

GenAI for Creatives: A Few Projects from Adobe Research, Lessons Learned and Avenues for Future Work

Sylvain Paris

Adobe

I live in Boston. Thanks a lot for inviting me!



Dall•E 2 Got People's Attention 4 Years Ago

*A photorealistic image
of an astronaut riding a
horse*



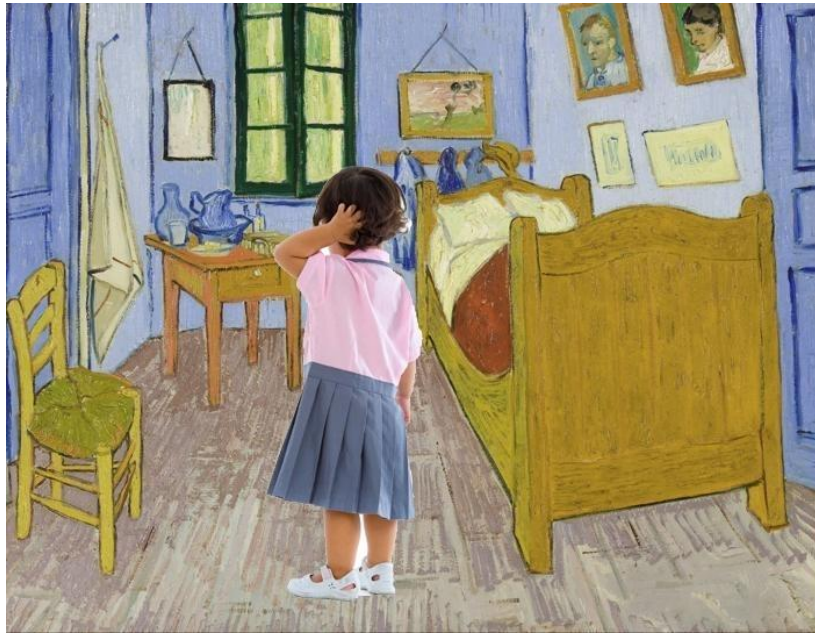
Progress Happens Faster Than Ever



*A photorealistic image
of an astronaut riding a
horse*

[Nano Banana 2]

Generative AI Models Solve Hard Problems Without Even Trying



Naïve compositing



Deep Painterly Harmonization
[Luan et al., EGSR 2018]

6 ~ 12 months of work
by a talented student
and 3 senior researchers



Nano Banana 2:
*Make the girl look like she
belongs to the painting*



Naïve compositing

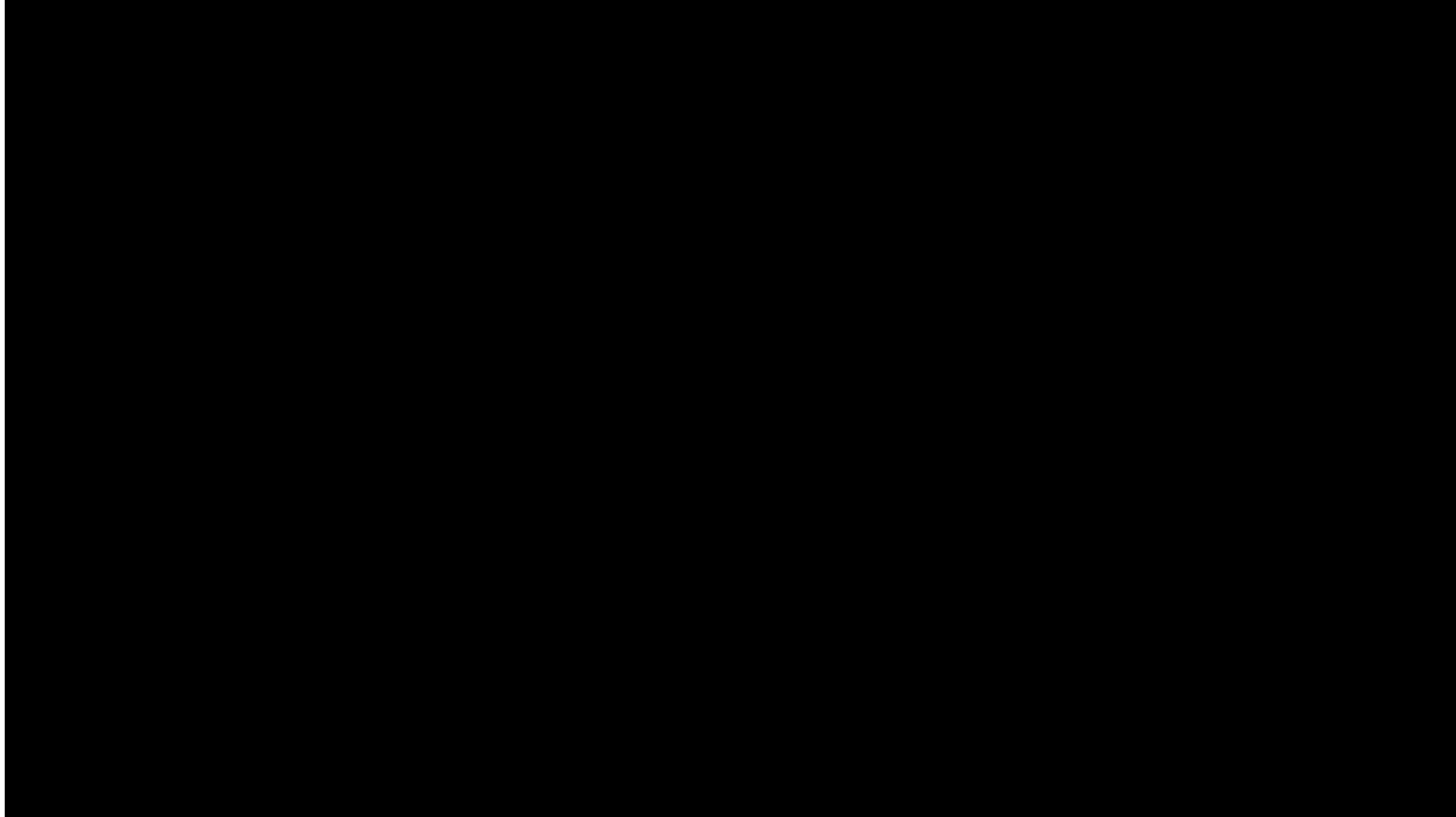


Deep Painterly
Harmonization



Nano Banana 2

Then Sora Did it for Video 2 Years Ago



**A large part of Adobe's business
is about tools for creatives.**

This is a momentous change.

This impacts Adobe Research.

A large part of Adobe's business
is about tools for creatives.

This is a momentous change.

This impacts Adobe Research.

Is there anything left to do?

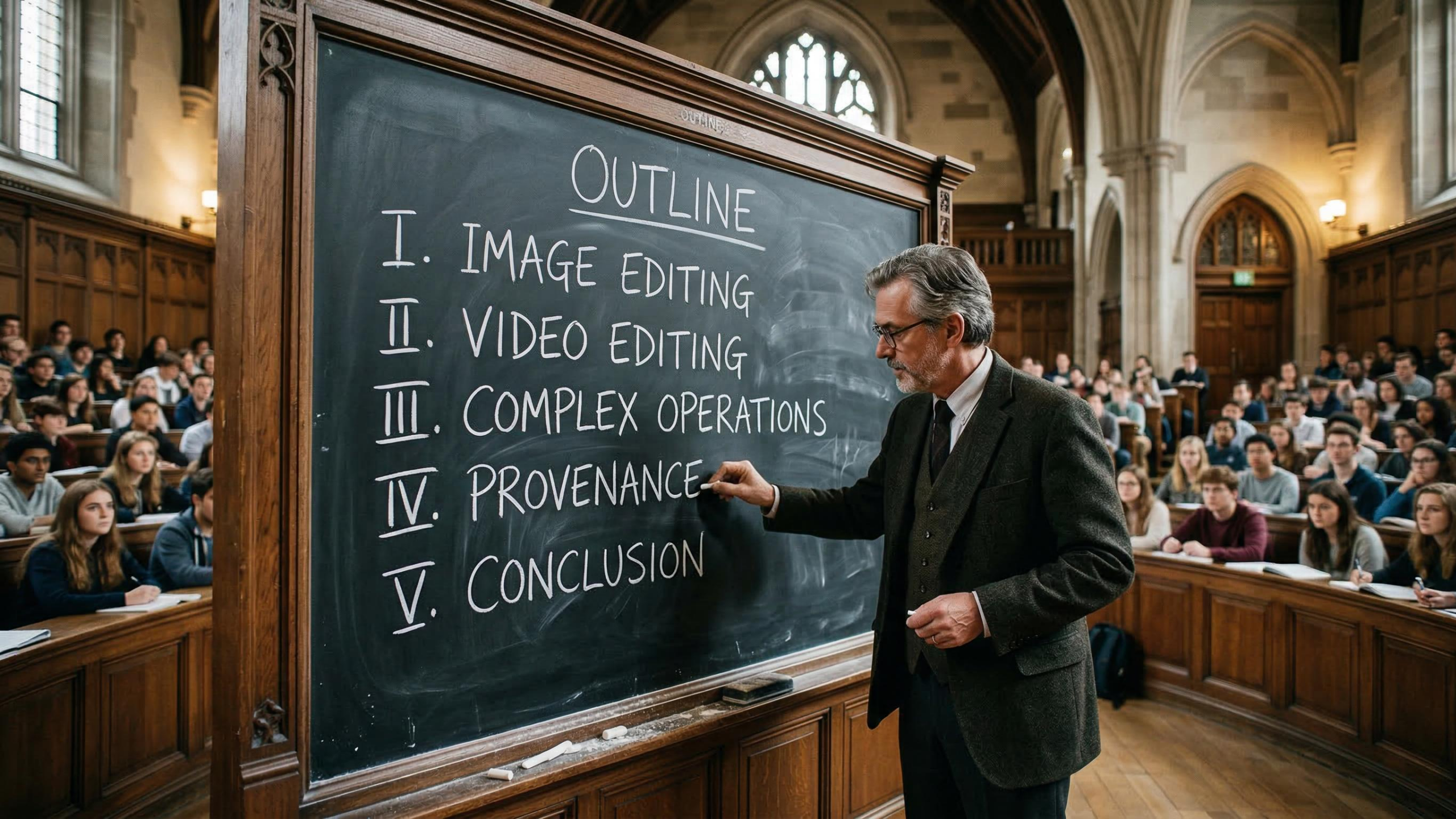
**What if this is not exactly
what I have in mind?**

Control matters.



OUTLINE

- I. IMAGE EDITING
- II. VIDEO EDITING
- III. COMPLEX OPERATIONS
- IV. PROVENANCE
- V. CONCLUSION



UniReal: Advanced Image Editing

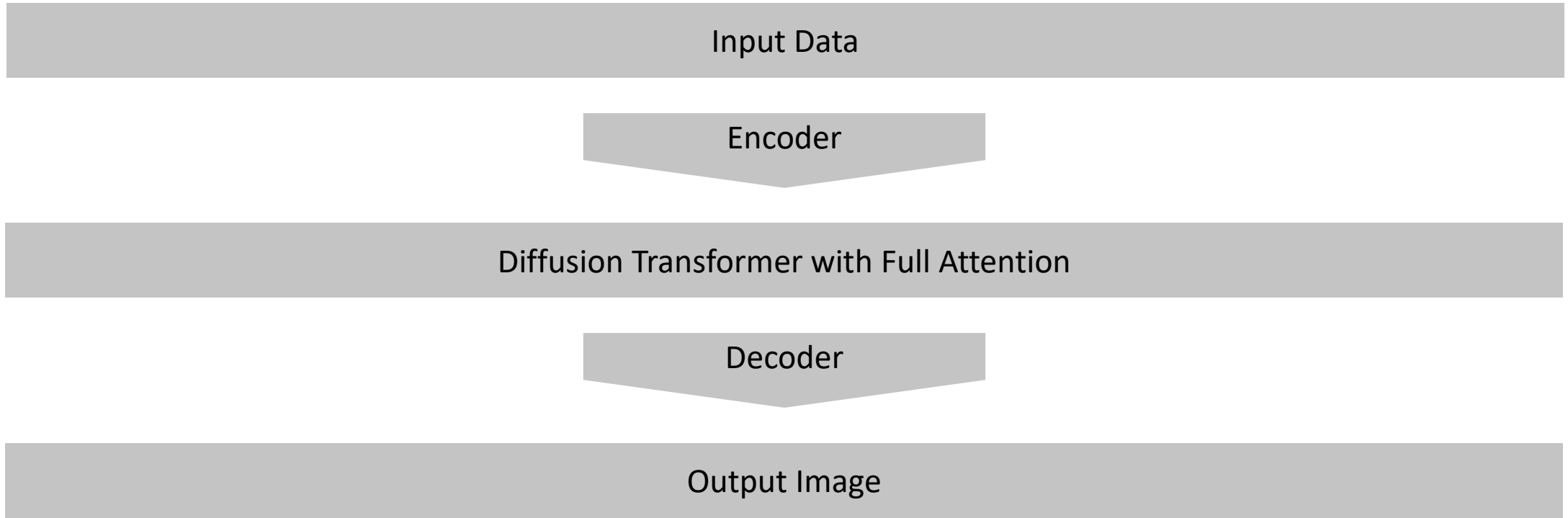
- Goal: support many advanced image edits with a **single** model



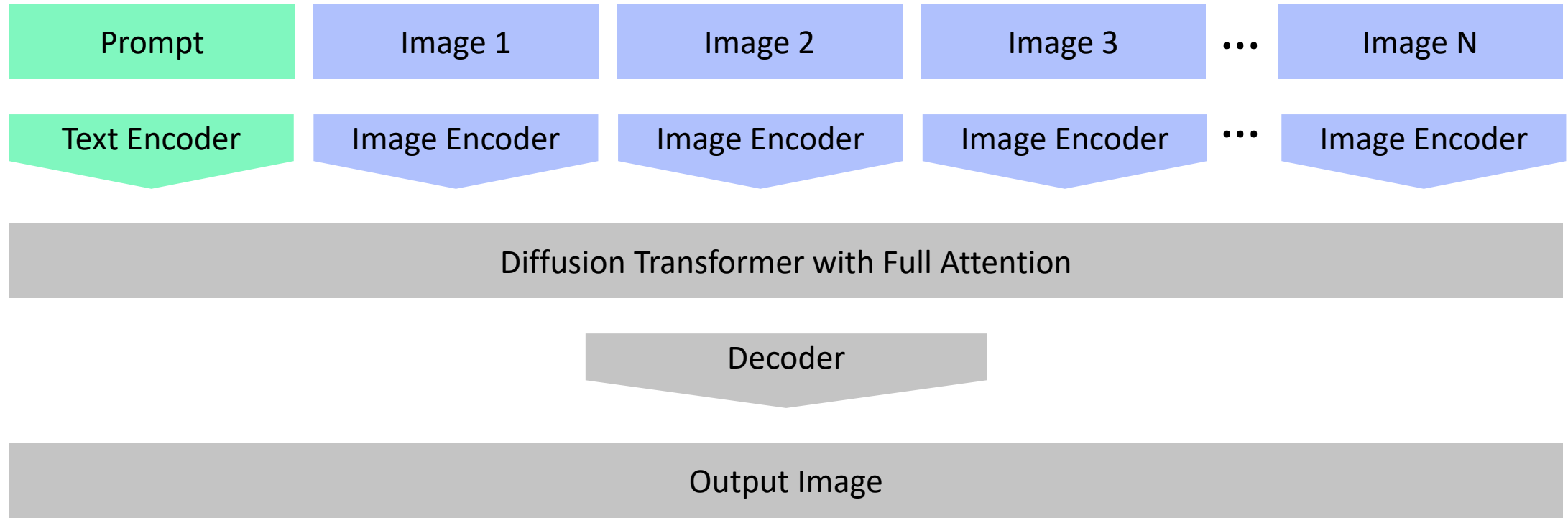
Xi Chen¹ Zhifei Zhang² He Zhang² Yuqian Zhou² Soo Ye Kim² Qing Liu² Yijun Li² Jianming Zhang² Nanxuan Zhao² Yilin Wang² Hui Ding² Zhe Lin² Hengshuang Zhao¹

¹The University of Hong Kong ²Adobe

Simplified View of the Model Architecture



Naïve Approach: Serialized All The Inputs



Naïve Approach: Serialized All The Inputs

It might work with enough data and compute resources.

But it is a hard task: learning general image editing in all its diversity and ambiguity.

For instance:

- An image may contain an element to be composited on another tone.
- Or it may be the background.
- Or it may be the mask indicated where that element is in the picture.

Most Tasks Have Some Structure

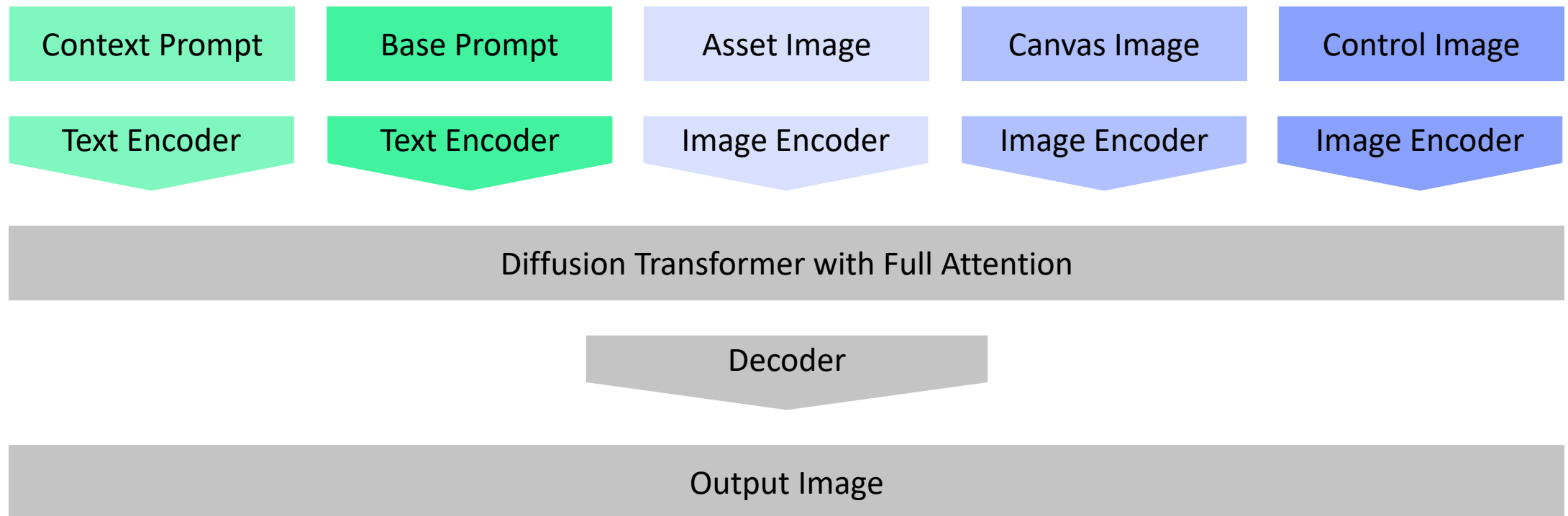
Prompt combines:

- Context, e.g., “realistic” and “uses a reference image”
- Instruction, e.g., description of the content or edits

Images can be:

- An “asset”: an element with an important role, e.g., the object to be added
- The “canvas”: the background or scene
- A “control”: complementary information to execute the prompt, e.g., a mask or a location

Specialize for the Most Popular Tasks and Exploit the Structure



Training Data from Video

- Video frames provide pseudo before-after pairs
- MLLM for generating labels



Make the person lift the weights closer together and shift the angle to focus more on the hands holding the dumbbells.



Change the background to a beach.



Change the background to a forest.



Make it snow.



Make it a sketch.



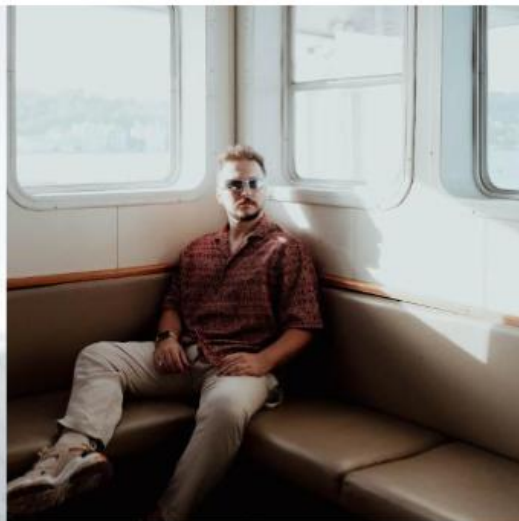
Make it watercolor.



Make it Picasso's style.



Add the sunglasses from IMG1 to the woman of IMG2.



Add the garment from IMG1 to the man of IMG2.

Take-Away Messages

Less in More

Videos as Training Data for Image Editing

The Path Towards Real-Time Video Editing

Adobe

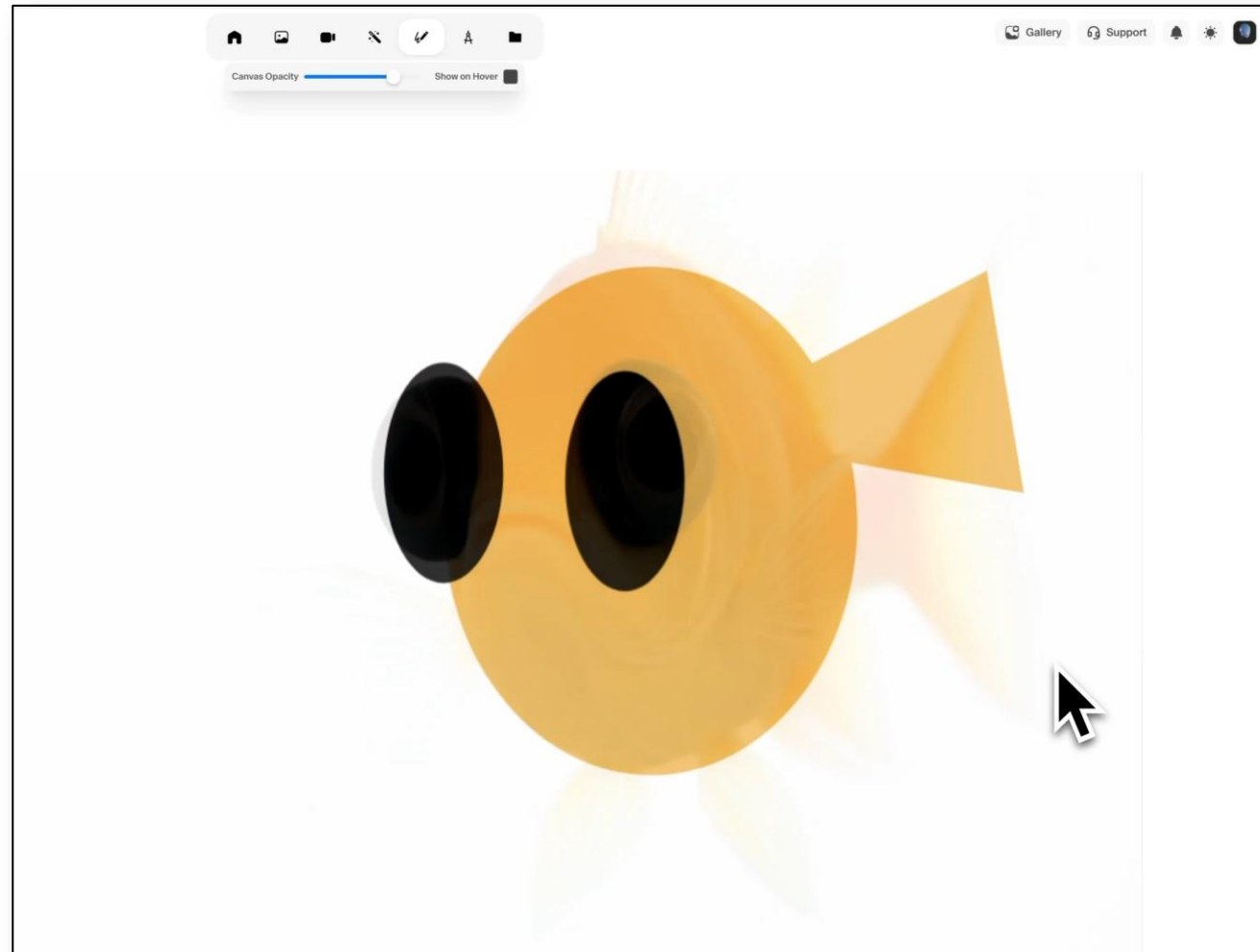
Example of Existing Solutions: Google's Genie 3

Fast, good visual quality but “frozen” world, limited controllability



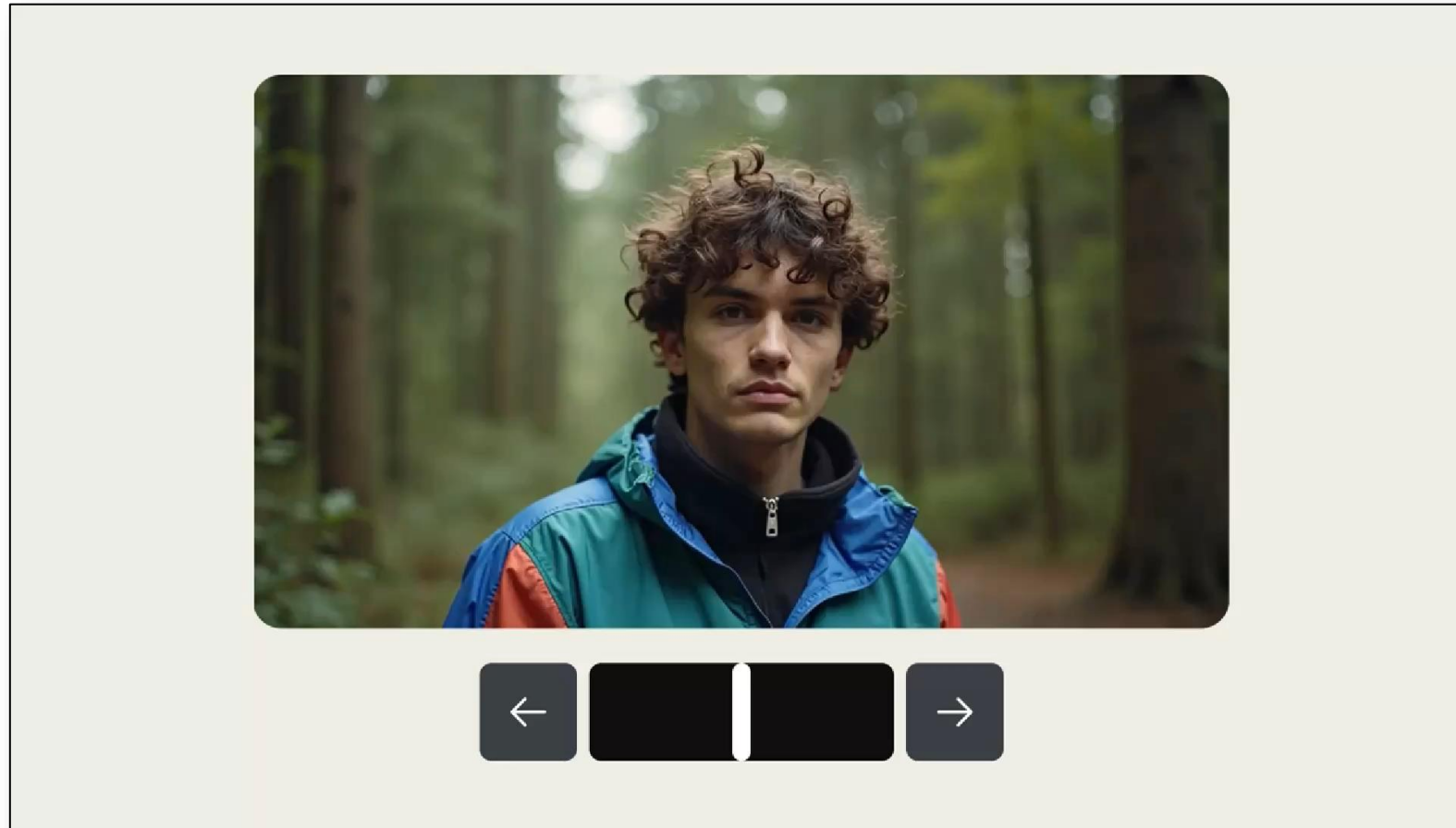
Example of Existing Solutions: Krea's Realtime Video

Fast and promising but limited, not at Adobe's level yet



Example of Existing Solutions: Runway's Realtime Controls

Great control but paused



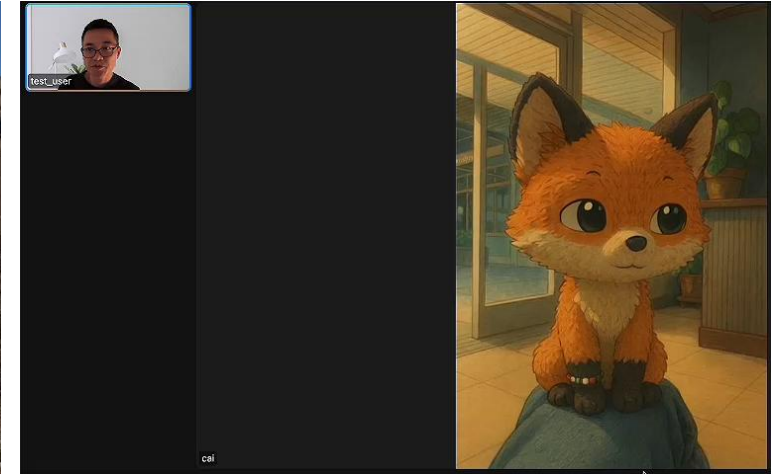
And there are more... All with similar pros & cons



Mirage, from Decart



Mirage2



Character.AI

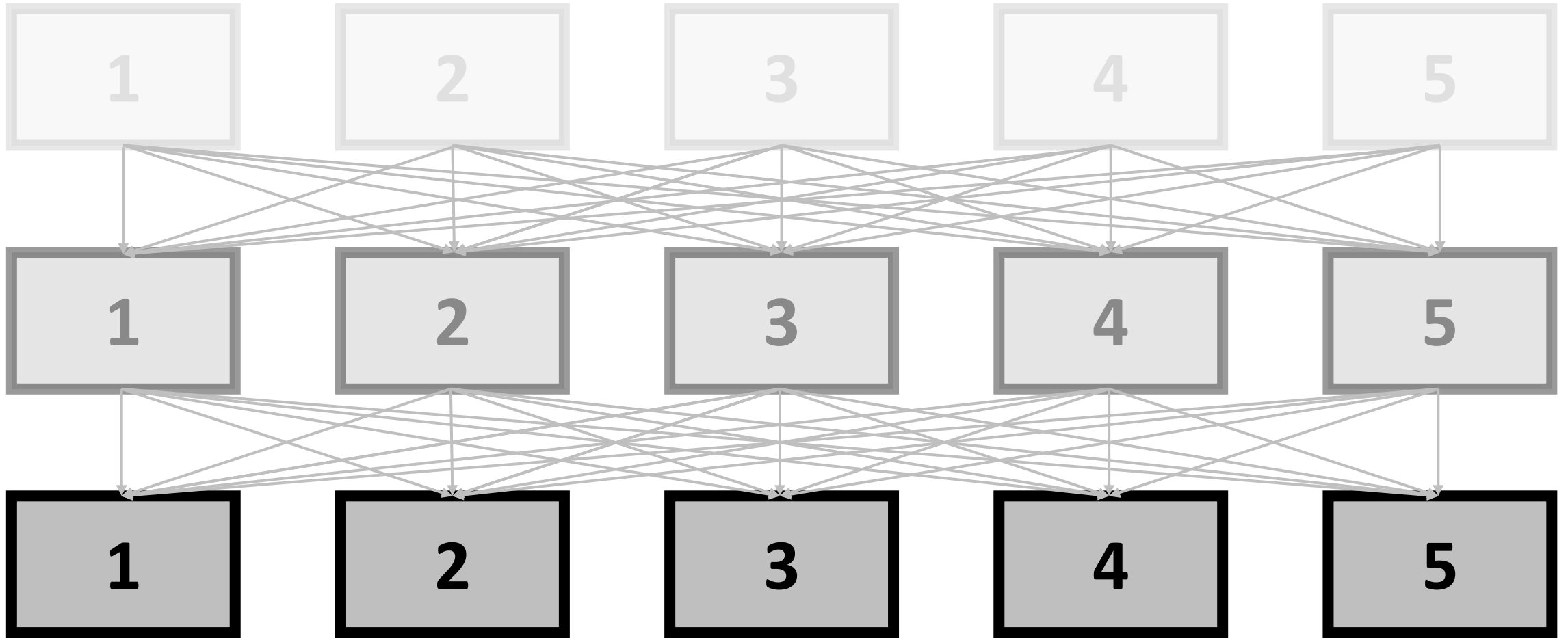


Yan, from Tencent



SeedWeed APT2, from ByteDance

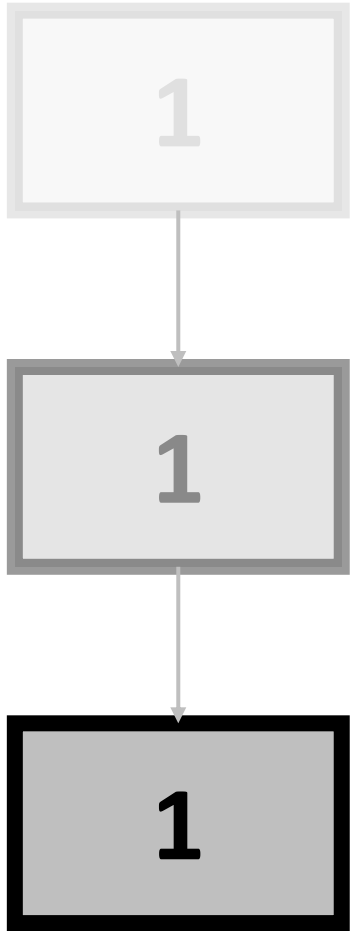
“Traditional” Video Generation with Bidirectional Attention



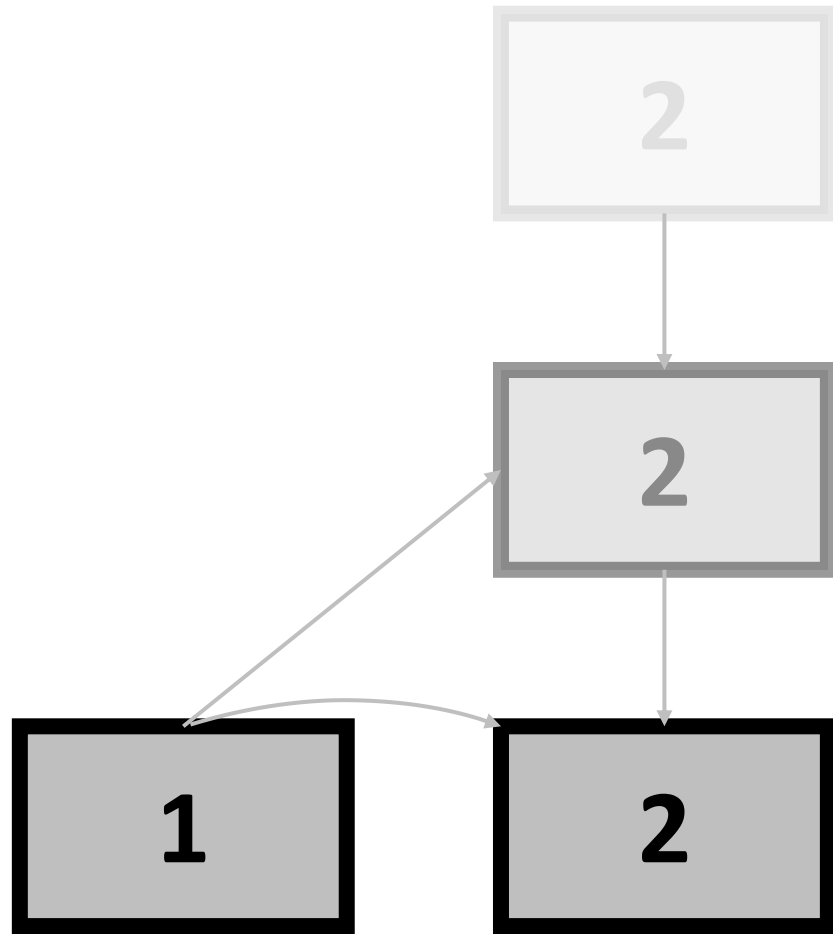
Good: maximal use of the available information

Bad: nothing is available until the very end → high latency

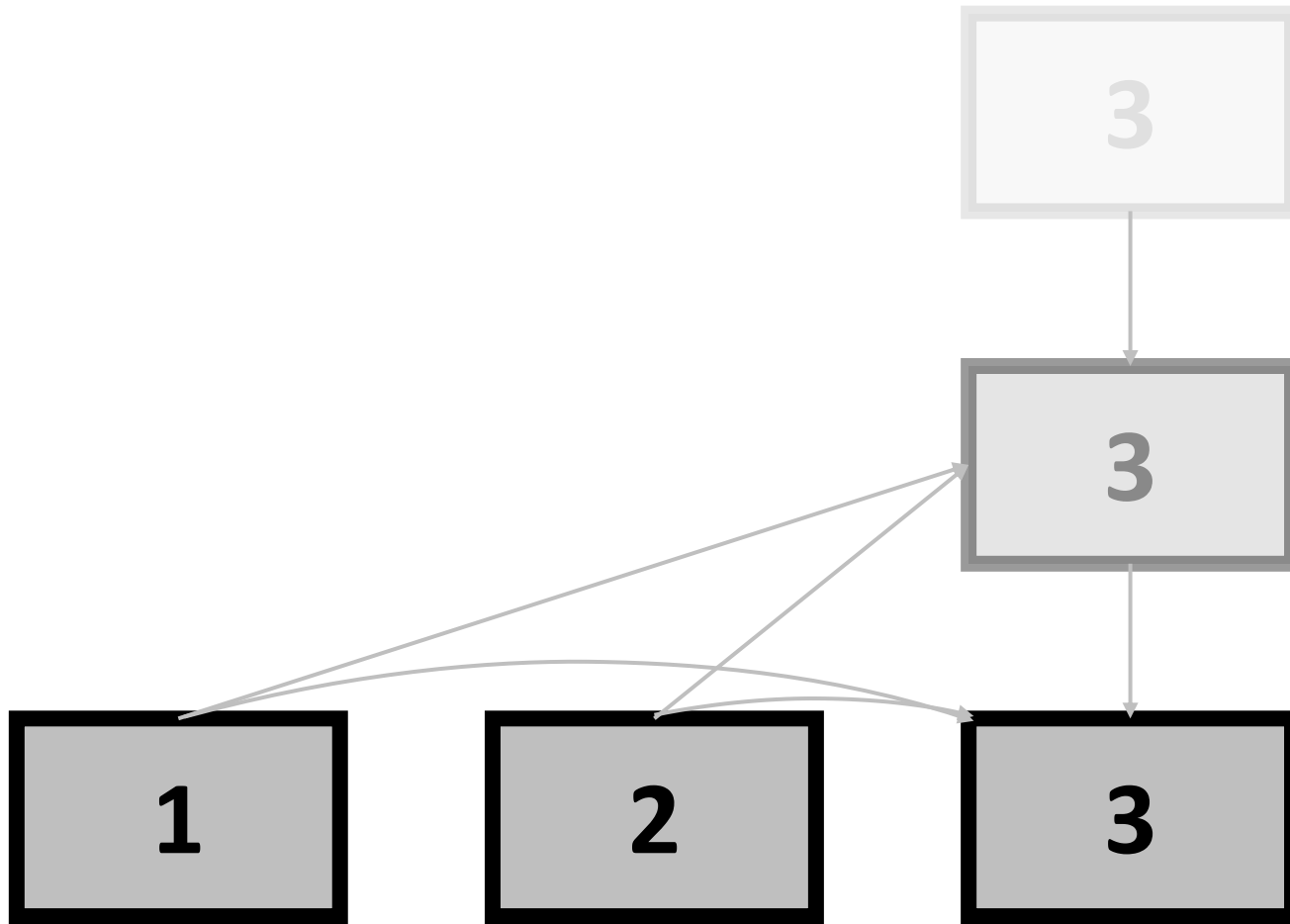
Autoregressive Video Generation with Causal Attention



Autoregressive Video Generation with Causal Attention



Autoregressive Video Generation with Causal Attention

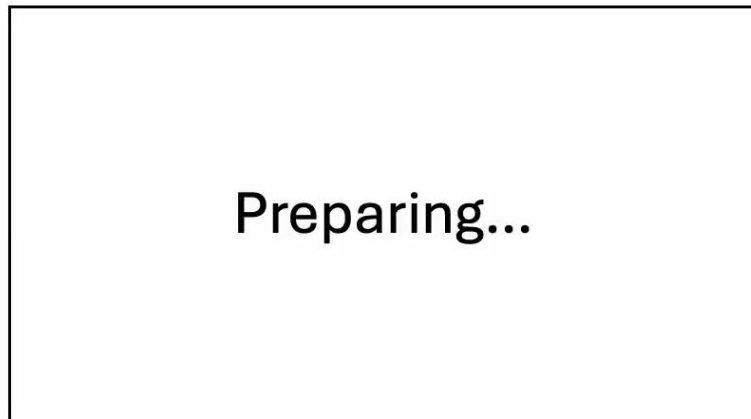


Frames are available right away one by one → real time is possible.

CausVid: A First Step Towards the New Paradigm

- CausVid model is distilled from a “traditional” slow teacher video diffusion model
- CausVid is two order magnitude faster than video diffusion model at 360P resolution on a single H100

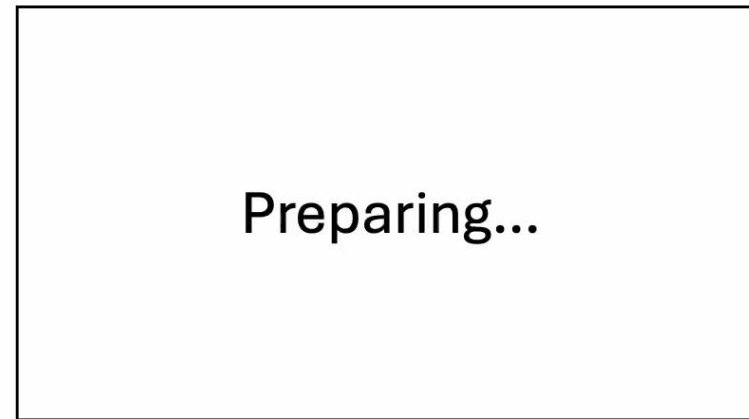
Slow Video Diffusion



Progress: 0/1

0.07 FPS

CausVid (Ours)



Progress: 0/15

10 FPS



00:00₀

Real-Time Streaming Video Generation

The camera flies across a beach.

Upload Image

Connecting...

Send

Disconnect

Video Output



Statistics

Target FPS: 8
Frames Processed: 0
Status: Connecting...
Queue Size: 0
Recording: Yes

Prompt History

6:20:35 PM
● Initial: The camera flies across a beach.

Logs

2024-11-08T00:00:05.500Z - Started recording

Self-Forcing: Bridging the Train-Test Gap in AR Video Diffusion

- Autoregressive roll-out process **during training** to mimic test-time inference



Generation from Self-Forcing



Generation from Self-Forcing

Comparison Self-Forcing with CausVid

- Oversaturation addressed
- Much less temporal flickering
- More natural motion



CausVid



Self-Forcing

Comparing Self-Forcing with CausVid

- Self-Forcing runs **17 FPS** at **480P** on a single H100 GPU



CausVid



Self-Forcing

MotionStream

[1000, 927, 771, 0]

Enter as list format, e.g., [1000, 0] or [1000, 927, 771, 0]

START ↵

PAUSE

STOP ↵

CLEAR

STOPPED



Start on click



Save on stop



Adaptive FPS

Interactive Canvas

stream_idx: 0 total_frames: 0



Generating video...

Generated Video



Enable "Save on stop" - videos will be saved locally!

MotionStream

START ↵

PAUSE

STOP ↵

CLEAR

STARTING GENERATION...

Start on click

Save on stop

Adaptive FPS

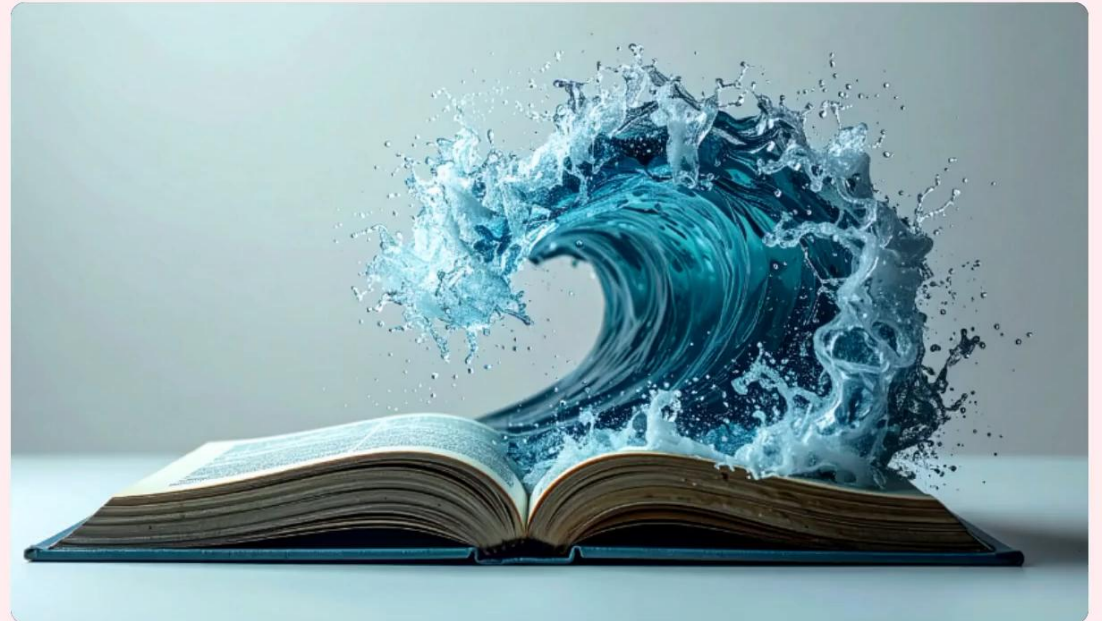
Interactive Canvas

stream_idx: 0 total_frames: 0



Generating video...

Generated Video



Enable "Save on stop" - videos will be saved locally!

MotionStream

START ↕ **PAUSE** **STOP** ↕ **CLEAR** **STOPPED** Start on click Save on stop Adaptive FPS

Interactive Canvas



Generating video...

Generated Video

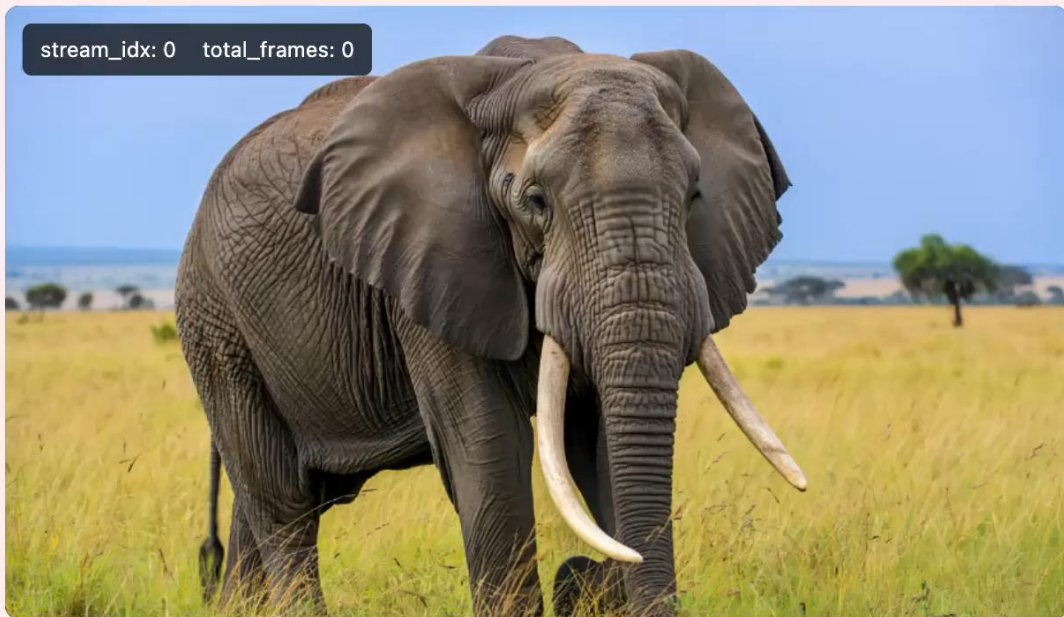


Enable "Save on stop" - videos will be saved locally!

MotionStream

START ↵ **PAUSE** **STOP** ↵ **CLEAR** **STOPPED** Start on click Save on stop Adaptive FPS

Interactive Canvas



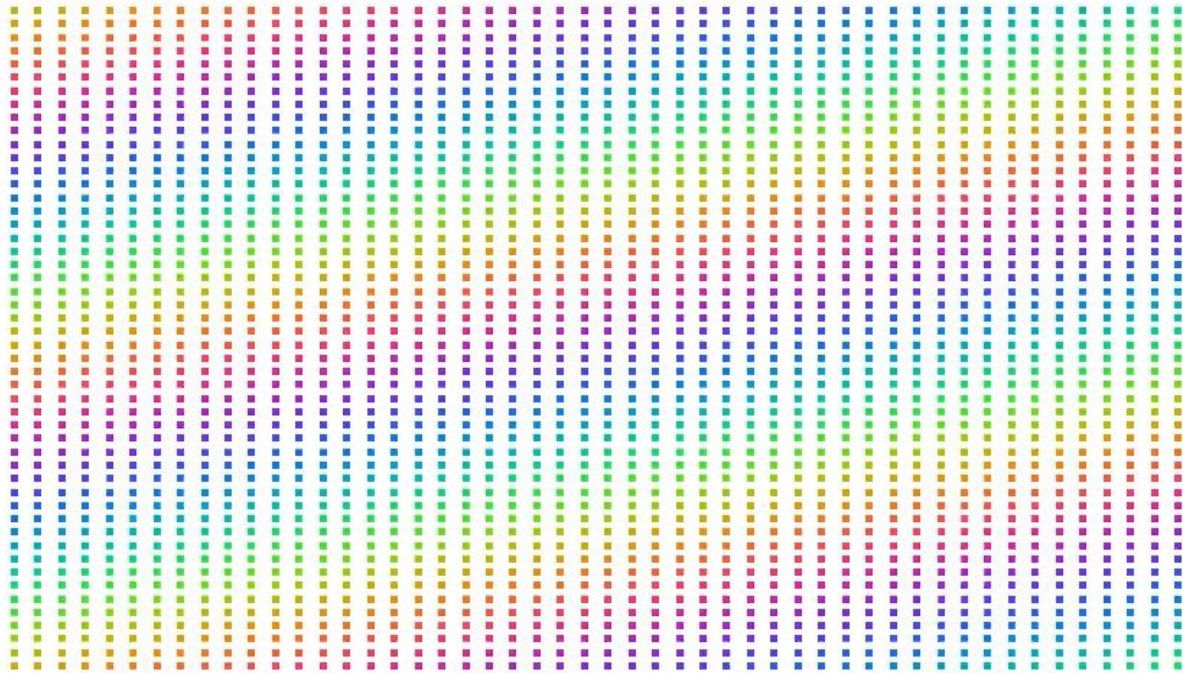
Click & drag to control

Generated Video



Enable "Save on stop" - videos will be saved locally!

Camera Control



Points “lifted” to 3d



Input image
Synthesized camera control

Pose Control

Driving video



Pose Control

Driving video



Foundation model for interaction

Take-Away Messages

Performance as a Feature

Simple Low-Level Controls Are Expressive

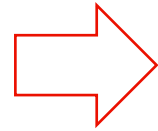
Beyond Simple Edits: X-Planner



Chun-Hsiao Yeh, Yilin Wang, Nanxuan Zhao, Richard Zhang, Yuheng Li, Yi Ma, Krishna Kumar Singh

Current Editing Models **Fall Short** on Complex Instruction-Based Image Editing

Input Image:



Edited Image:



MGIE
(ICLR'24)



LEDITS++
(CVPR'24)

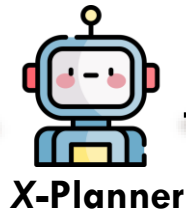


InstructPix2Pix

*“Could you make this image
Christmas?”*

- **[Planning]** MLLMs (e.g., GPT4) is needed for analyzing the complex instruction
- **[Automatic Control Guidance]** Localizing the area of interest (mask, box) is necessary for editing

Qualitative Results of X-Planner



insertion: `<box>` Add Christmas ornaments around the cat 

Make all animals look like they are celebrating Christmas?



local texture: Change the dog to have a red and white Christmas suit



background: Make the background look like a cozy snowy Christmas setting

Qualitative Results of X-Planner



Make all animals look like they are celebrating Christmas?

KUDE

insertion: `<box>` Add Christmas ornaments around the cat

Qualitative Results of X-Planner



insertion: `<box>` Add Christmas ornaments around the cat

local texture: Change the dog to have a red and white Christmas suit

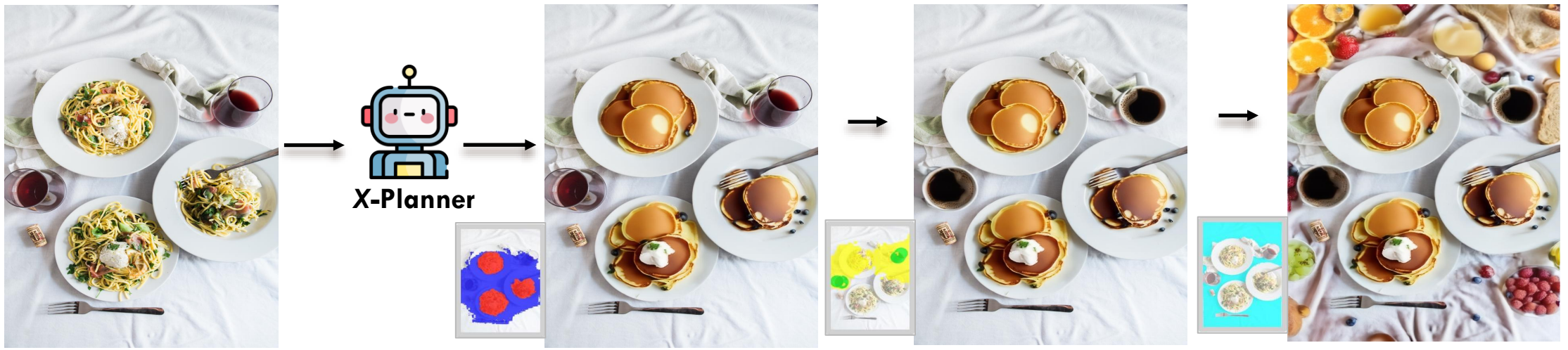
Qualitative Results of X-Planner



local texture: Change the dog to have a red and white Christmas suit

background: Make the background look like a cozy snowy Christmas setting

Qualitative Results of X-Planner



Could you change this meal setting to a breakfast theme?

replace: Replace the pasta with pancakes

replace: Change the wine glass to a coffee mug

background: adjust the background to a cozy morning setting

Qualitative Results of X-Planner



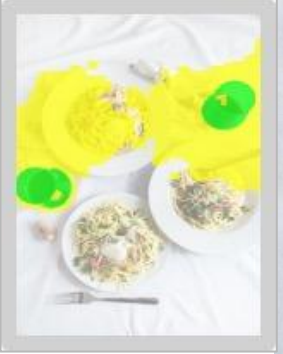
Could you change this meal setting to a breakfast theme?

replace: Replace the pasta with pancakes

Qualitative Results of X-Planner



replace: Replace the pasta with pancakes



replace: Change the wine glass to a coffee mug



Qualitative Results of X-Planner



replace: Change the wine glass to a coffee mug

background: adjust the background to a cozy morning setting

Under the Hood

- MLLM trained with automatically created training data

Additional Benefits

- Transparent process: edit list helps user understand what was done
- More control: each step can be individually adjusted or removed

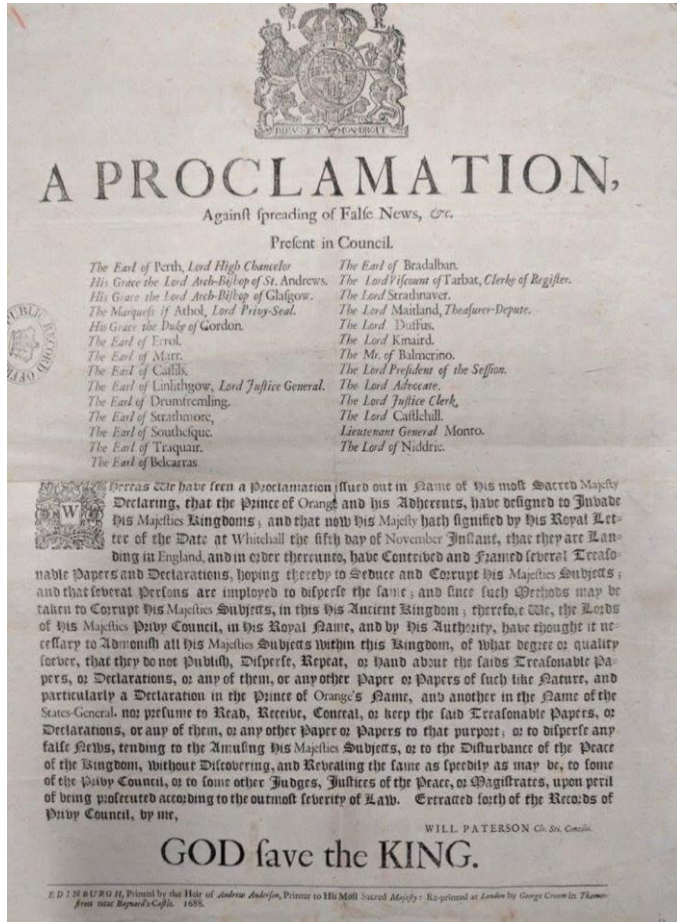


Take-Away Messages

Breaking Down Complex Tasks
Into Simple Ones Helps

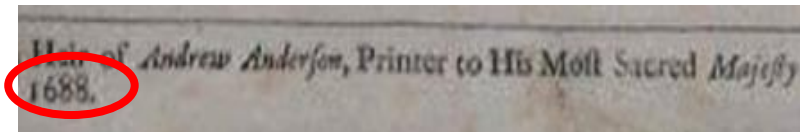
Provenance

Fake News Has Always Been a Problem...



Fake news in the time of James II

The UK National Archives



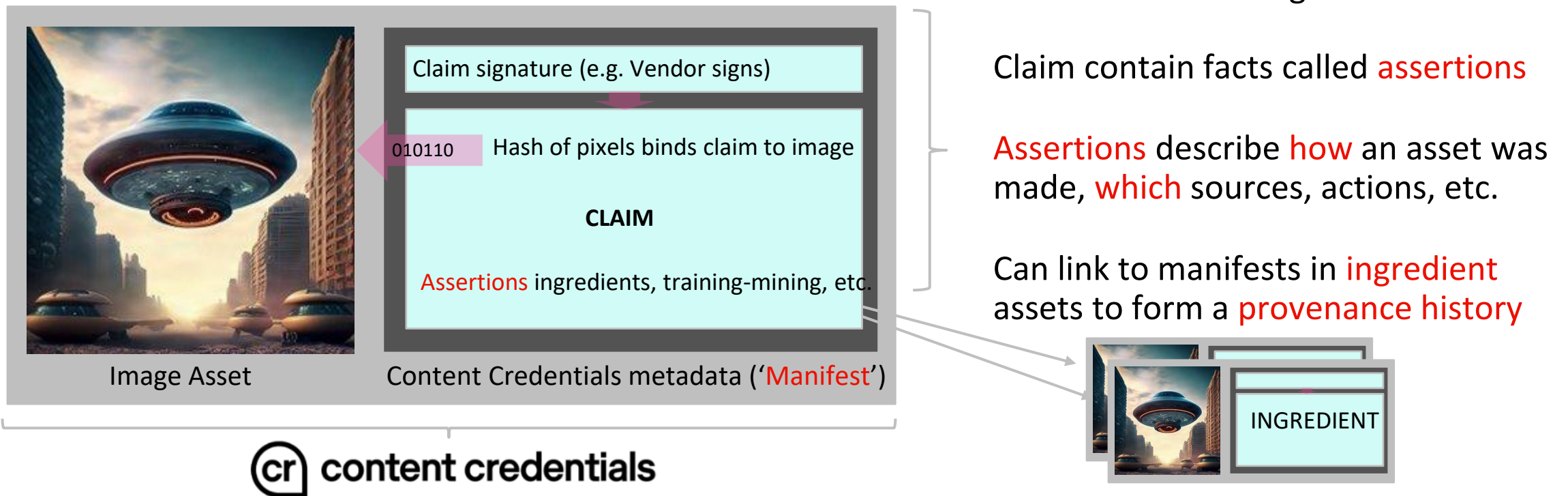
Cottingley Fairies

Photographs c. 1917 (Elsie Wright, UK)

How Can Provenance Help?

Content Credentials establish content provenance and authenticity at scale to give publishers, creators, and consumers the ability to trace the origin of media.

C2PA is an open cross-industry standard for specifying provenance of media



Why Do We Need Watermarking?

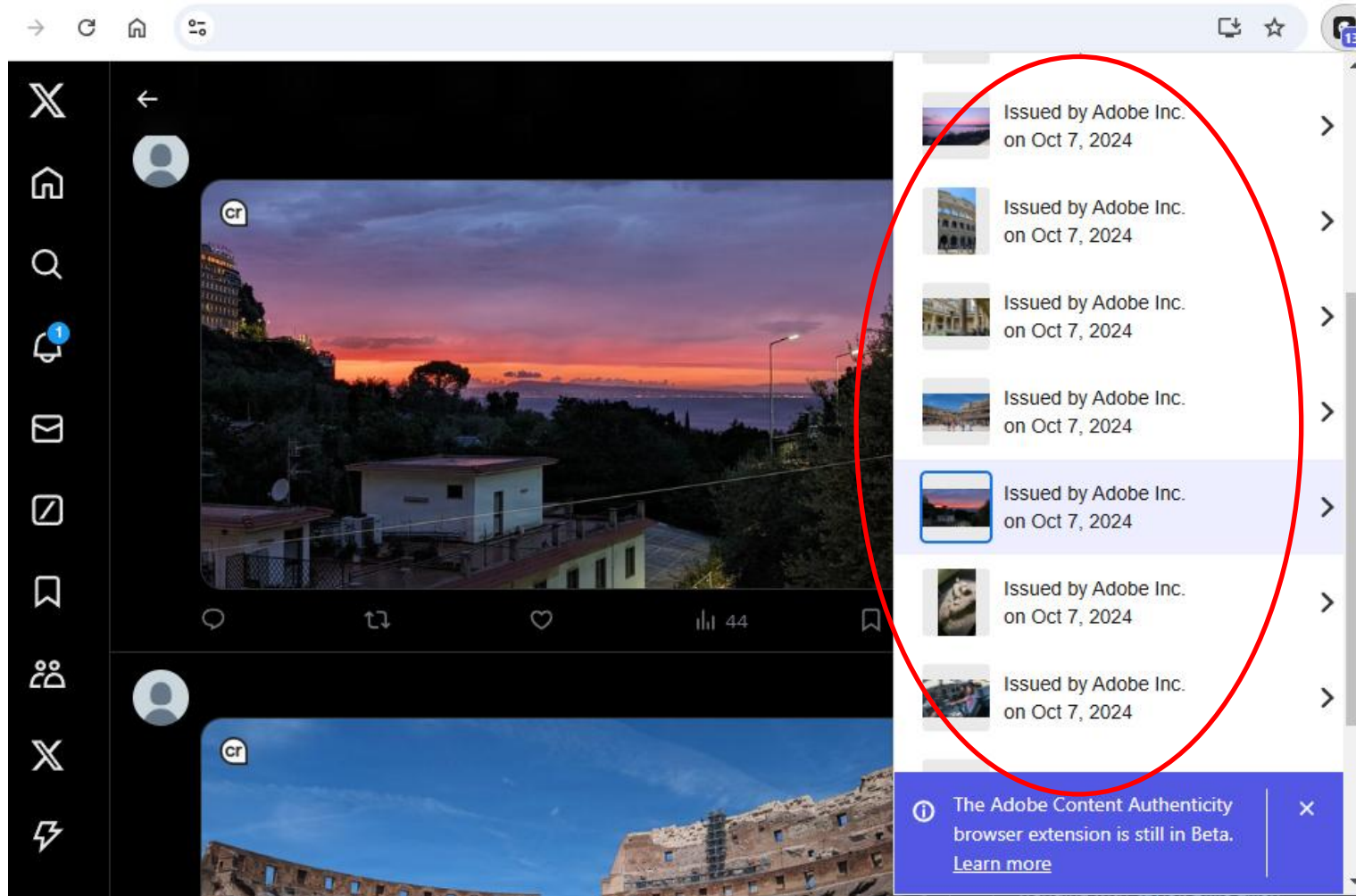
Content Credentials are often **stripped** by non-compliant content platforms or when screenshotted.



Currently all social platforms (even those that display presence of credentials) strip the metadata

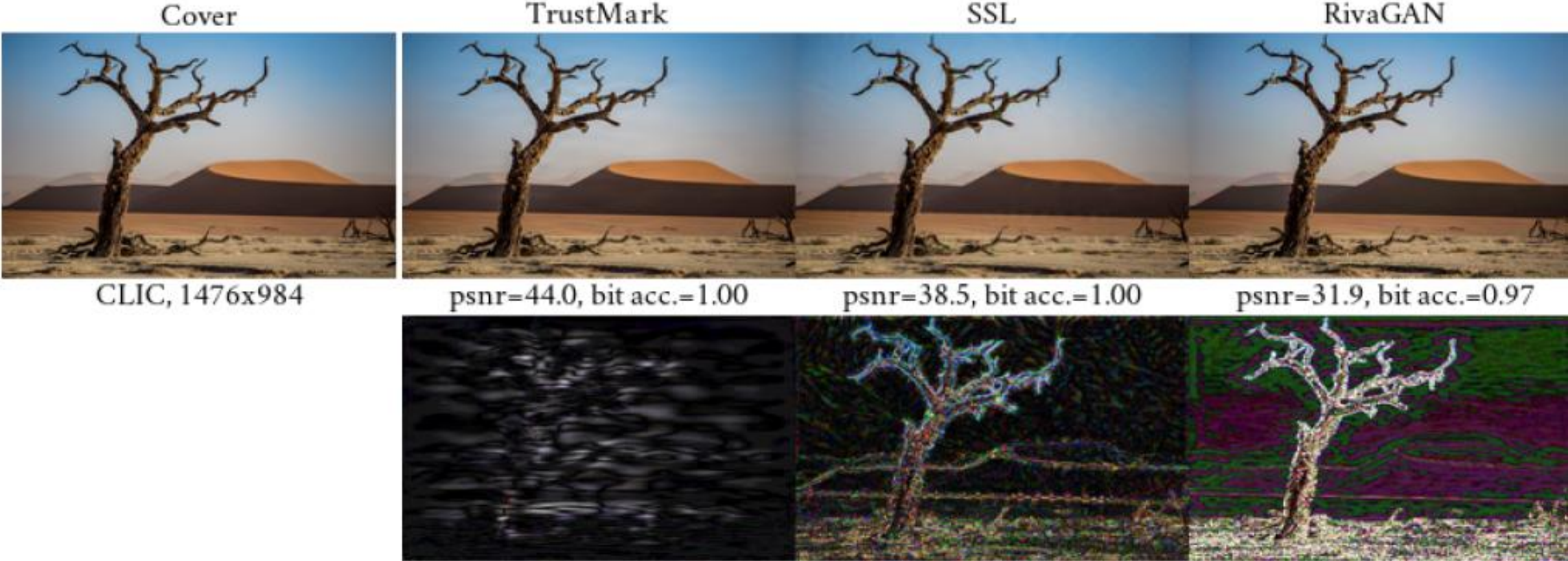
Watermarking Makes Content Credentials Stick, Anywhere on the Web

Browser detects invisible watermarks and recovers stripped content credentials

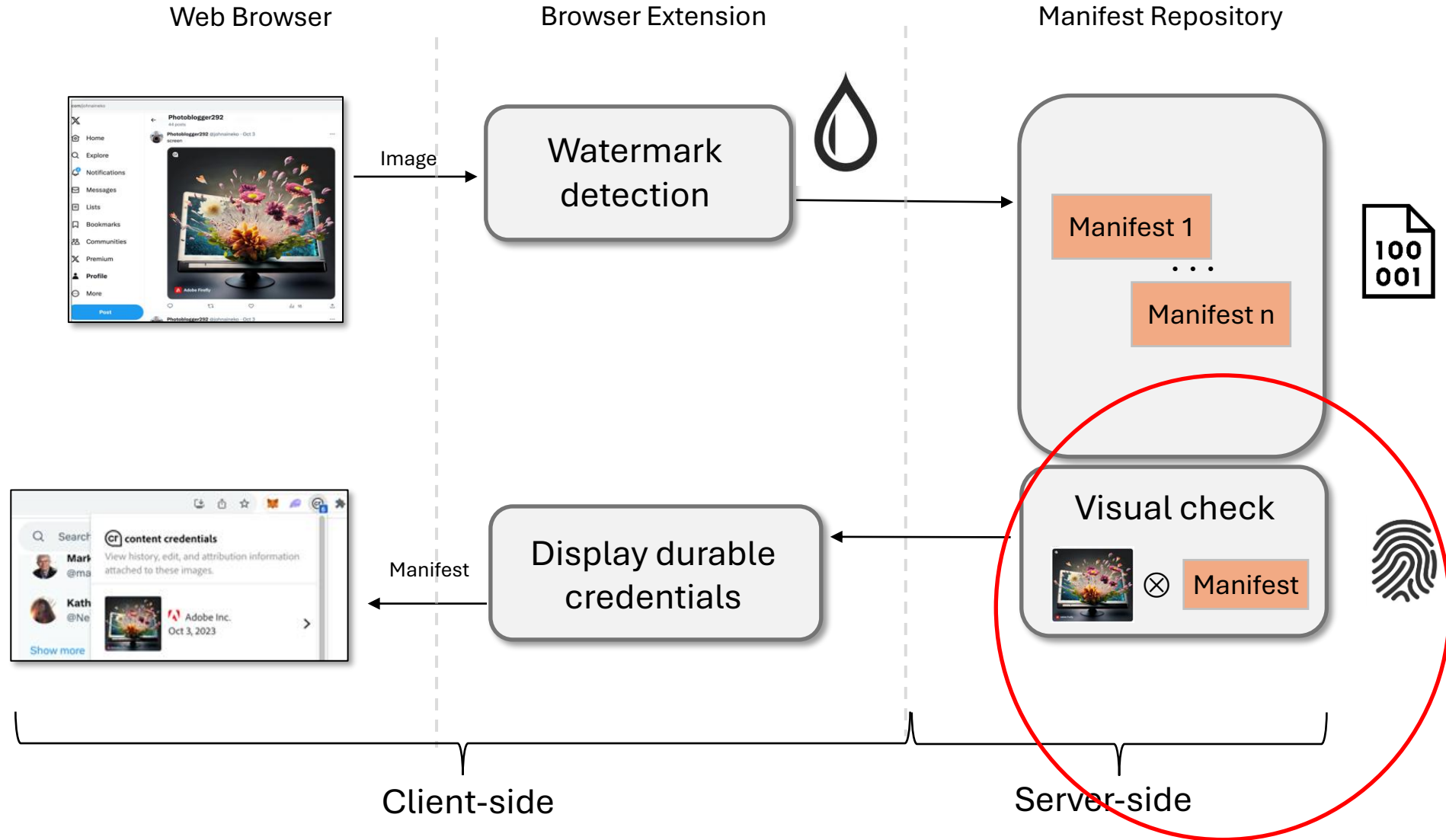


Content Credentials recovered

TrustMark: Open Source Watermarking for Provenance



Watermark Resolution with Visual Fingerprint Check



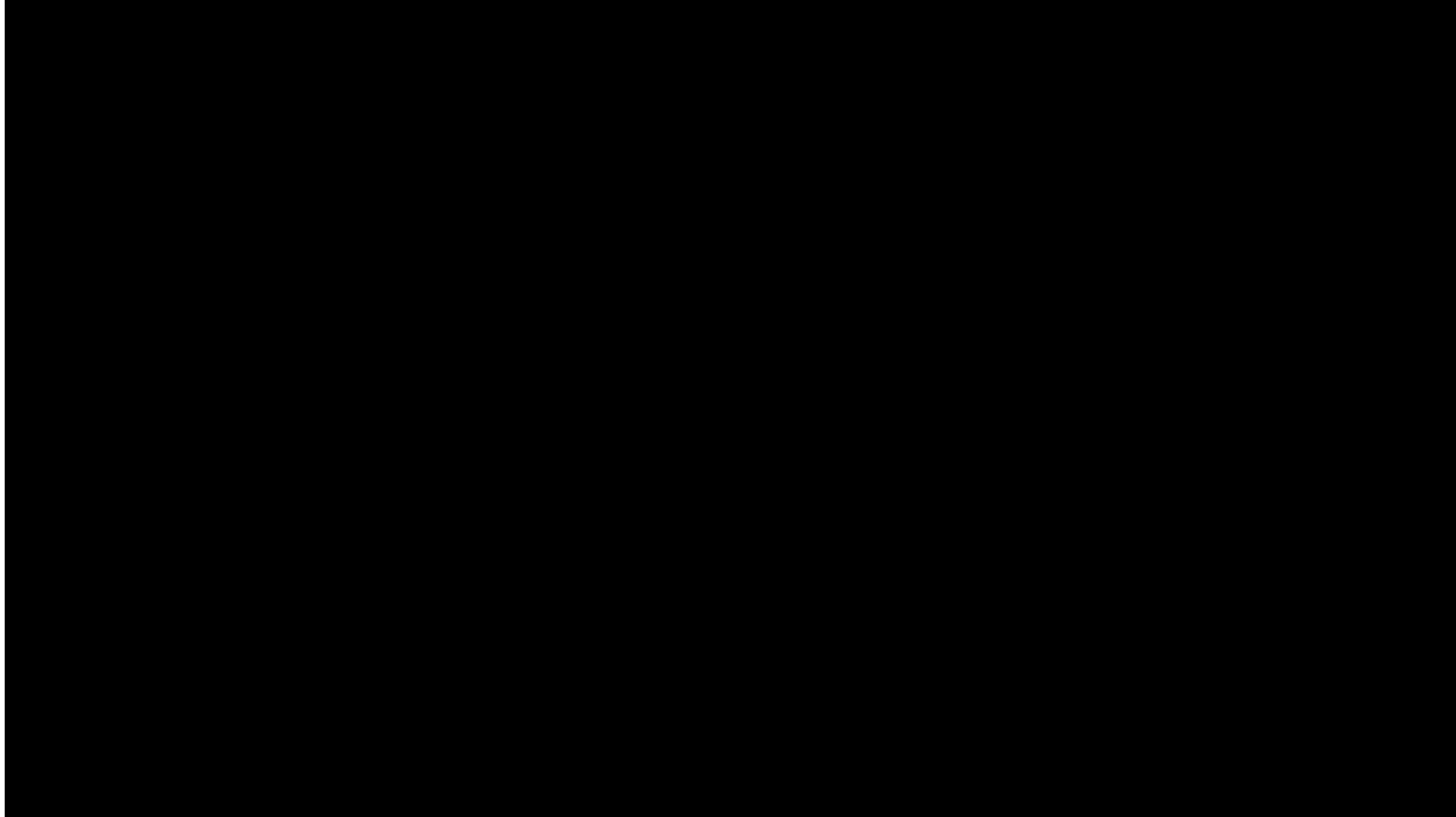
Security check against spoofing (transfer) of watermark

Take-Away Messages

Helping Good Actors Is Possible

What's Missing and What's Next?

We Used to Have Rich, Fine-Grain Controls.



These Controls Came From Analysis.

- There is a large set of studies dedicated to understanding the latent spaces of GANs.

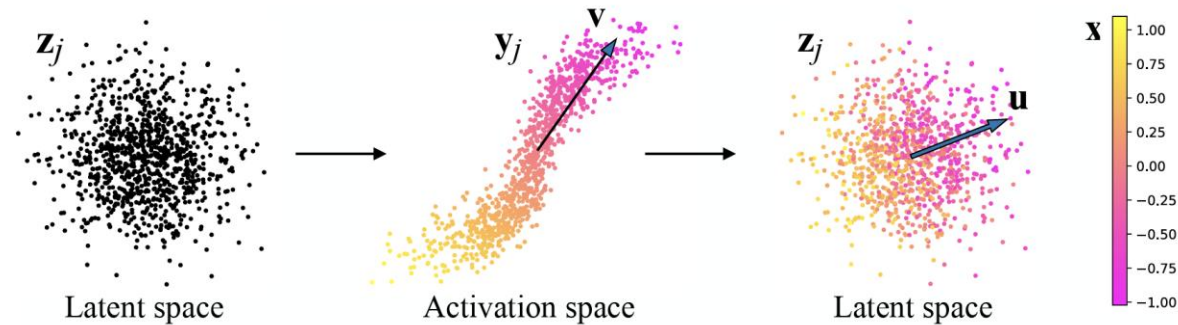


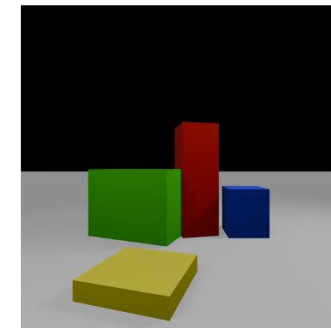
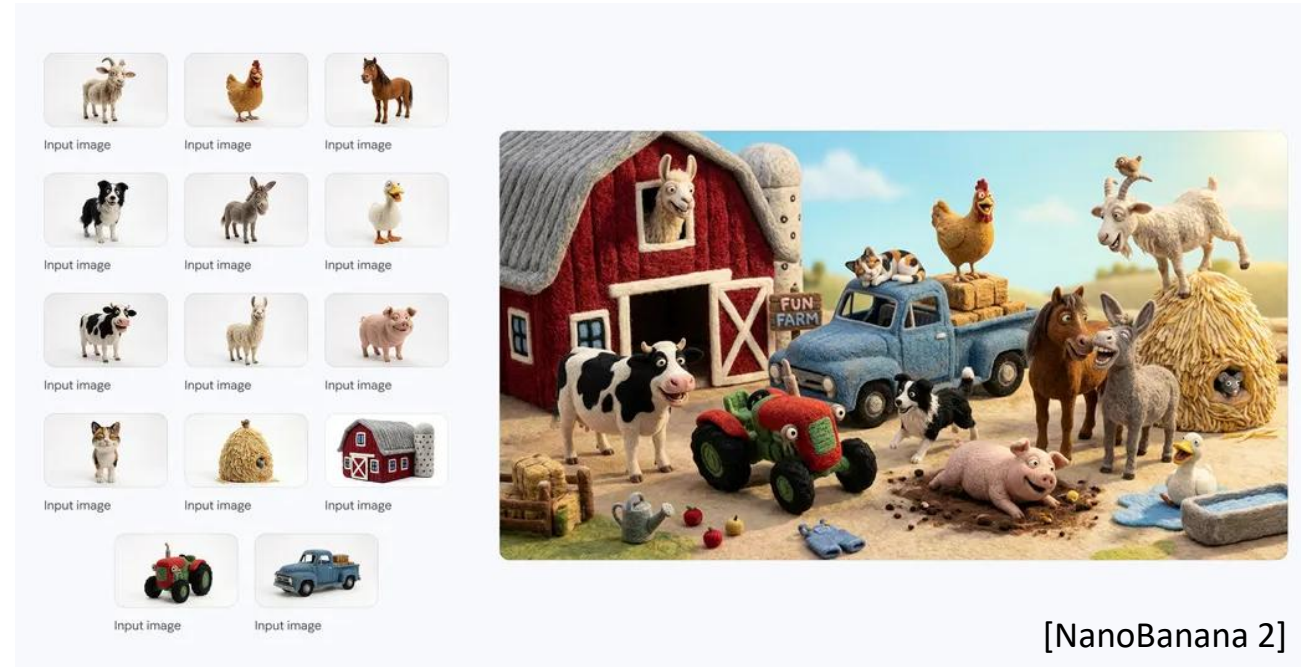
Figure 2: 2D Illustration of identifying a principal activation direction for BigGAN. Random latent vectors z_j are sampled, and converted to activations y_j . The PCA direction v is computed from the samples, and PCA coordinates x_j computed, shown here by color-coding. Finally, back in the latent space, the direction u is computed by regression from z_j to x_j .

- Challenge: activation space of AI models is huge
- Mitigation: use CLIP space

Similar Approach Is Possible with GenAI, but Identity Can Vary

More Broadly: How to Control GenAI Model?

- Prompts are limited, “prompt engineering” is particularly awful.
- Dragging points à la MotionStream, elements & characters (see fig. on the right), 3D proxies (see bottom fig.)...
- What else?

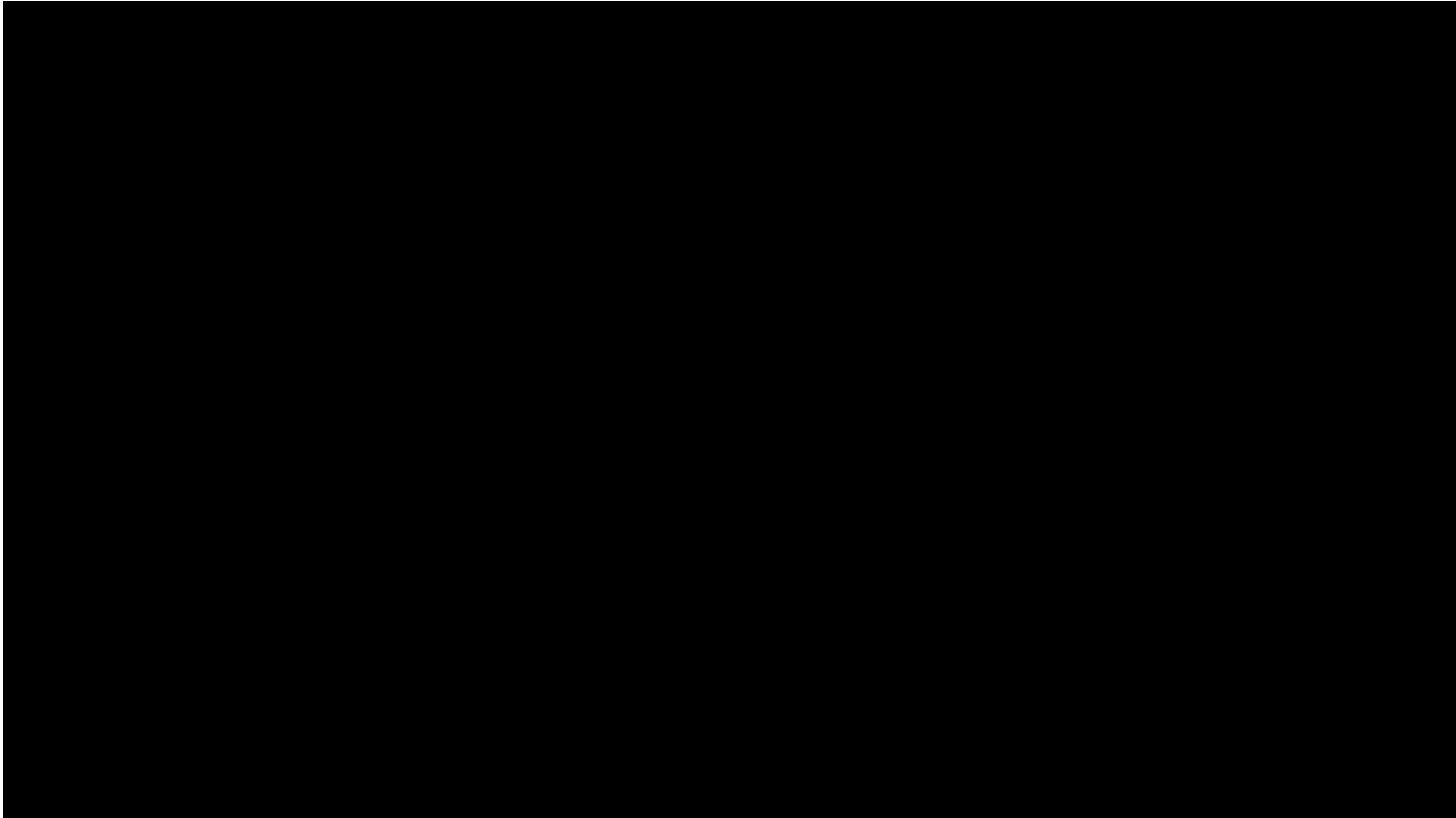


3D boxes



[SIGMA-GEN, Saha et al. 2025]

Remaining Challenge: Combining Everything Without Constrain



Thanks!

- The authors of the projects: they did all the work.
- Aaron Hertzmann and Richard Zhang for helping prepare the presentation.
- Daniel, Stefan, the INSTICC team, and all the organizers.
- You!



Wrapping Up

AI opens a new frontier
for research on creativity.

Control matters.

A lot to explore.