

## #TLDR

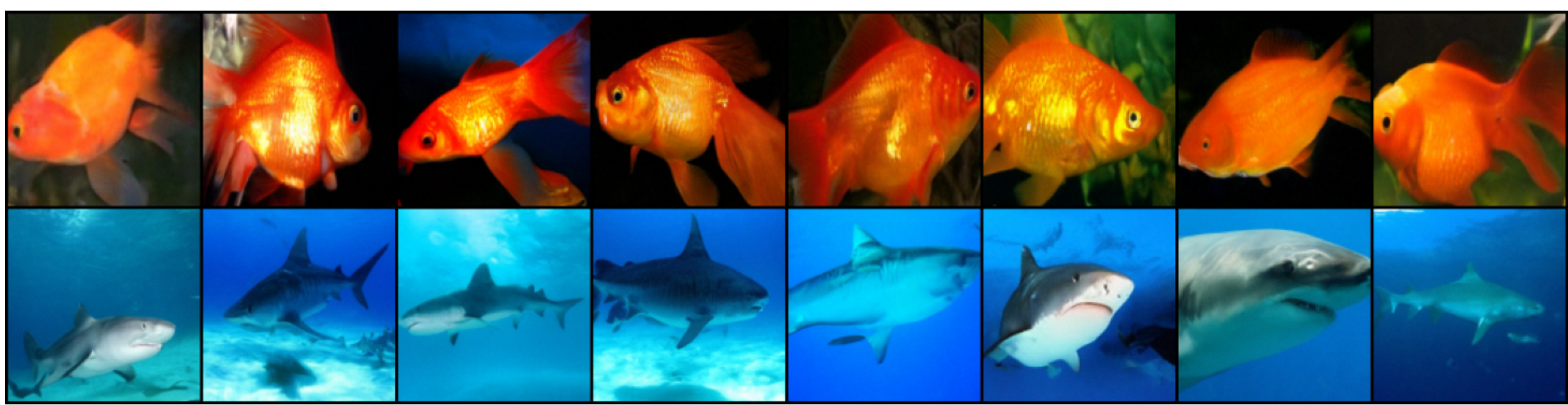
- We introduce **GUIDE** - a generative replay method that benefits from classifier guidance to generate **rehearsal data samples prone to be forgotten**.
- Our sampling approach moves replay examples closer to the classifier's decision boundary, making them highly valuable for class-incremental learning.

## Motivation

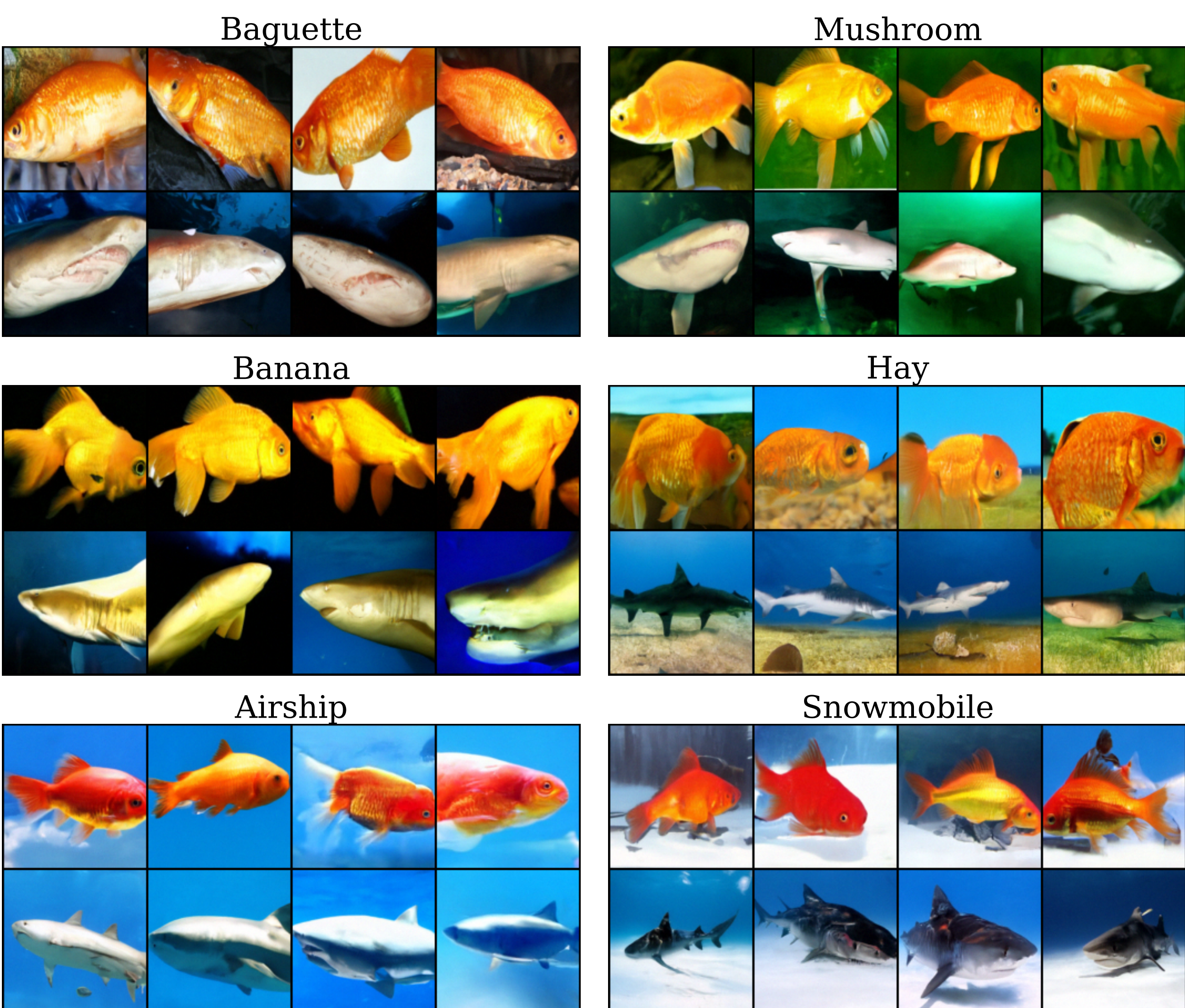
Existing generative replay methods randomly sample rehearsal examples from a generative model, which contrasts with buffer-based approaches where the **sampling strategy plays an important role**.

## Guidance towards unknown classes

- The combination of guidance signals produces samples that align with the training data distribution of the diffusion model yet exhibit features from classes unknown to the model.
- In **GUIDE**, we use this observation in the continual training of a classifier.



(a) Samples generated **without guidance** towards any unknown class.

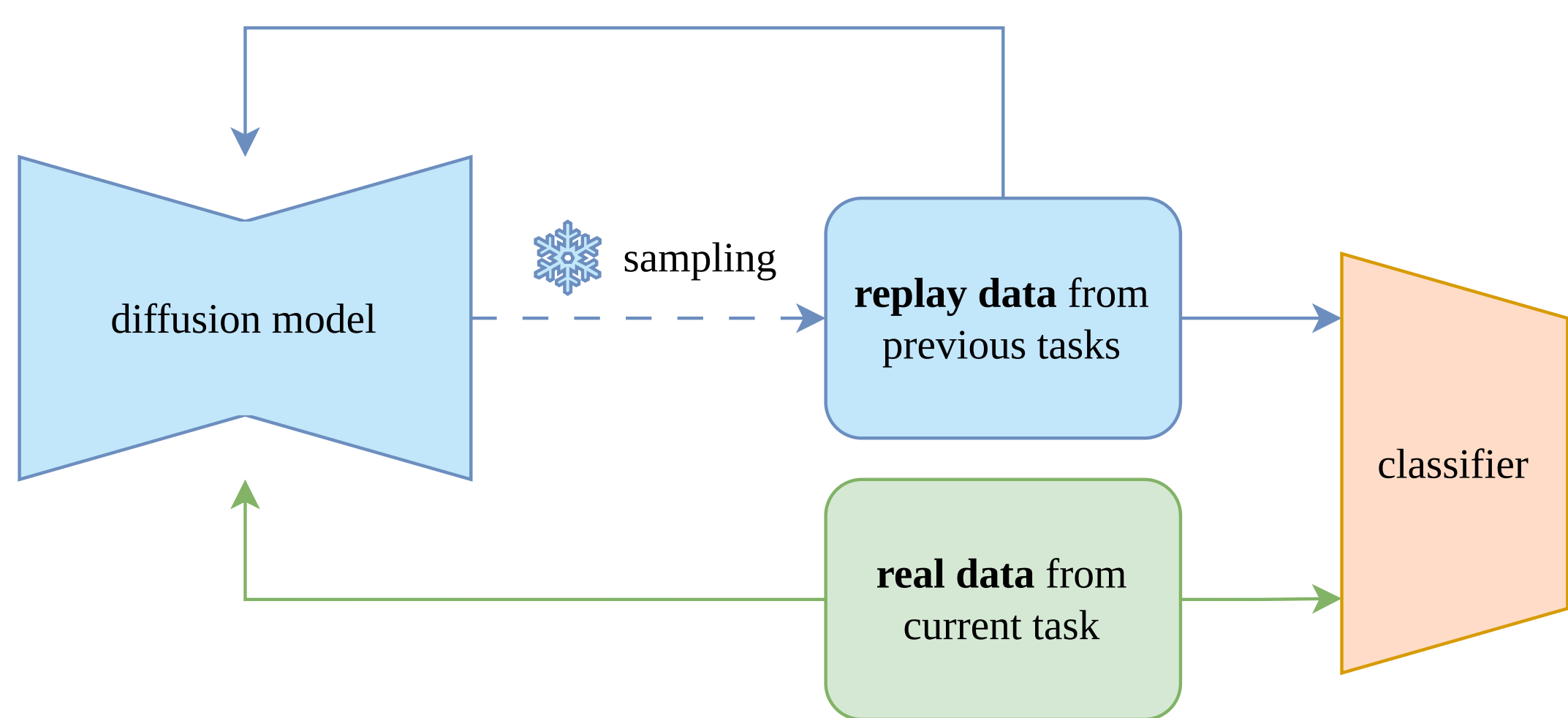


(b) Samples generated **with guidance** towards class unknown to diffusion model depicted above each plot.

Figure 1. Samples from our unconditional diffusion model trained only on *goldfish* and *tiger shark* classes.

## Background – generative replay

In each task  $i$ , we train a classifier  $f_{\phi_i}(y|x)$  and a class-conditional diffusion model  $\epsilon_{\theta_i}(x_t, t, y)$  on currently available data along with **synthetic data samples** from preceding tasks generated by the previous frozen diffusion model  $\epsilon_{\theta_{i-1}}(x_t, t, y)$ .



## Results – average accuracy after the final task

Method	CIFAR-10		CIFAR-100		ImageNet100
	$T=2$	$T=5$	$T=5$	$T=10$	$T=5$
Joint	93.14 ± 0.16		72.32 ± 0.24		66.85 ± 2.25
Continual Joint	85.63 ± 0.39	86.41 ± 0.32	73.07 ± 0.01	64.15 ± 0.98	50.59 ± 0.35
Fine-tuning	47.22 ± 0.06	18.95 ± 0.20	16.92 ± 0.03	9.12 ± 0.04	13.49 ± 0.18
DGR VAE	60.24 ± 1.53	28.23 ± 3.84	19.66 ± 0.27	10.04 ± 0.17	9.54 ± 0.26
DGR+distill	52.40 ± 2.58	27.83 ± 1.20	21.38 ± 0.61	13.94 ± 0.13	11.77 ± 0.47
RTF	51.80 ± 2.56	30.36 ± 1.40	17.45 ± 0.28	12.80 ± 0.78	8.03 ± 0.05
MeRGAN	50.54 ± 0.08	51.65 ± 0.40	9.65 ± 0.14	12.34 ± 0.15	-
BIR	53.97 ± 0.97	36.41 ± 0.82	21.75 ± 0.08	15.26 ± 0.49	8.63 ± 0.19
GFR	64.13 ± 0.88	26.70 ± 1.90	34.80 ± 0.26	21.90 ± 0.14	32.95 ± 0.35
DDGR	80.03 ± 0.65	43.69 ± 2.60	28.11 ± 2.58	15.99 ± 1.08	25.59 ± 2.29
DGR diffusion	77.43 ± 0.60	59.00 ± 0.57	28.25 ± 0.22	15.90 ± 1.01	23.92 ± 0.92
<b>GUIDE</b>	<b>81.29 ± 0.75</b>	<b>64.47 ± 0.45</b>	<b>41.66 ± 0.40</b>	<b>26.13 ± 0.29</b>	<b>39.07 ± 1.37</b>

## Method

Modified diffusion denoising process during generation of a sample from class  $y_{i-1}$ :

$$\hat{\epsilon}_{\theta_{i-1}}(x_t, t, y_{i-1}) = \epsilon_{\theta_{i-1}}(x_t, t, y_{i-1}) + \underbrace{s \nabla_{x_t} \ell(f_{\phi_i}(y|\hat{z}_0(x_t), y_i))}_{\text{guidance towards class from current task } i}$$

- Generated rehearsal samples are positioned near the classifier's decision boundary and are more likely to be forgotten.
- Since we use the previous frozen diffusion model, which is **not trained on the classes from the current task**, we consistently obtain samples from the desired class  $y_{i-1}$ .

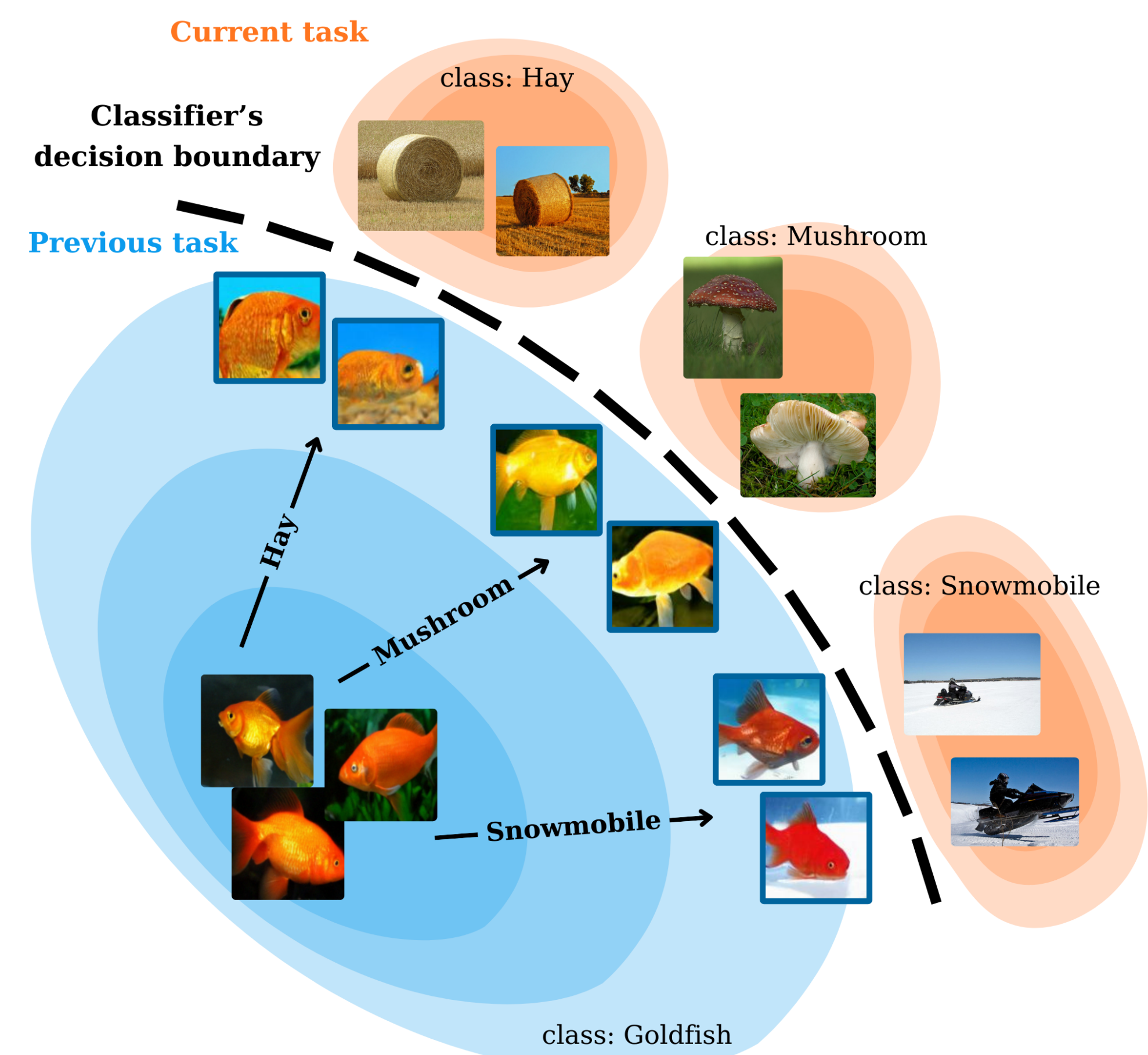


Figure 2. Rehearsal sampling in **GUIDE**. The replay samples, highlighted with **blue borders**, share features with the examples from the current task, which may be related to characteristics such as color or background.

## GUIDE mitigates forgetting in a classifier

Task number	DGR Diffusion					GFR					GUIDE				
	1	2	3	4	5	1	2	3	4	5	1	2	3	4	5
1	85.1	41.72	12.37	9.2	5.78	87.38	61.35	56.83	55.12	52.66	85.1	50.55	24.2	14.3	10.55
2		86.43	41.42	24.47	11.27		59.33	45.18	38.85	34.55		82.45	68.82	47.38	24.07
3			82.43	42.4	18.18			39.35	31.5	27.43			75.82	62.05	43.25
4				78.23	22.3				34.48	29.23				68.23	56.52
5					83.7					30.25					73.93

## Where are the rehearsal samples located?

We launch an adversarial attack on the generated rehearsal samples to check how easy it is to change their class to the one from the current task.

	Misclassified examples	Prediction confidence	
		Previous classifier	Current classifier
DGR diffusion	55.13%	99.60%	90.03%
<b>GUIDE</b>	72.66%	86.42%	61.61%

Properties of rehearsal samples in **GUIDE**:

- It is easy to fool a classifier with them
- Yield lower outputs in the currently trained classifier than in the classifier from the previous task

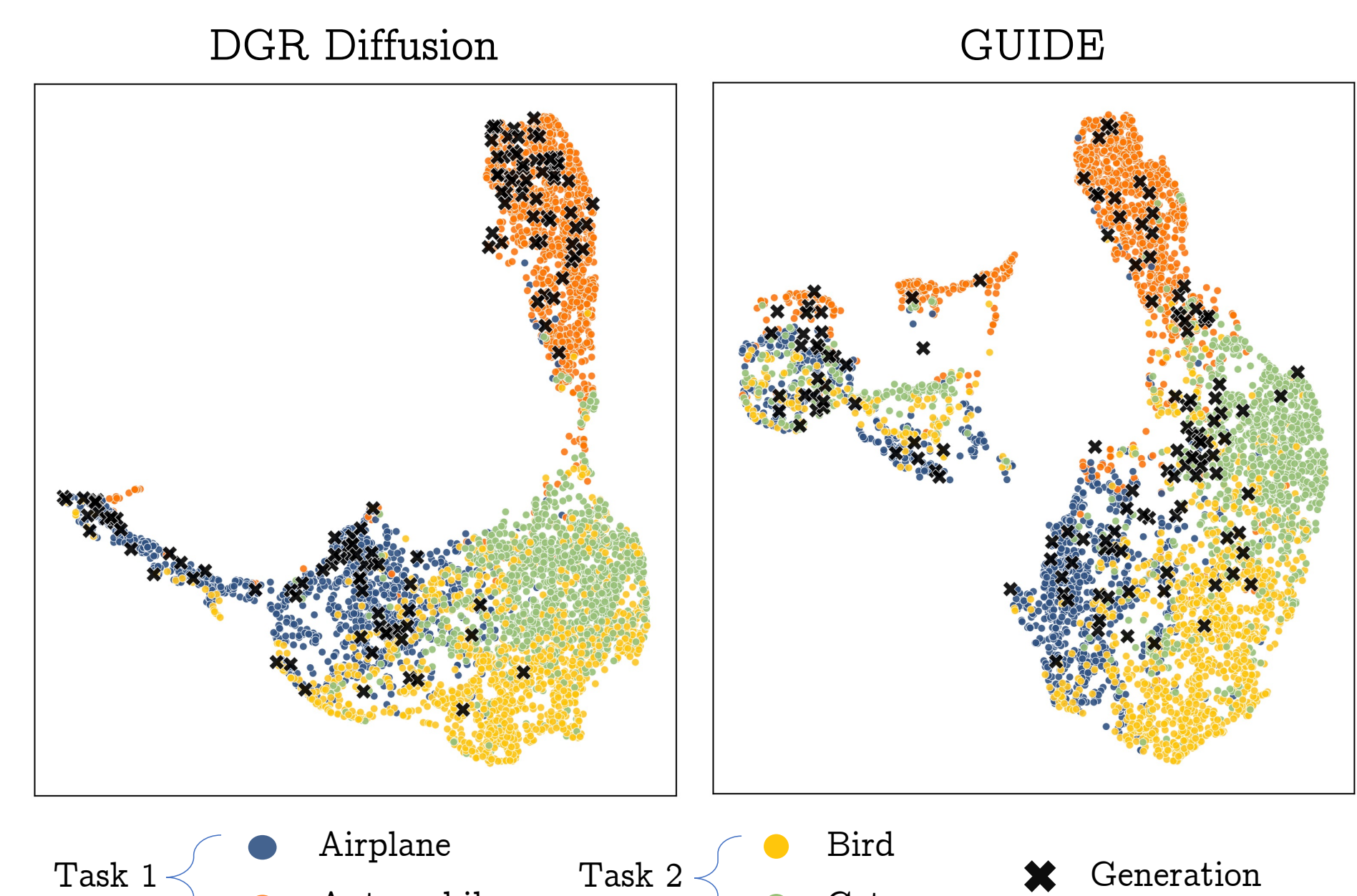


Figure 3. Rehearsal samples in the standard generative replay scenario predominantly originate from high-density regions of class manifolds, while **GUIDE** yields generations that are more similar to the examples from the new task.

Let our paper **GUIDE** you!

