

## A Details of the Intervention Procedure

In Algorithm 1 we describe the standard process of performing interventions in concept-based models using an intervention policy.

---

### Algorithm 1 Intervention Algorithm

---

```

1: Inputs:
2:    $T$  (total number of interventions)
3:    $\pi$  (intervention policy, which takes the concepts as input)
4:    $\hat{c}$  (concepts predicted by the concept encoder)
5:    $\tilde{c} \leftarrow \hat{c}$  ▷ output of the concept encoder,  $g$ 
6: for  $t \in \{0, \dots, T - 1\}$  do
7:    $i \leftarrow \pi(\tilde{c})$  ▷  $i$  is the concept that we want the user to intervene on
8:    $\tilde{c}_i \leftarrow c_i$  ▷ replace the  $i$ th concept in  $\tilde{c}$  with its ground truth value  $c_i$ 
9: end for
10: return  $\tilde{y} = f(\tilde{c})$  ▷ updated class prediction after all interventions have been performed
    
```

---

In Algorithm 2 we describe the procedure used in our setup, which realigns unintervened concepts following an intervention step. We use this algorithm to compute the loss for training the realignment model.

---

### Algorithm 2 Realignment Model Training Loss

---

```

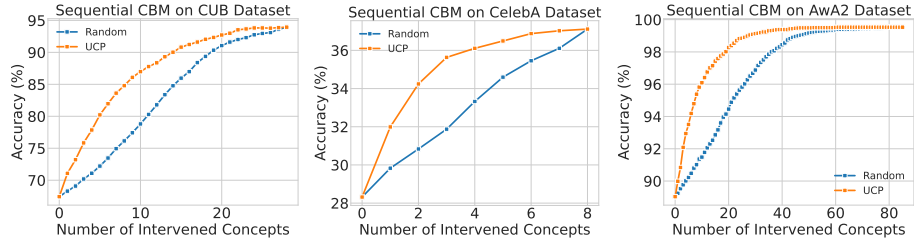
1: Inputs:
2:    $T$  (total number of interventions)
3:    $\pi$  (intervention policy, which takes the concepts as input)
4:    $\hat{c}$  (concepts predicted by the concept encoder)
5:    $\tilde{c} \leftarrow \hat{c}$  ▷ output of the concept encoder,  $g$ 
6:    $\kappa_{-1} \leftarrow \hat{c}$  ▷ initialize realigned concepts
7:    $\mathcal{L} \leftarrow 0$  ▷ initialize loss
8: for  $t \in \{0, \dots, T - 1\}$  do
9:    $i \leftarrow \pi(\kappa_{t-1})$  ▷  $i$  is the concept that we want the user to intervene on
10:   $\tilde{c}_i \leftarrow c_i$  ▷ replace the  $i$ th concept in  $\tilde{c}$  with its ground truth value  $c_i$ 
11:   $\kappa_t \leftarrow u(\tilde{c})$  ▷ output of realignment model
12:   $\mathcal{L} \leftarrow \mathcal{L} + \text{CE}(\kappa_t, c)$  ▷ aggregate loss
13: end for
14: return  $\mathcal{L}/T$  ▷ average loss across all intervention steps
    
```

---

## B Comparison Between Random and UCP Policies

In this section, we compare the classification accuracies achieved by following the random and UCP intervention policies on the three datasets, respectively.

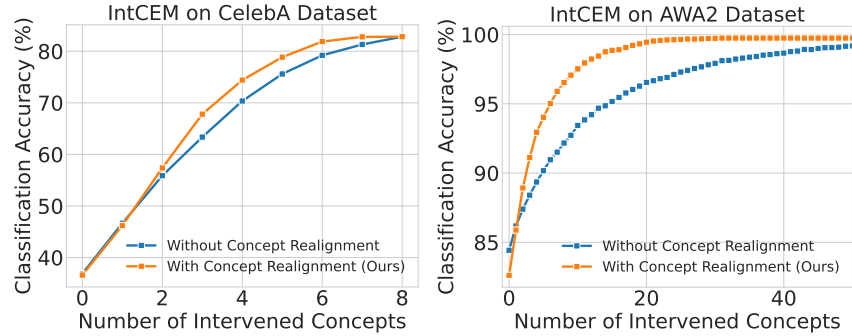
In Fig. 10 we show that the UCP policy is superior across all datasets, and is therefore our default policy across all experiments in this study.



**Fig. 10:** Comparison between accuracy under UCP and Random intervention policies. UCP is superior in all three datasets.

## C Additional Results on IntCEMs

In this section, we report the performance of posthoc concept realignment on the intervention-aware CEMs (IntCEMs) on the CelebA and AWA2 datasets to supplement the results in Section 4.3. In Fig. 11 we show that concept realignment improves the performance of the SoTA approach in both datasets.



**Fig. 11:** Classification accuracy with and without posthoc concept realignment in intervention-aware CEMs. In both cases, concept realignment improves performance of the base IntCEM model.

## D Additional Results

Here, we report the area under the curve of concept prediction loss and classification accuracy using three different random seeds. It can be seen in Tables 2 and 3 that concept realignment consistently improves performance on both metrics.

**Table 2:** Area Under Curve (AUC) of Concept Prediction Loss and Classification Accuracy with/without CIRM for three random seeds on the CUB dataset. We use the same backbone for sequential and independent CBMs. CIRM improves performance across all models and runs.

Base Model	Realigned	Concept Loss AUC ↓			Accuracy AUC ↑		
		Run 1	Run 2	Run 3	Run 1	Run 2	Run 3
Sequential CBM	×	6.72	7.11	6.77	2460.86	2394.1	2444.08
	✓	<b>3.16</b>	<b>3.16</b>	<b>3.24</b>	<b>2510.48</b>	<b>2460.41</b>	<b>2501.08</b>
Independent CBM	×	6.72	7.11	6.77	2653.37	2652.75	2652.47
	✓	<b>3.16</b>	<b>3.16</b>	<b>3.24</b>	<b>2678.09</b>	<b>2675.04</b>	<b>2675.48</b>
Joint CBM	×	5.93	5.84	5.89	2580.28	2533.56	2591.32
	✓	<b>3.67</b>	<b>3.49</b>	<b>3.58</b>	<b>2608.89</b>	<b>2559.93</b>	<b>2622.53</b>
CEM	×	5.99	13.19	6.50	2521.31	1681.97	2579.84
	✓	<b>3.21</b>	<b>6.66</b>	<b>3.43</b>	<b>2558.07</b>	<b>1762.58</b>	<b>2617.25</b>

**Table 3:** Area Under Curve (AUC) of Concept Prediction Loss and Classification Accuracy with/without CIRM for three random seeds on the CelebA dataset. We use the same backbone for sequential and independent CBMs. CIRM improves performance across all models and runs.

Base Model	Realigned	Concept Loss AUC ↓			Accuracy AUC ↑		
		Run 1	Run 2	Run 3	Run 1	Run 2	Run 3
Sequential CBM	×	1.59	1.64	1.65	281.09	279.64	279.47
	✓	<b>1.51</b>	<b>1.53</b>	<b>1.55</b>	<b>284.76</b>	<b>284.00</b>	<b>284.21</b>
Independent CBM	×	1.59	1.64	1.65	280.86	308.38	310.57
	✓	<b>1.51</b>	<b>1.53</b>	<b>1.55</b>	<b>282.48</b>	<b>312.72</b>	<b>316.45</b>
Joint CBM	×	2.88	3.23	3.10	273.06	236.22	296.80
	✓	<b>1.75</b>	<b>1.77</b>	<b>1.74</b>	<b>273.76</b>	<b>246.09</b>	<b>303.76</b>
CEM	×	1.65	1.90	1.83	396.70	366.87	361.60
	✓	<b>1.49</b>	<b>1.66</b>	<b>1.58</b>	<b>401.84</b>	<b>370.88</b>	<b>363.57</b>