



SAEPer: Sparse AutoEncoders Projection for Reliable Personalization in Diffusion Models



Marek Skiba^{1*}, Jowita Drozdowicz^{1*}, Vladimir Zaigrajew²,
Przemysław Biecek^{1,2}, Piotr Sankowski¹

¹ Uniwersytet Warszawski, ² Politechnika Warszawska

Personalization in Diffusion Models

One image, many concepts, endless possibilities

Personalization teaches diffusion models to **recreate a subject or attribute** from just one reference photo while keeping prompts fully open-ended. Instead of one monolithic image, we break it into **concepts: identity, material, style, lighting, pose**, and make each reusable in new scenes.

Our **lightweight controller** leaves the base model untouched: it modulates existing words (like "hat," "pose," "lighting") at inference. During training, it binds each visual concept to its corresponding word, letting you **recombine pieces of a reference photo through natural language**.

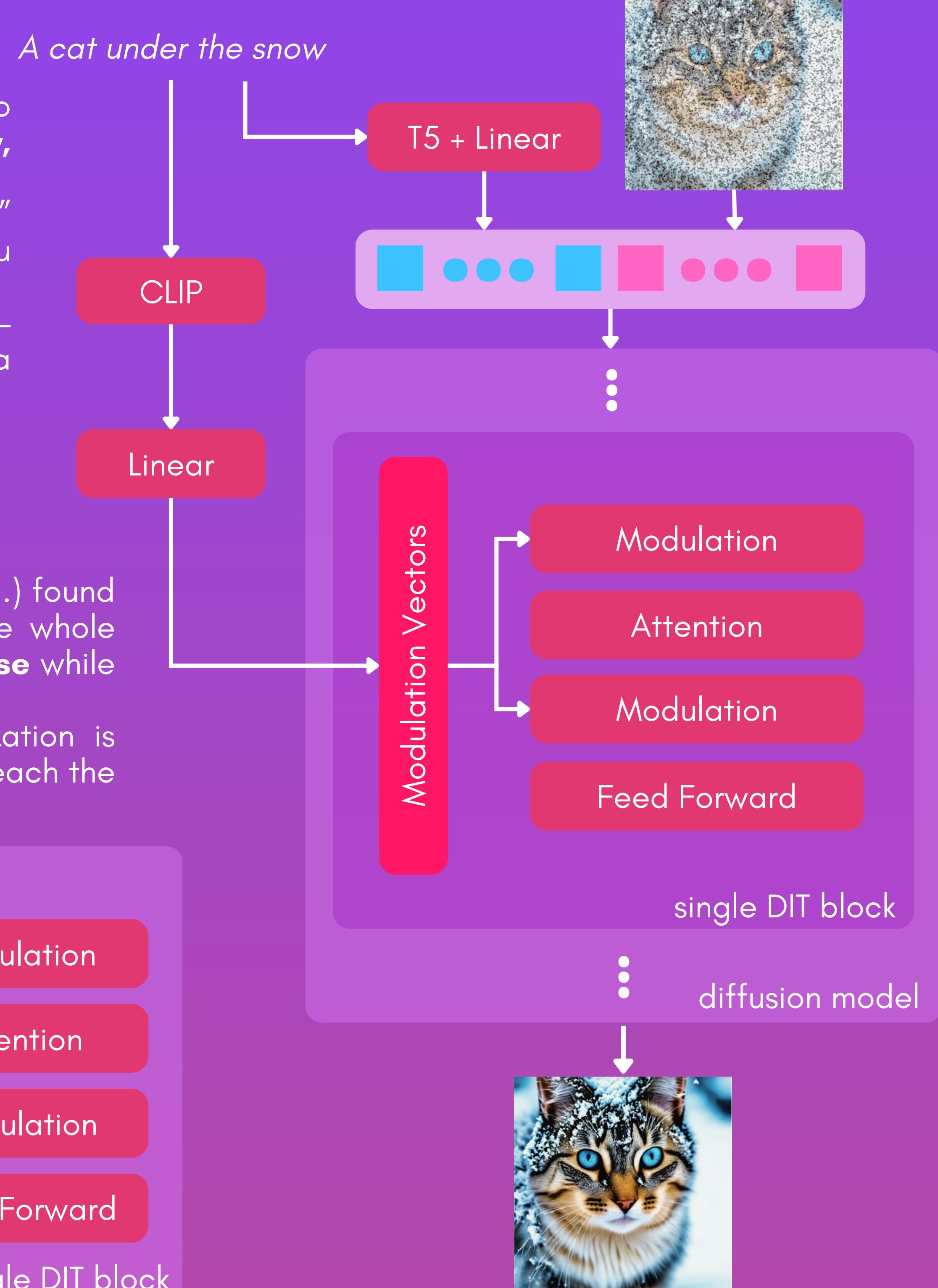
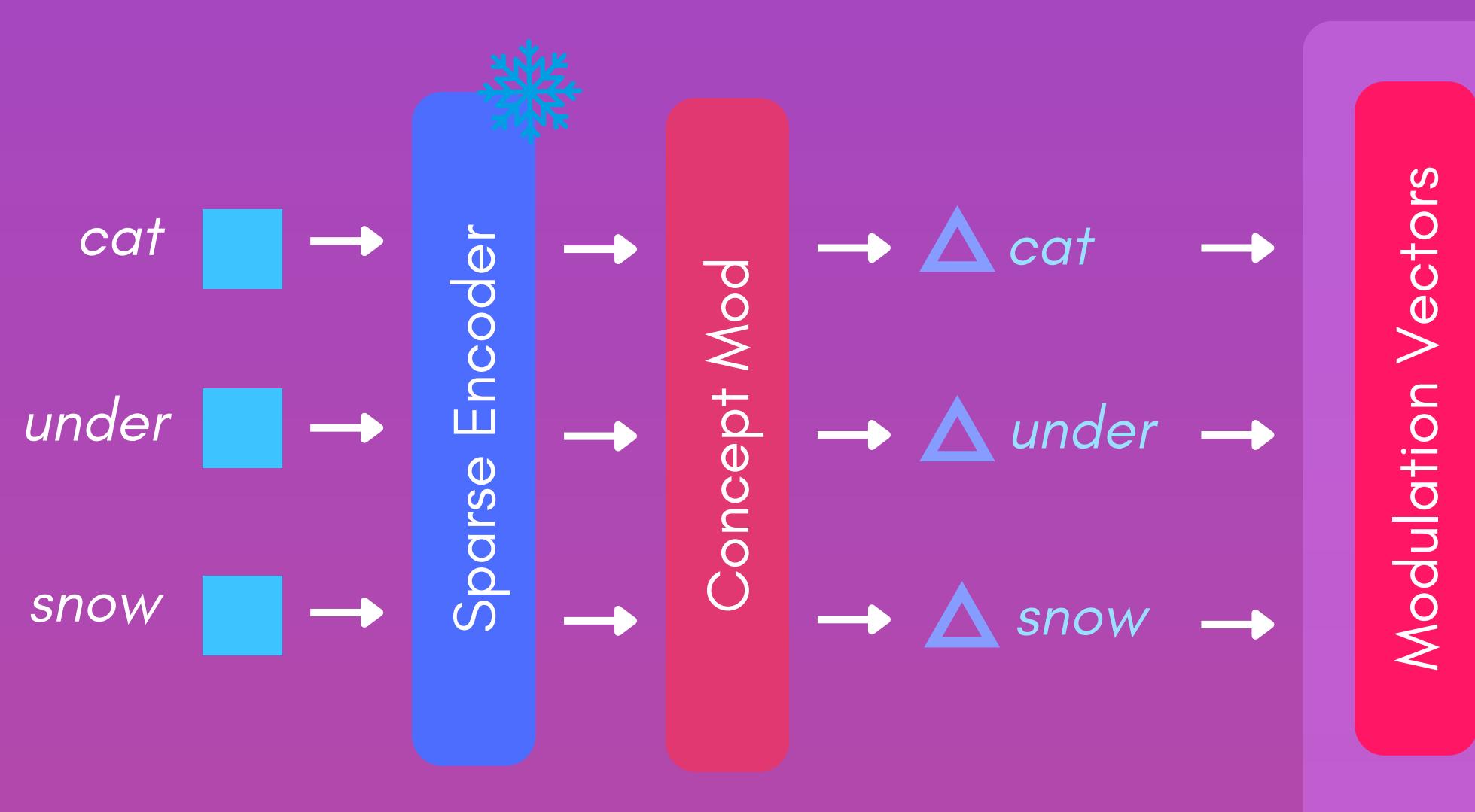
Generate a portrait with the same hat, move the pose to a beach, or keep the lighting but switch to watercolor—all without finetuning or new embeddings. It's **fast, compositional, and robust**, enabling faithful reuse from a **single image**.

TokenVerse – Inspiration

One token, precise control, full freedom

Modern diffusion transformers guide images through **cross-attention** and **modulation**. TokenVerse (Garibi et al.) found that steering the modulation path per word token yields **precise, localized control**. Instead of altering the whole embedding, a **ConceptMod** module learns **per-token shifts** that adjust attributes like **identity, material, or pose** while keeping the rest untouched.

At inference, these **token-level steering directions** are plug-and-play—apply them only where personalization is desired. This enables **fine-grained edits, flexible composition**, and **zero base-model tuning**, a clean way to teach the model what to change and what to leave alone.



Reference image



A girl wearing a **hat**, a **coat** and a **shirt**, in the **city**



A **hat** resting on a **chair**

Generated images - SAEPer (ours)



A **girl** playing with a **puppy**



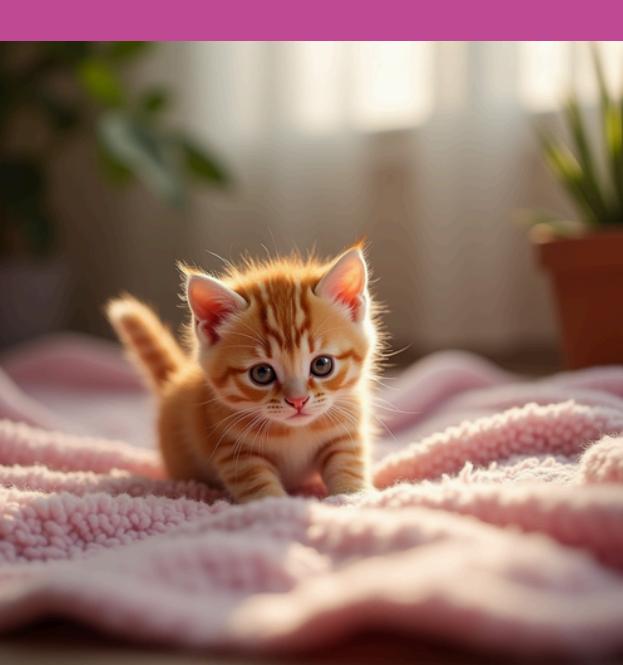
A **coat** hanging by a **door**



A **folded shirt** on the **table**



A **backpack** decorated with **stars** with a **kitten** and a **puppy** on its front, **northern lights**



A **backpack** decorated with **stars** on a plain **background**



A **kitten** on a **blanket**



A **puppy** running through **water**



A **beach** with **northern lights** in **background**

SAEPer – Our Method

Sparse codes, clean separation, stronger control

While TokenVerse offers control, **concept leakage** often remains: the "dog" still drags along its "hat." We address this by introducing **SAEPer**, combining TokenVerse with **Sparse Autoencoders (SAEs)** to **explicitly disentangle factors**.

We first encode tokens into **sparse, semantically independent features**, then let ConceptMod learn on this cleaner space. A **targeted pruning step** keeps only the most relevant neurons, ensuring compactness and focus. This sparse representation **removes the need for isolation loss**, naturally enforcing **per-token separation** and minimizing leakage: **sharper control, cleaner composition, and higher visual fidelity** without sacrificing flexibility.

Experiments & Results

Sharper separation, faithful detail, robust composition

Our experiments show that **SAEPer** delivers a **significant boost in disentanglement and concept isolation**. Even **fine patterns** like stars on a backpack, jewelry details, or small fabric prints are correctly separated and can be **reused independently** in new prompts. Unlike previous methods, SAEPer **preserves fine visual fidelity** while maintaining **clean per-token separation**, preventing attributes from leaking across concepts.

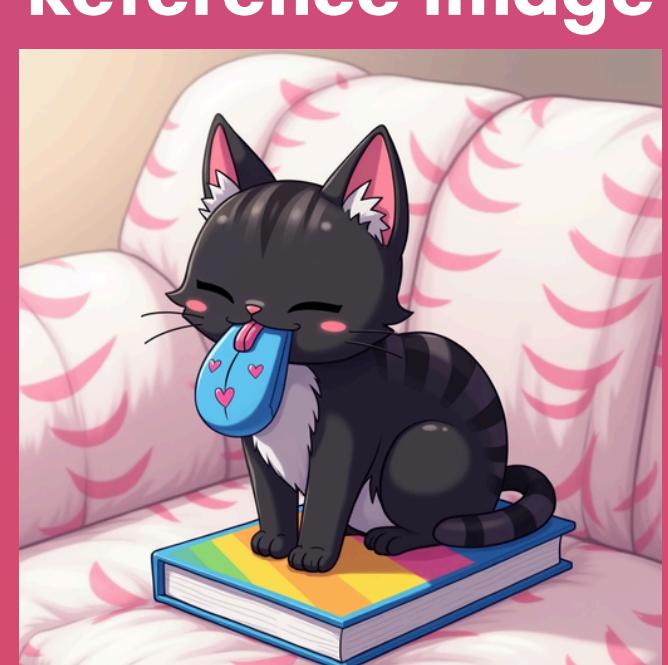
Empirically, we observe **robust multi-image fusion** – concepts from several references can be combined into a single coherent scene, with each element retaining its distinctive look. SAEPer thus achieves **precise control and faithful reuse**, bridging creative flexibility with **consistent personalization**.

Infinite Composability

Independent concepts, infinite combinations composability

With clean separation of identity, clothing, and environment, SAEPer composes new scenes by **mixing elements across examples**: a hat, a pose, a pattern, a place, all staying faithful to their sources.

Reference image



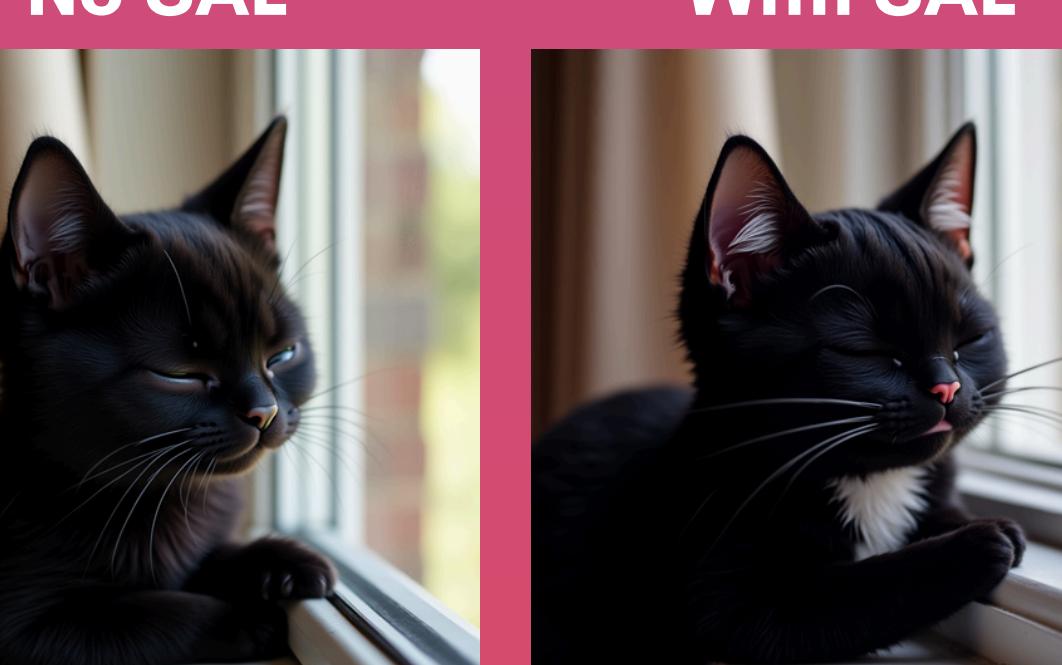
A **cat** biting a **computer mouse**, sitting on the couch with **stripes**



A **computer mouse** placed beside a **laptop** in an **office**



No SAE



A **cat** sleeping beside a **window**



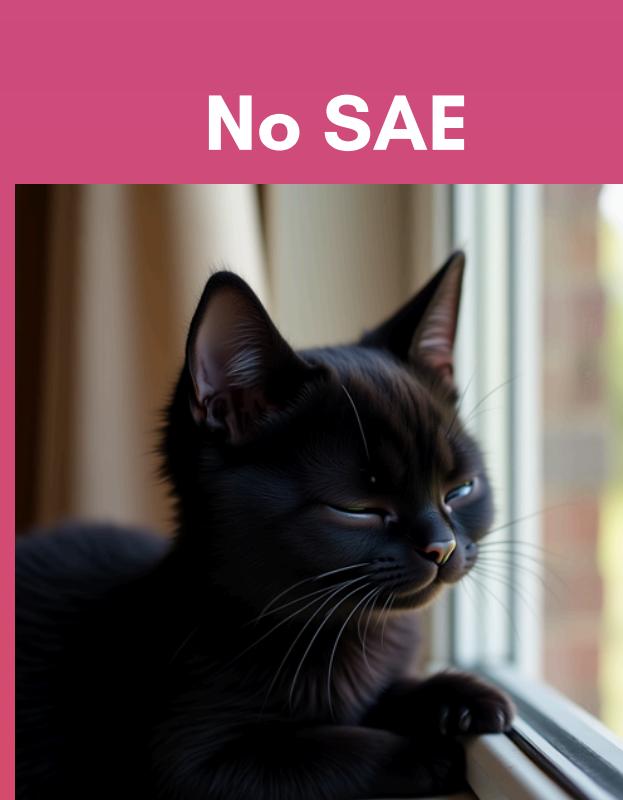
A **cat** biting a **computer mouse**, sitting on the couch with **stripes**



A **man** wearing a **hat**, a **necklace** and a **shirt** sitting in the **armchair**



A **hat** resting on a **concrete block** under **soft light**



A **man** playing **guitar** on a **stage** with **colorful light**



With SAE



A **puppy** wearing a **hat**, eating a **computer mouse**, sitting on street of the **city**



A **puppy** wearing a **hat** and a **coat**, sitting in the **armchair** with **stripes**



A **girl** wearing a **hat** decorated with **stars** and a **shirt**, in the **city**