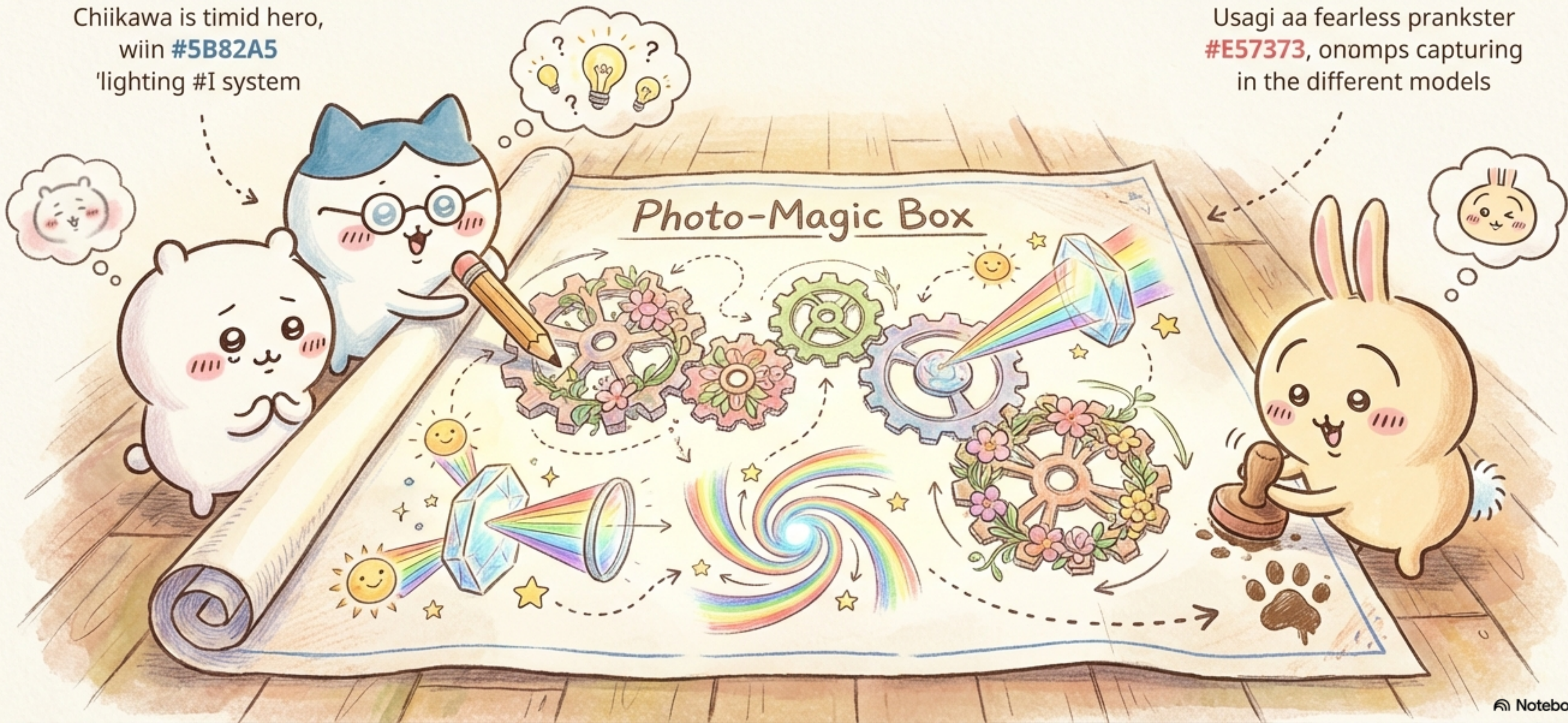


The Quest for the Photo-Magic Box

A Story of Learning to See in 3D without Instructions

Chiikawa is timid hero,
wiin **#5B82A5**
'lighting #I system

Usagi aa fearless prankster
#E57373, onomps capturing
in the different models



What if we could explore a memory from any angle?

Imagine a machine that takes just a handful of 2D photos from a place... and magically creates a new picture from any viewpoint you choose! This is the challenge of Novel View Synthesis ([#5B82A5](#)) (or NVS for short). It's like turning a few flat pit into a living, explorable 3D scene.

A Few Sparse Photos



A Brand New View!



The Old Way: Building with an Instruction Manual

Early attempts to build the Photo-Magic Box relied on lots of “Helper Tools”—extra information that told the machine exactly how the world should look. This is like giving it a strict instruction manual. This form of **camera poses**, sometimes (NeRF and 3DGS).

Known Camera Poses

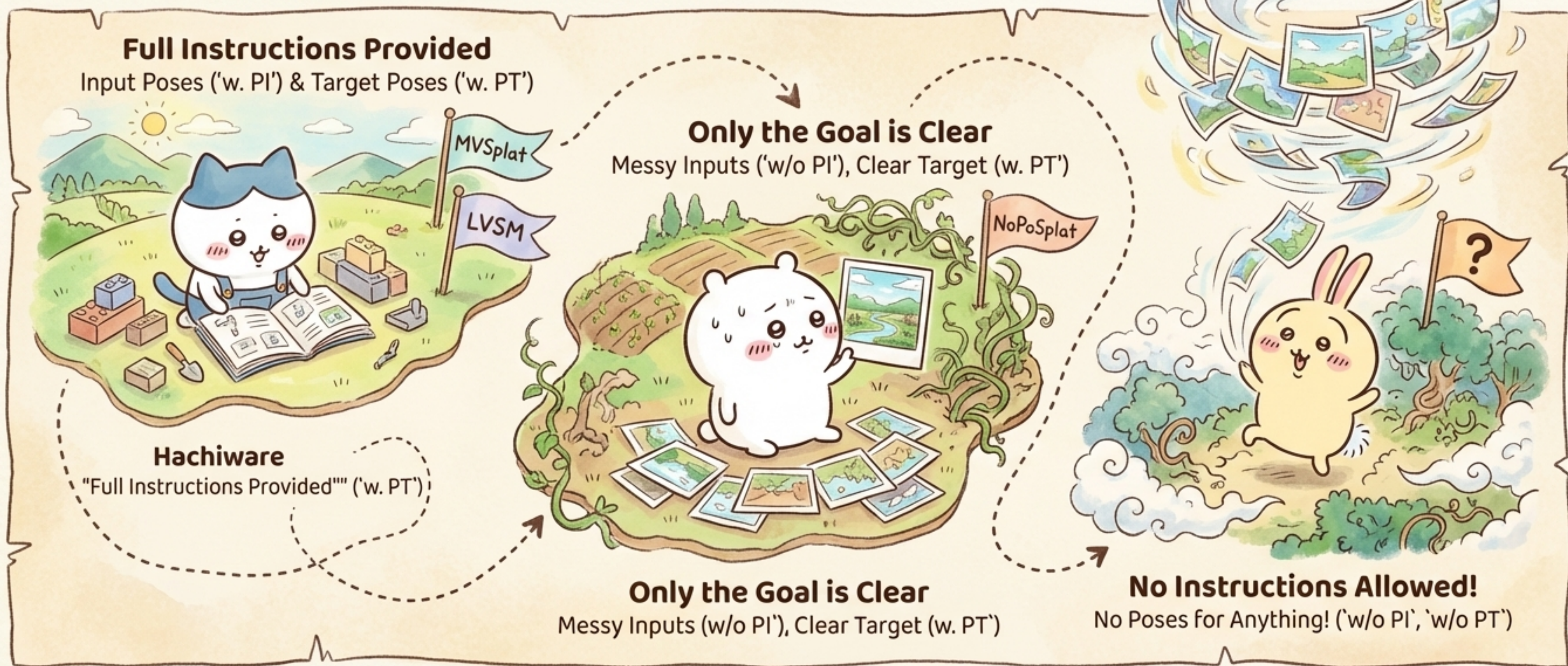


Built-in 3D Rules



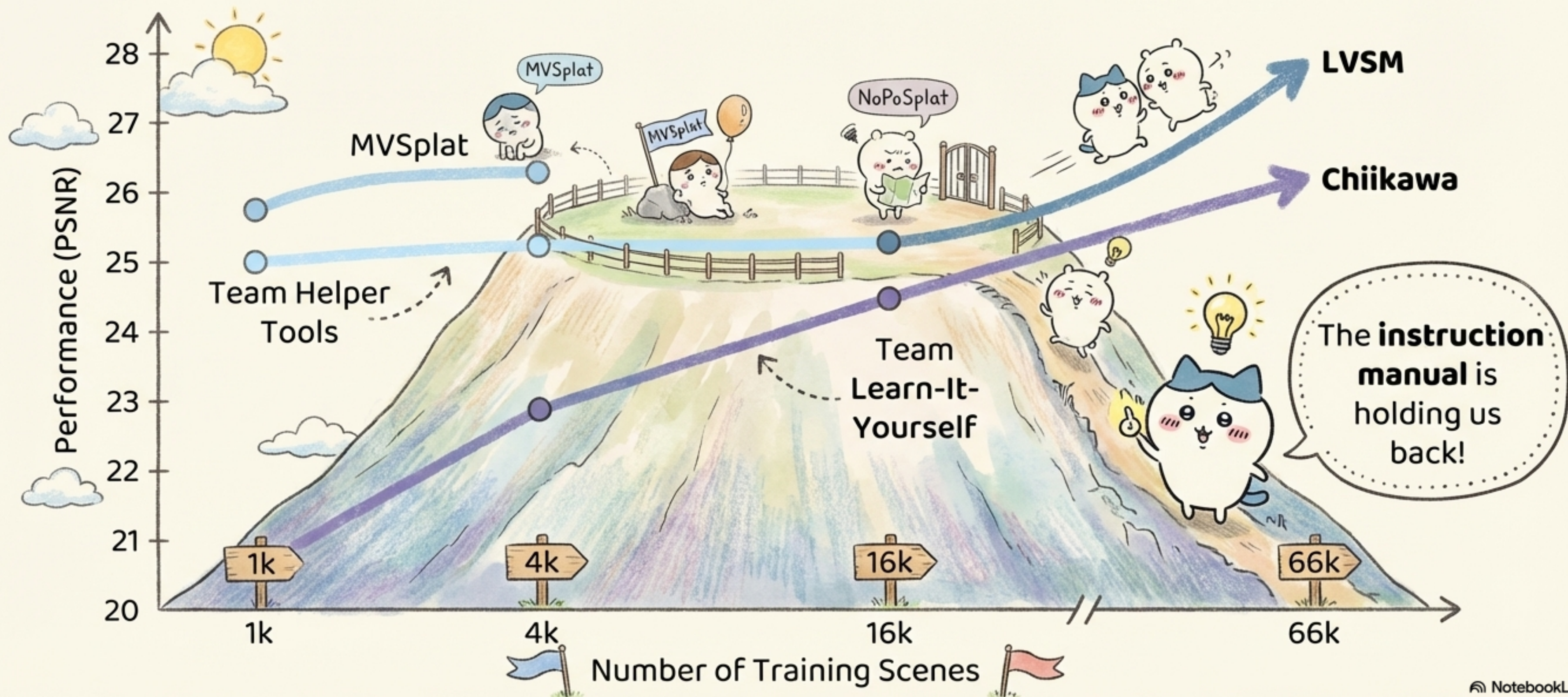
A Map of Different Building Strategies

Depending on how many 'Helper Tools' you use, you're in a different part of the NVS world.



An Unexpected Discovery in the Great Data Race!

Our heroes noticed something strange. When they gave their machines more and more photos to learn from, the teams using **fewer "Helper Tools"** started getting much, much better—and faster!



The Secret Principle: The Less You Depend, The More You Learn

The 'instruction manual' (**3D inductive bias**) is a crutch. It helps when you have very little data, but it prevents the machine from discovering deeper, more flexible rules on its own.

By removing the strict rules, the machine is forced to learn true **3D awareness** directly from the visual data. The performance of these **'data-centric'** methods accelerates more as data scales.



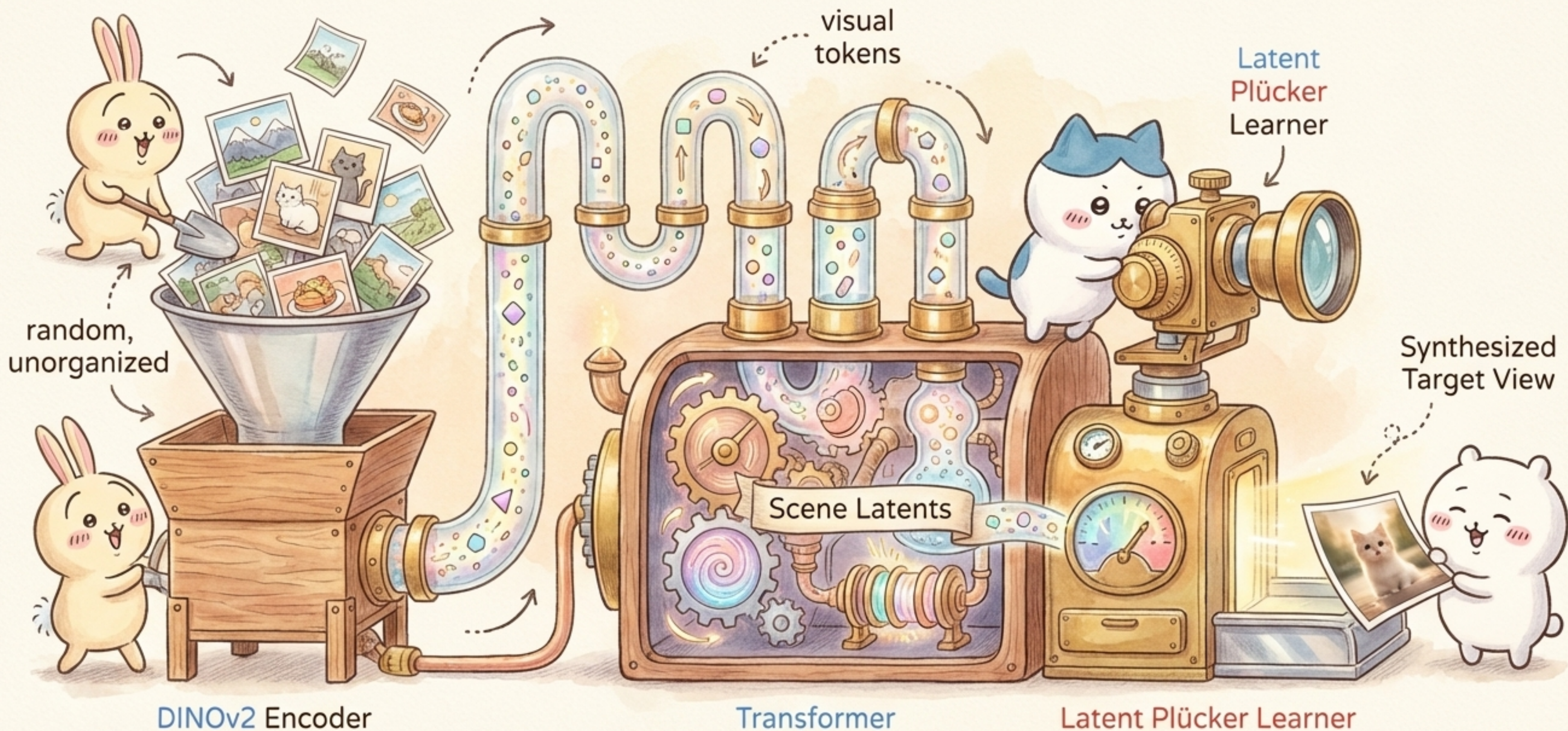
Our Ultimate Quest: Build a Box for the Unposed World

Motivated by their discovery, our heroes set out to build a new Photo-Magic Box for the most challenging scenario: the **unposed setting** (highlighted in blue). This machine would need to work with completely random, disorganized photos, without any camera information for the inputs or the final target view.

It would have to learn everything from scratch.

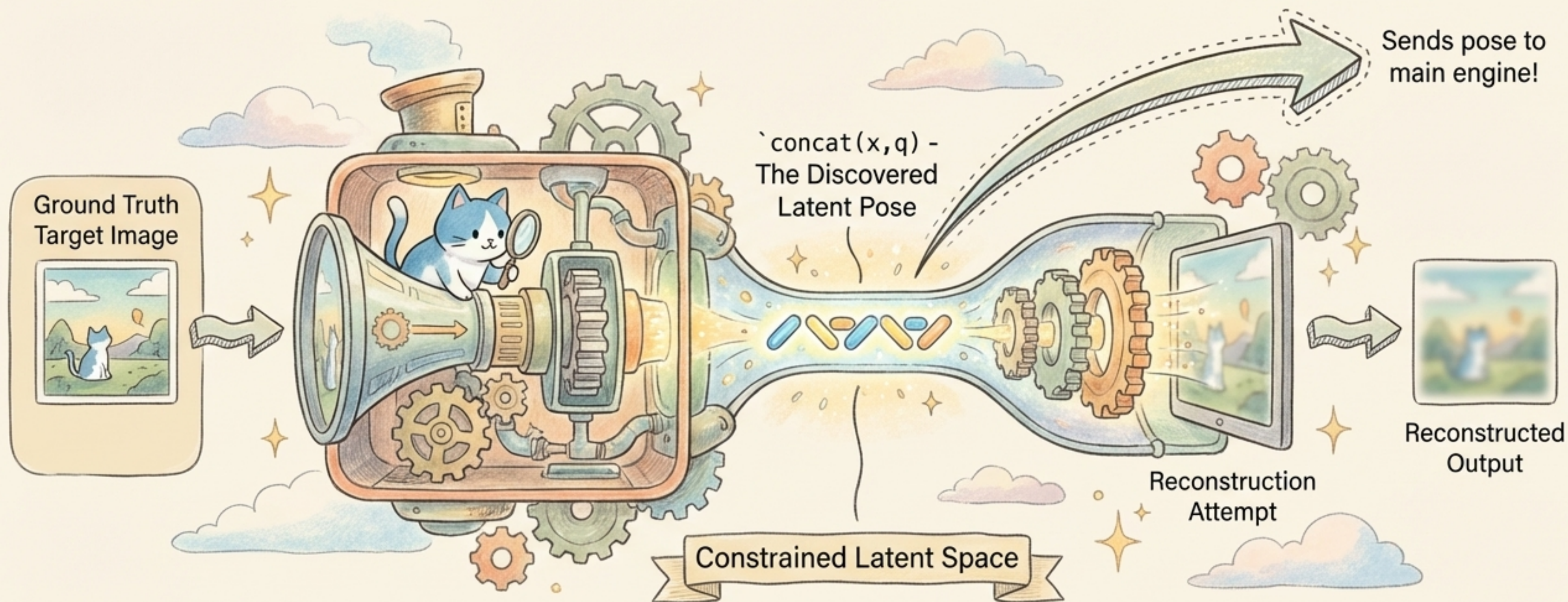


Introducing: The UP-LVSM! (Unposed Large View Synthesis Model)



The Secret to the Machine: A Magic Viewfinder

How does the machine know the camera angle for a picture if we don't tell it? It uses a clever trick! The '**Latent Plücker Learner**' looks at the target picture and automatically figures out the correct viewpoint all by itself. It learns to build its own internal 'pose space' without any 3D supervision, like an autoencoder.



The Photo-Magic Box is a Success!

The results are in! Trained without *any* 3D supervision, our UP-LVSM creates **photorealistic** and **3D-consistent** views. It achieves performance comparable to—and in some cases better than—methods that require camera poses.

Performance on RealEstate10K

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
MVSplat	26.45	0.874	0.123
LVSM	27.60	0.874	0.117
NoPoSplat	25.46	0.854	0.137
UP-LVSM (Ours)	28.82	0.891	0.104

Qualitative Comparison

Input Views



MVSplat



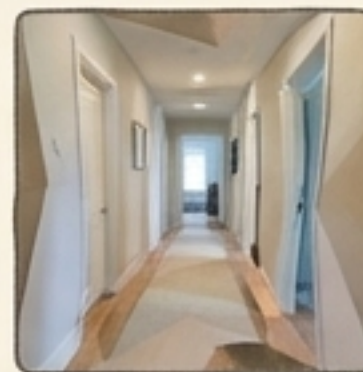
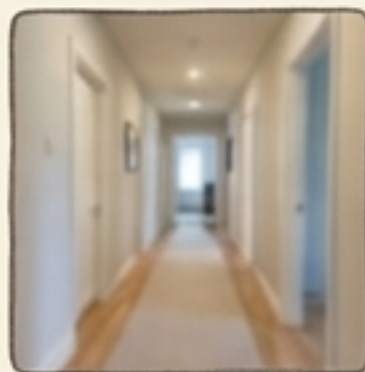
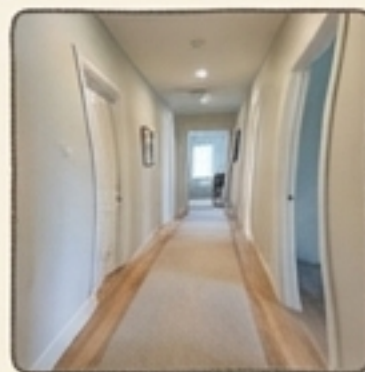
LVSM



NoPoSplat

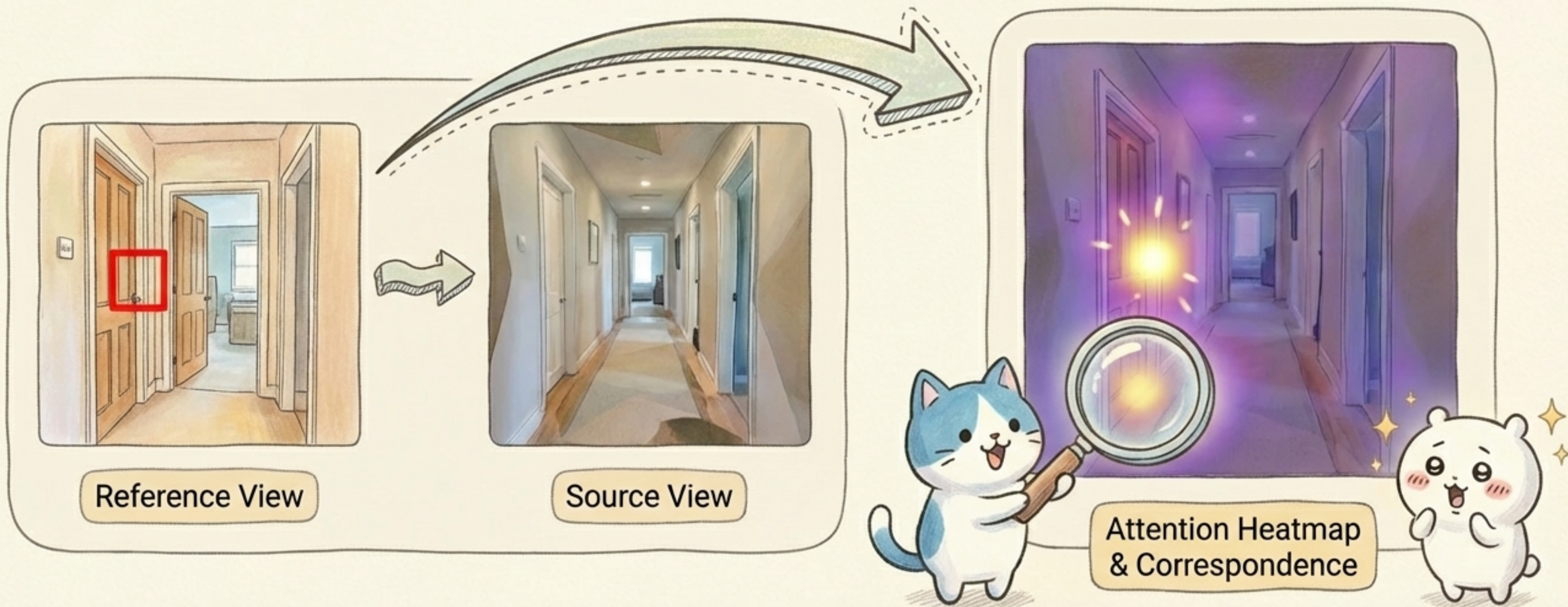


UP-LVSM (Ours)



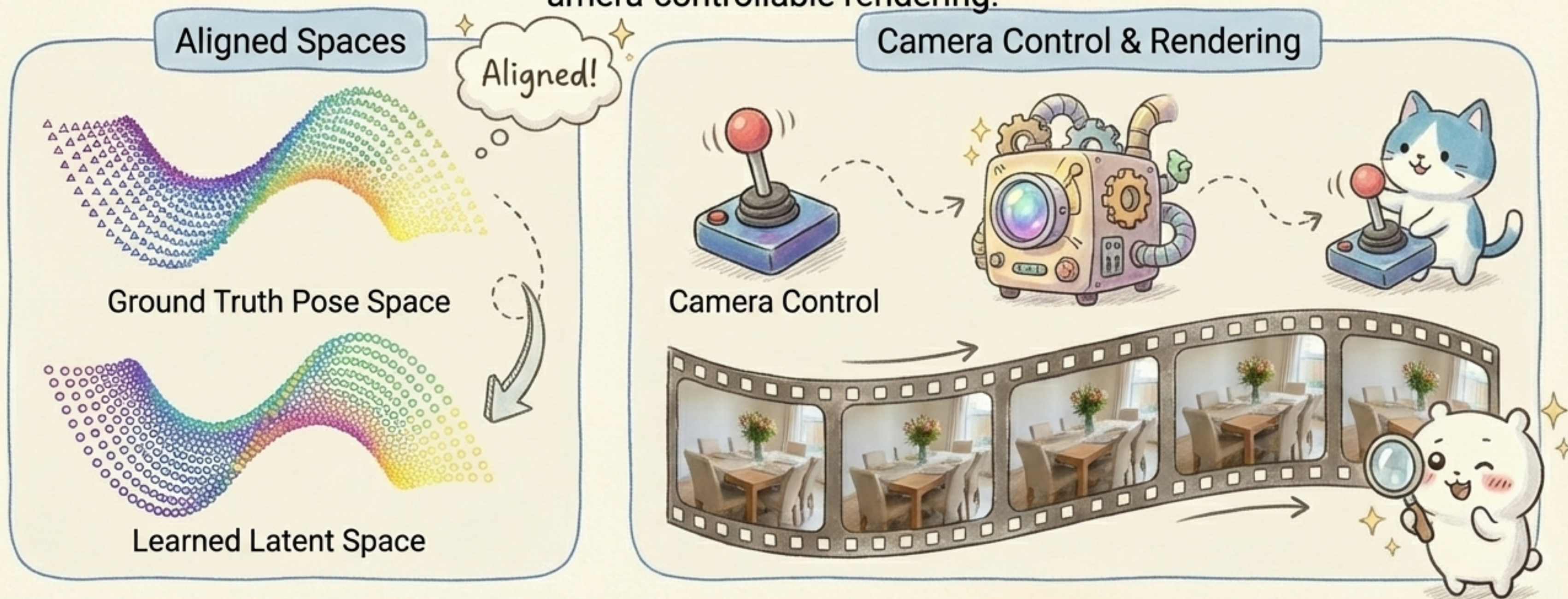
Proof: The Machine is Truly Understanding 3D Space

We can prove the machine is learning by peeking inside its “brain.” When we ask it which parts of two different images correspond to the same real-world object, the attention maps light up in exactly the right places! This shows it’s learning **3D spatial correspondence** from 2D images alone.



We Can Even Steer the Magic!

Because the "Magic Viewfinder" learns a meaningful and organized map of camera poses (latent space - highlighted in blue), we can give it real-world controls! By fine-tuning with a tiny amount of posed data, we can add a simple "linear mapper" that acts like a joystick, allowing for explicit, camera-controllable rendering.



The Adventure Showed Us: True Learning Comes from Freedom

Our heroes' quest revealed a powerful truth. By removing the constraints of pre-defined rules and 3D knowledge (**highlighted in blue**), they built a machine that could learn a deeper, more flexible understanding of our world, just from looking at pictures. This **data-centric (in orange)** paradigm opens a new path for teaching machines to see.



What else can we learn by letting go of our assumptions?