 counterparts in Appendix D. All implementation details are described in Appendix E. For test-time
 103 interventions, concepts are selected based on their associated prediction uncertainty, i.e. the concept
 104 with predicted probability closest to 0.5 is chosen. Unlike CBM where the intervened probabilities
 105 are set to 0/1, in SCBMs and PSCBMs concept logits are set to thus of ε or $1 - \varepsilon$, respectively, where
 106 ε is small for stability. We evaluate our methods on the Caltech-UCSD Birds-200-2011 dataset [Wah
 107 et al., 2011], composed of photographs of birds of 200 different classes. We use the variant introduced
 108 in Koh et al. [2020] for CBMs with 112 binary per-class concepts and, as evaluation metrics, we
 109 focus on concept and target accuracy.

110 **Test performance** In Table 1 we present concept and target accuracy of the evaluated models
 111 without interventions. The results show that adding our covariance learning module improves
 112 performance. Both variants of PSCBM outperform SCBM on both concept and target accuracy. For
 113 target accuracy, PSCBM trained without interventions clearly surpasses the regular CBM, while on
 114 concept accuracy the two achieve comparable results. The training times presented in the last column
 115 show that training a PSCBM without interventions is computationally inexpensive due to the fact that
 only one module is trained. Training with interventions however requires significantly more time.

Table 1: Test-time target and concept accuracy without interventions and their AUCs under inter-
 ventions. AUC is normalised to stay within the interval [0,1]. We report the mean and standard
 deviation over three runs for baselines and nine runs for PSCBM. The best scores for each metric are
 highlighted with **bold font**, and those within a standard deviation from the best one are underlined.

Method	Target Accuracy (%)	Concept Accuracy (%)	Target Accuracy AUC	Concept Accuracy AUC	Training time (s)
CBM	67.3973 ± 0.5722	94.9403 ± 0.1059	0.9551 ± 0.0000	0.9825 ± 0.0000	7204 ± 247
SCBM	65.5173 ± 0.8539	94.4569 ± 0.1328	0.9671 ± 0.0001	0.9870 ± 0.0000	8134 ± 767
PSCBM	68.3983 ± 0.1992	<u>94.9285</u> ± 0.0206	0.9680 ± 0.0003	0.9859 ± 0.0001	740 ± 94
PSCBMi	<u>68.1970</u> ± 0.1274	<u>94.9026</u> ± 0.0350	0.9704 ± 0.0002	0.9866 ± 0.0001	14084 ± 267

117 **Intervention performance** To evaluate intervention performance we rely on two approaches: visual
 118 inspection of the intervention curves, and the area under the intervention curve (AUC) [Singhi et al.,
 119 2024], which summarizes each curve into a single value for easier comparison. The last two columns
 120 of Table 1 report the AUC scores. SCBM achieves the highest concept AUC, but both PSCBM
 121 variants outperform it on target AUC, with PSCBMi being the strongest overall. This highlights the
 122 benefit of training with interventions. Both PSCBM models also outperform the regular CBM on
 123 both metrics, showing that adding a covariance matrix post-hoc improves intervention performance.
 124 For a more detailed comparison, we examine the intervention trajectories under the uncertainty
 125 policy in Figure 2. We additionally include the corresponding results when selecting concepts with
 126 a random policy in Appendix F. Both PSCBM variants consistently outperform CBM, with the
 127 intervention-trained variant performing best. However, for a few interventions, PSCBM without
 128 interventions during training lags behind SCBM, indicating that joint training of the covariance
 129 is more effective. PSCBM with interventions surpasses SCBM on target accuracy after about 20
 130 interventions, though it does not reach SCBM’s performance on concept accuracy. We note that
 131 SCBM could also be trained with interventions, which might further boost its results. It is clear from
 132 these that adding a covariance on a trained CBM can significantly improve its performance.

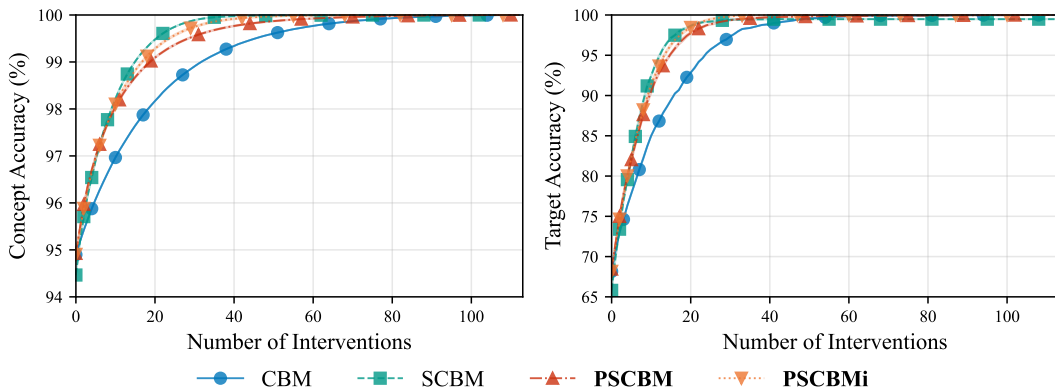


Figure 2: Intervention curves of Concept and Target Accuracy for PSCBM models and baselines when concept uncertainty policy is used. Confidence intervals are thinner than the lines.

133 4 Conclusion

134 In this paper, we introduce PSCBM, a lightweight extension to CBMs that models concept dependen-
 135 cies via a multivariate normal distribution. Crucially, PSCBM does not require retraining the full
 136 model: it only adds a compact covariance module. This makes it far less computationally demanding
 137 than training an SCBM from scratch, while retaining the benefits of modeling concept dependen-
 138 cies. We proposed two training procedures for the covariance module—optimizing the loss without
 139 interventions, and optimizing it after intervening on a random subset of concepts. Both variants
 140 outperform standard CBMs in test-time accuracy, and under interventions they adapt significantly
 141 faster than regular CBMs, though not as rapidly as a fully retrained SCBM. Importantly, training-time
 142 interventions improve intervention efficiency without harming accuracy when no interventions are in
 143 place, which demonstrates the efficacy of our approach.

144 Apart from its lightweight nature and efficacy, PSCBM also guarantees compatibility with the original
 145 CBM. By disabling the covariance module, PSCBM can revert to identical predictions as the baseline
 146 CBM. This property is particularly valuable in regulated domains such as healthcare, where a CBM
 147 may already have undergone approval (e.g., FDA testing). In such cases, retraining an SCBM might
 148 be prohibited or undermine user trust, whereas PSCBM preserves the validated predictions while still
 149 enabling stronger interventions when permitted.

150 The main limitation of our study is its evaluation on a single dataset; extending the analysis to
 151 additional benchmarks will be important for robustness. Future work should explore richer training-
 152 with-interventions schemes (e.g., varying the number of intervened concepts or using non-random
 153 policies), and investigate their applicability in regular SCBMs. Nonetheless, our results show that
 154 when retraining from scratch is infeasible, augmenting an existing CBM with a learned covariance
 155 matrix provides a simple and efficient way to substantially improve intervention effectiveness.

References

- 156
- 157 J. Ansel, E. Yang, H. He, N. Gimelshein, A. Jain, M. Voznesensky, et al. Pytorch 2: Faster machine
158 learning through dynamic python bytecode transformation and graph compilation. In Proceedings
159 of the 29th ACM International Conference on Architectural Support for Programming Languages
160 and Operating Systems, volume 2, pages 929–947, 2024.
- 161 K. Chauhan, R. Tiwari, J. Freyberg, P. Shenoy, and K. Dvijotham. Interactive concept bottleneck
162 models. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pages
163 5948–5955, 2023.
- 164 F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. arXiv
165 preprint arXiv:1702.08608, March 2017. doi: 10.48550/arXiv.1702.08608.
- 166 M. Espinosa Zarlenga, P. Barbiero, G. Ciravegna, G. Marra, F. Giannini, M. Diligenti, et al. Con-
167 cept embedding models: Beyond the accuracy-explainability trade-off. In Advances in Neural
168 Information Processing Systems, volume 35, pages 21400–21413, 2022.
- 169 M. Espinosa Zarlenga, K. Collins, K. Dvijotham, A. Weller, Z. Shams, and M. Jamnik. Learning to
170 receive help: Intervention-aware concept embedding models. In Advances in Neural Information
171 Processing Systems, volume 36, 2024.
- 172 M. Havasi, S. Parbhoo, and F. Doshi-Velez. Addressing leakage in concept bottleneck models. In
173 Advances in Neural Information Processing Systems, 2022.
- 174 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In 3rd International Conference
175 on Learning Representations, ICLR 2015, San Diego, CA, USA, 2015.
- 176 P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, and P. Liang. Concept bottleneck
177 models. In H. D. III and A. Singh, editors, Proceedings of the 37th International Conference on
178 Machine Learning, volume 119, pages 5338–5348. PMLR, 2020.
- 179 Sonia Laguna, Ričards Marcinkevičs, Moritz Vandenhirtz, and Julia Vogt. Beyond concept bottleneck
180 models: How to make black boxes intervenable? Advances in neural information processing
181 systems, 37:85006–85044, 2024.
- 182 Z. C. Lipton. The mythos of model interpretability. Communications of the ACM, 61(10):35–43,
183 June 2016.
- 184 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. International Conference
185 on Learning Representations, 2019.
- 186 A. Mahinpei, J. Clark, I. Lage, F. Doshi-Velez, and W. Pan. Promises and pitfalls of black-box
187 concept learning models. 38 th International Conference on Machine Learning, 2021.
- 188 Mikael Makonnen, Moritz Vandenhirtz, Sonia Laguna, and Julia E Vogt. Measuring leakage in
189 concept-based methods: An information theoretic approach. ICLR 2025 Workshop: XAI4Science,
190 2025.
- 191 A. Margeloiu, M. Ashman, U. Bhatt, Y. Chen, M. Jamnik, and A. Weller. Do concept bottleneck
192 models learn as intended? Workshop paper at ICLR 2021, 2021.
- 193 T. Oikarinen, S. Das, L. M. Nguyen, and T.-W. Weng. Label-free concept bottleneck models. In The
194 11th International Conference on Learning Representations, 2023.
- 195 A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, et al. Learning transferable visual
196 models from natural language supervision. In International Conference on Machine Learning,
197 pages 8748–8763, 2021.
- 198 S. Shin, Y. Jo, S. Ahn, and N. Lee. A closer look at the intervention procedure of concept bottleneck
199 models. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors,
200 Proceedings of the 40th International Conference on Machine Learning, volume 202, pages 31504–
201 31520. PMLR, 2023.
- 202 S. Silvey. Statistical Inference. Taylor & Francis, 1975.

- 203 Nishad Singhi, Jae Myung Kim, Karsten Roth, and Zeynep Akata. Improving intervention efficacy
204 via concept realignment in concept bottleneck models. In European Conference on Computer
205 Vision, pages 422–438. Springer, 2024.
- 206 Moritz Vandenhirtz, Sonia Laguna, Ričards Marcinkevičs, and Julia E. Vogt. Stochastic concept
207 bottleneck models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak,
208 and C. Zhang, editors, Advances in Neural Information Processing Systems, volume 37, pages
209 51787–51810. Curran Associates, Inc., 2024.
- 210 C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011
211 dataset, 2011. Dataset.
- 212 M. Yuksekogonul, M. Wang, and J. Zou. Post-hoc concept bottleneck models. In The 11th International
213 Conference on Learning Representations, 2023.

214 A Related Work

215 **Concept Bottleneck Models** Concept Bottleneck Models (CBM) have originally been proposed by
216 Koh et al. [2020]. The idea is to transform a regular neural network into an interpretable model by
217 converting one of its layers into a concept bottleneck, where each neuron corresponds to a high-level
218 concept that can be interpreted by a human. Such a model can be framed as a composite function
219 $f(\mathbf{x}) = g(h(\mathbf{x}))$, where $h(\cdot)$ is a concept predictor that calculates the concepts \hat{c} , and $g(\cdot)$ is the
220 target predictor that predicts the final label from the concepts. Since all information that goes from
221 the input to the output passes through the concept bottleneck, the network’s prediction should be well
222 explained by the concepts. Koh et al. [2020] use what has been called *soft* concept encoding: the
223 target predictor uses real-valued concept logits or probabilities produced by the encoder. However, it
224 has been shown that this method is prone to information leakage, where undesired information can
225 pass through the bottleneck [Margeloiu et al., 2021, Mahinpei et al., 2021, Makonnen et al., 2025].
226 As a remedy, Havasi et al. [2022] propose *hard* concepts: the concept encoder predicts concept
227 probabilities, which parametrize Bernoulli distributions, from which binary inputs for the target
228 predictor are sampled. Additionally, they introduce in their work a side channel which allows the
229 target predictor to use information that cannot be expressed in terms of the high-level concepts. This
230 modification can increase target accuracy if the concept set is incomplete. An alternative approach to
231 modeling the concept bottleneck has been proposed by Espinosa Zarlenga et al. [2022] who introduce
232 Concept Embedding Models (CEM), where each concept is modeled not by a single neuron but
233 instead by a pair of learnable vector embeddings corresponding to the presence or absence of a
234 concept in the input. This allows passing information not captured by the concepts through the
235 bottleneck without the need for a side channel. Several methods have been proposed to transform
236 regular models into Concept Bottleneck Models. Yuksekogonul et al. [2023] propose a method to
237 convert one layer of a neural network trained without concepts into a concept bottleneck by using
238 only a small concept-annotated subset. Oikarinen et al. [2023] demonstrate that concepts can also be
239 generated with the help of a Vision-Language Model such as CLIP [Radford et al., 2021], reducing
240 the human effort associated with annotating data.

241 **Concept interventions** A particular advantage of CBMs is the possibility of affecting downstream
242 predictions by modifying incorrectly predicted concepts. This process is called a concept intervention.
243 Koh et al. [2020] only consider interventions on randomly selected concepts. On the other hand,
244 Chauhan et al. [2023], Shin et al. [2023] propose and evaluate more efficient concept selection
245 policies, such as choosing the concept with the highest associated prediction uncertainty, that is,
246 with probability closest to 0.5. This policy allows reducing target prediction error by half after 12
247 concept interventions on average, while with random concept selection an average of 43 interventions
248 is necessary to achieve it. In Espinosa Zarlenga et al. [2022], a regularization technique called
249 RandInt is introduced, whereby randomly selected concepts are set to their true value during training.
250 This is then shown to increase the model’s responsiveness to test-time concept interventions. As
251 an alternative for an explicit concept selection policy, Espinosa Zarlenga et al. [2024], Singhi et al.
252 [2024] parameterize it by a neural network which can be trained jointly with the model or in a
253 post-hoc manner. A further extension of the interventions framework is done by Laguna et al. [2024],
254 introducing a method for intervening on a black-box model without altering its architecture. They
255 also define a measure of intervenability and show that fine-tuning a model for this quantity can
256 improve the performance of interventions.

257 **Modeling concept dependencies** Since correlated concepts may require multiple interventions,
258 several methods model dependencies. Havasi et al. [2022] use an autoregressive structure, by which
259 the value of every concept depends on the previous concepts. It increases the model’s accuracy, but
260 has the disadvantage that concepts have to be evaluated sequentially, which increases evaluation time.
261 Singhi et al. [2024] augment the CBM with a neural network that is trained to update all concepts in
262 response to an intervention made to one of them. This can be combined with the trainable concept
263 selection policy mentioned above. These elements can be either trained jointly with the full model
264 or added to an existing CBM post-hoc. However, unlike our method, it does not provide an *explicit*
265 representation of concept correlations. Finally, Vandenhirtz et al. [2024] introduce Stochastic Concept
266 Bottleneck Models (SCBM), in which concept correlations are captured by a multivariate normal
267 distribution. It requires predicting not only the expected value of the concepts, but also a covariance
268 matrix. A special advantage of this model is the explicit representation of concept dependencies,
269 which is used in an intervention strategy based on the distribution’s confidence region.

270 **B Comparing Different Intervention Strategies and Policies**

271 In Section 2, we decomposed a concept intervention into the selection of concepts for intervention and
 272 updating their values. The former part we denote as *intervention policy* and the latter - as *intervention*
 273 *strategy*. In the following we describe and evaluate different possible choices. We consider two
 274 intervention policies: *random* concept selection (as in Koh et al. [2020]) and *concept uncertainty*
 275 policy, which selects the concept whose predicted probability is closest to 0.5 (that is, the most
 276 uncertain concept). For intervention strategies we consider four options:

- 277 1. *Empirical Percentile Strategy* sets the value of the intervened-on concept to the 5th (if it is
 278 0) or 95th (if it is 1) percentile of the empirical distribution of concept predictions made by
 279 the model. Koh et al. [2020] use it for intervening in CBM with soft concept encoding.
- 280 2. *Hard Strategy* introduced by Havasi et al. [2022] sets concept probabilities to 0 or 1. In the
 281 case of PSCBM, where concepts are represented by their logits, one cannot directly use the
 282 logits of 0-1 probabilities, as they are infinite. Instead, we set concepts to the logits of ε for
 283 absent concepts and $1 - \varepsilon$ for present concepts, where ε is a small positive number.
- 284 3. *Simple Percentile Strategy* is a softened version of the hard strategy, where concept probabili-
 285 ties of absent or present are set to 0.05 or 0.95, respectively.
- 286 4. *Confidence Region Strategy* is a strategy that directly uses the concepts’ multivariate normal
 287 distribution. If a subset $\mathcal{S} \subset \{1, \dots, \mathcal{C}\}$ of concepts has been selected for intervention, the
 288 updated concept logits $\eta'_\mathcal{S}$ are calculated by solving the following optimization problem:

$$\begin{aligned} \eta'_\mathcal{S} &= \arg \max_{\eta_\mathcal{S}} \log p(\mathbf{c}_\mathcal{S} | \eta_\mathcal{S}), \\ \text{s.t. } &-2 (\log p(\eta_\mathcal{S} | \boldsymbol{\mu}_\mathcal{S}, \boldsymbol{\Sigma}_{\mathcal{S},\mathcal{S}}) - \log p(\boldsymbol{\mu}_\mathcal{S} | \boldsymbol{\mu}_\mathcal{S}, \boldsymbol{\Sigma}_{\mathcal{S},\mathcal{S}})) \leq \chi_{d,1-\alpha}^2 \quad (1) \\ &\eta_i - \mu_i \geq 0 \quad \text{if } c_i = 1, \quad \forall i \in \mathcal{S} \\ &\eta_i - \mu_i \leq 0 \quad \text{if } c_i = 0, \quad \forall i \in \mathcal{S} \end{aligned}$$

289 where \mathcal{S} is the subset of concepts that have been selected for intervention, $\boldsymbol{\mu}_\mathcal{S}$ and $\boldsymbol{\Sigma}_{\mathcal{S},\mathcal{S}}$ are
 290 the predicted concept mean and covariance of this concept subset, c_i are the concept values
 291 selected by the user, d is the dimensionality of the distribution, i.e. the number of concepts,
 292 and α is the requested confidence of the confidence region. In this problem we seek to
 293 find concept logits $\eta'_\mathcal{S}$ which maximize the log-probability of intervened concept values,
 294 and stay within the α -confidence region of the predicted distribution, that is, the model’s
 295 prediction shouldn’t be completely disregarded (first constraint). The log-probability of this
 296 distribution, after being multiplied by -2, asymptotically follows the χ^2 distribution [Silvey,
 297 1975], which gives rise to the first inequality. The other two inequalities follow from the
 298 requirement that the new concept logits shouldn’t move away from the ground truth: if
 299 $c_i = 1$, η_i should increase, and if $c_i = 0$, it should decrease. This last strategy can only be
 300 used in the context of SCBM, while the three others can be used by both SCBM and CBM.

301 In Table 2 we report the AUC for concept and target accuracy for all models and all intervention
 302 policies and strategies that can be used with them. The strategy based on confidence region cannot be
 303 used with a regular CBM, because it requires the covariance matrix.

304 Interestingly, there is no single best strategy. The *Hard* one typically produces the best or close-to-the-
 305 best results for target accuracy, but the *Confidence Region* strategy is often better on concept accuracy,
 306 especially with models using an amortized covariance. Hence, the choice of the best strategy is not
 307 an obvious one, and we encourage the reader to evaluate different policies in relation to the specific
 308 data and setup at hand.

Table 2: Comparison of AUC for all combinations of strategy and policy pairs for all baseline and PSCBM models. The values are normalized to lie in the interval $[0, 1]$. For each model, policy pair, the best value is **bold**, and those lying within a standard deviation are underlined. Likewise, the selected strategy is also **bold** and the one following up is underlined.

Method	Policy	Strategy	Target Accuracy AUC	Concept Accuracy AUC
CBM	Concept Uncertainty	Hard	0.9551 \pm 0.0000	0.9825 \pm 0.0000
		Simple Percentile	0.9513 \pm 0.0000	<u>0.9825</u> \pm 0.0000
		Empirical Percentile	0.9550 \pm 0.0000	0.9798 \pm 0.0002
	Random	Hard	0.8927 \pm 0.0003	<u>0.9658</u> \pm 0.0000
		Simple Percentile	0.8847 \pm 0.0015	0.9659 \pm 0.0001
		Empirical Percentile	<u>0.8927</u> \pm 0.0007	0.9642 \pm 0.0001
SCBM (global covariance)	Concept Uncertainty	Hard	0.9671 \pm 0.0001	0.9870 \pm 0.0000
		Simple Percentile	0.9453 \pm 0.0056	0.9835 \pm 0.0002
		Empirical Percentile	0.9570 \pm 0.0060	0.9864 \pm 0.0001
		Confidence Region	0.9630 \pm 0.0023	0.9866 \pm 0.0001
	Random	Hard	0.9085 \pm 0.0010	0.9697 \pm 0.0001
		Simple Percentile	0.8390 \pm 0.0029	0.9147 \pm 0.0002
SCBM (amortized covariance)	Concept Uncertainty	Hard	0.9646 \pm 0.0000	0.9860 \pm 0.0000
		Simple Percentile	0.9653 \pm 0.0001	0.9873 \pm 0.0000
		Empirical Percentile	0.9653 \pm 0.0002	0.9863 \pm 0.0000
		Confidence Region	0.9664 \pm 0.0001	0.9871 \pm 0.0000
	Random	Hard	0.9234 \pm 0.0005	0.9785 \pm 0.0000
		Simple Percentile	0.9099 \pm 0.0019	0.9509 \pm 0.0005
PSCBM (global covariance)	Concept Uncertainty	Hard	0.9680 \pm 0.0003	0.9859 \pm 0.0001
		Simple Percentile	<u>0.9680</u> \pm 0.0002	0.9859 \pm 0.0001
		Empirical Percentile	0.9665 \pm 0.0005	0.9849 \pm 0.0001
		Confidence Region	0.9657 \pm 0.0002	0.9853 \pm 0.0001
	Random	Hard	0.9081 \pm 0.0014	0.9705 \pm 0.0002
		Simple Percentile	<u>0.9078</u> \pm 0.0012	0.9705 \pm 0.0002
PSCBM (amortized covariance)	Concept Uncertainty	Hard	0.9665 \pm 0.0018	0.9859 \pm 0.0004
		Simple Percentile	0.9665 \pm 0.0018	0.9859 \pm 0.0004
		Empirical Percentile	0.9591 \pm 0.0025	0.9832 \pm 0.0007
		Confidence Region	<u>0.9661</u> \pm 0.0013	<u>0.9855</u> \pm 0.0002
	Random	Hard	<u>0.9129</u> \pm 0.0018	0.9638 \pm 0.0046
		Simple Percentile	<u>0.9127</u> \pm 0.0016	0.9638 \pm 0.0045
PSCBMi (global covariance)	Concept Uncertainty	Hard	<u>0.9708</u> \pm 0.0001	0.9857 \pm 0.0000
		Simple Percentile	0.9708 \pm 0.0002	0.9857 \pm 0.0000
		Empirical Percentile	0.9036 \pm 0.0006	0.9666 \pm 0.0001
		Confidence Region	0.9379 \pm 0.0003	0.9767 \pm 0.0000
	Random	Hard	<u>0.8846</u> \pm 0.0014	0.9552 \pm 0.0003
		Simple Percentile	0.8851 \pm 0.0007	<u>0.9552</u> \pm 0.0002
	Empirical Percentile	0.7797 \pm 0.0021	0.9321 \pm 0.0002	
	Confidence Region	0.8321 \pm 0.0013	0.9475 \pm 0.0001	

309 **C Pseudocode for Interventions Training**

310 In Algorithm 1 we present the pseudocode for one iteration of the intervention-aware training
 311 procedure.

Algorithm 1 Loss calculation for one iteration of training with random interventions

```

1: Inputs:
2:    $g$  (target predictor)
3:    $\mathcal{L}$  (loss function)
4:    $L$  (number of concepts for intervention)
5:    $N$  (number of random interventions per datapoint)
6:    $\tau$  (intervention strategy, which modifies concept values)
7:    $(\mathbf{c}, y)$  (concept and target labels)
8:    $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$  (predicted expected value and covariance of concepts)
9: Algorithm:
10:   $l = 0$  (Initialize the loss to 0)
11: for  $i$  in  $1 \dots N$  do
12:    $S = \text{Random}(L)$  (Randomly select  $L$  concepts for intervention)
13:    $\boldsymbol{\mu}', \boldsymbol{\Sigma}' = \tau(S, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$  (Use strategy  $\tau$  to update the concept distribution)
14:    $\eta' \sim \mathcal{N}(\boldsymbol{\mu}', \boldsymbol{\Sigma}')$  (Sample concept logits from the updated distribution)
15:    $\mathbf{c}' = \sigma(\eta')$  (Pass concept logits through a sigmoid to obtain a vector of probabilities)
16:    $\mathbf{c}_{1:M} = \text{Bern}(\mathbf{c}')$  (Sample  $M$  concept vectors from the Bernoulli distribution defined by  $\mathbf{c}'$ )

17:    $y' = \frac{1}{M} \sum_{m=1}^M g(\mathbf{c}^{(m)})$  (Calculate the target)
18:    $l = l + \mathcal{L}(y', y, \mathbf{c}', \mathbf{c})$  (Loss after intervention)
19: end for
20: return  $\frac{l}{N}$  (Average loss)

```

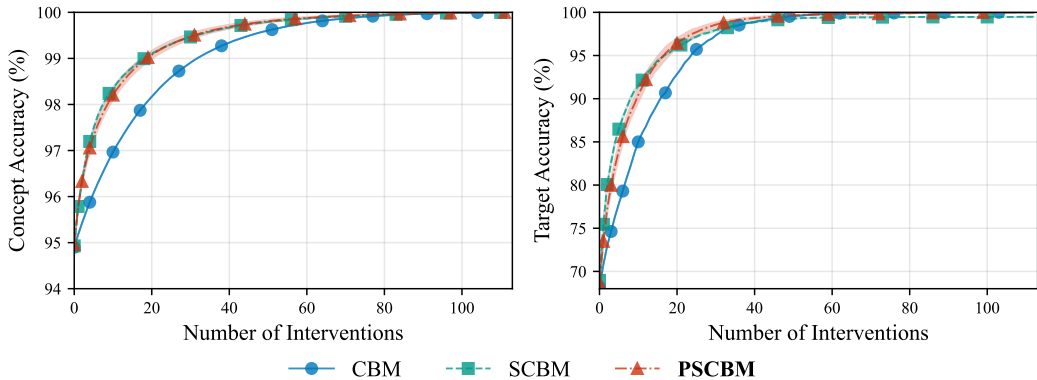
312 **D Results with amortized covariance**

313 In this appendix, we show additional results for the studied models using an *amortized* covariance
 314 matrix instead of a globally learnt one as in the main manuscript. The amortized covariance, as
 315 introduced in SCBM Vandenhirtz et al. [2024], makes the μ and Σ learnt per instance, i.e. amortized
 316 by the input x . We show results for CBM and SCBM as baselines and for a PSCBM trained without
 317 interventions. In Table 3 we show test-time results for target and concept accuracy, as well as target
 318 and concept accuracy AUC for interventions with probability uncertainty policy. In Figure 3, we
 319 show the respective intervention curves. We observe that without interventions, the performance
 320 of PSCBM is better than thus of the baselines. With interventions, it scores the highest on target
 321 accuracy AUC, and on par with the rest on concept accuracy AUC. Intervention curves confirm these
 322 results. However, in the curve for target accuracy, one can observe that after the first few interventions,
 323 the SCBM improved faster than the PSCBM. This is important because in a real-world scenario it is
 324 expected that only a few concept interventions would be made.

Table 3: Test-time target and concept accuracy without interventions and their AUCs during interventions in PSCBM and baselines with an *amortized* covariance matrix. AUC is normalised to stay within the interval [0,1]. We report the mean and standard deviation over three runs. The best scores for each metric are emphasized with **bold** font, and those within a standard deviation from the best one are underlined.

Method	Target Accuracy (%)	Concept Accuracy (%)	Target Accuracy AUC	Concept Accuracy AUC	Training time (s)
CBM	67.3973 ± 0.5722	94.9403 ± 0.1059	0.9551 ± 0.0000	0.9825 ± 0.0000	7204 ± 247
SCBM	<u>68.6342</u> ± 0.3790	<u>94.9190</u> ± 0.0074	0.9646 ± 0.0000	0.9860 ± 0.0000	8130 ± 1073
PSCBM	68.6611 ± 0.2067	94.9605 ± 0.0639	0.9665 ± 0.0018	<u>0.9859</u> ± 0.0004	670 ± 55

Figure 3: Concept and Target Accuracy for amortized PSCBM models and baselines when Concept Uncertainty Policy is used. The shadings represent the standard deviation.



325 **E Implementation Details**

326 All models have been implemented in PyTorch [Ansel et al., 2024] and optimized using Adam
 327 [Kingma and Ba, 2015]. To apply weight decay, we used the optimizer AdamW, which decouples
 328 weight decay from gradient computation [Loshchilov and Hutter, 2019]. All models use a ResNet-18
 329 encoder. The remaining parts including the concept mean and covariances, and the target predictor
 330 are all linear layers. All baselines have been trained with three random initializations. All PSCBM
 331 versions have been evaluated with each of the three CBMs with three random seeds, which amounts
 332 to a total number of nine evaluations. The computations have been done on a cluster composed
 333 mostly of NVIDIA GeForce RTX 2080 GPUs, and every job used a single GPU and two CPUs. When
 334 we train PSCBM with interventions, we use 20 random masks per datapoint per epoch, and found
 335 that a larger number does not lead to better performance. We intervene on 20% of the total number
 336 of concepts, and their values at training time are set according to the *Hard* Strategy, introduced in
 337 Appendix B. We performed hyperparameter optimization with respect to learning rate scheduling,
 338 initial learning rate, and weight decay factor. In Table 4 we report the optimal values based on
 339 validation loss after the last training epoch. It should be noted that PSCBMa is not very sensitive
 340 to weight decay and that for PSCBMg the advantage of step learning rate over cosine schedule was
 341 minimal, albeit statistically significant. Similarly, for CBM and for both SCBM variants, all four
 342 tested hyperparameter combinations achieved similar performance.

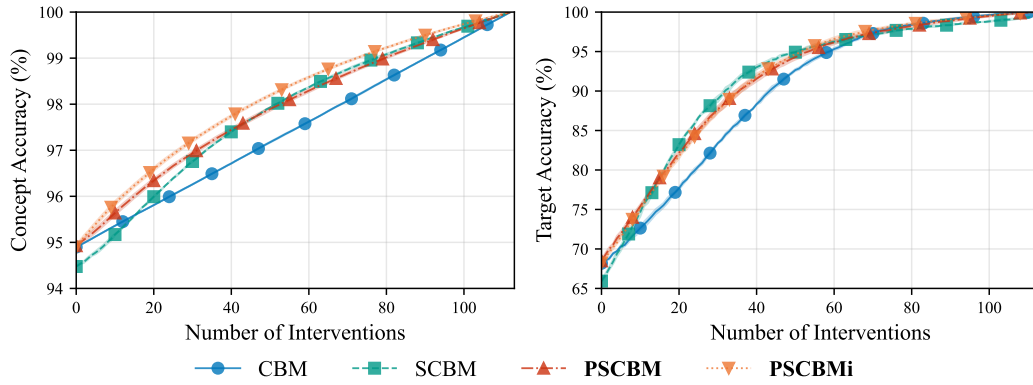
Table 4: Optimal hyperparameter values for each model

Model	Learning rate scheduler	Initial learning rate	Weight decay
CBM	Step-wise	10^{-4}	0.1
SCBMg	Step-wise	10^{-4}	0.1
SCBMa	Step-wise	10^{-4}	0.1
PSCBMg	Step-wise	10^{-3}	1
PSCBMa	Cosine	10^{-4}	0
PSCBMg _i	Cosine	10^{-4}	4

343 **F Results with Random Policy**

344 In Figure 4 we show the intervention curves when concepts for intervention are chosen randomly.
345 One can observe generally lower improvement rates in comparison to the selection policy based
346 on concept uncertainty. On concept accuracy, both PSCBM variants outperform SCBM. On target
accuracy, they outperform the regular CBM but fall slightly behind the SCBM.

Figure 4: Concept and Target Accuracy for PSCBM models and baselines when Random policy is used. The shadings represent the standard deviation.



347