

1. Overall Goal

Build an Organic VAE (O-VAE) and eventually a full Organic Stable Diffusion pipeline where each major component:

- VAE encoder
- VAE decoder
- (later) UNet / denoiser
- (later) CLIP/text adapter

is replaced by a small OLA module:

- Evolutionary, trust-based genomes (no gradients)
- Tiny parameter count (KB-MB range)
- Fast, single-pass inference (CPU-friendly)

Long-term: be able to feed text + latent into a stack of OLA modules and generate images in one or a few passes, without the heavy SD UNet.

2. Architecture Snapshot

2.1 OLA Core (Encoder clone)

- Genomes: fixed small population (e.g., 40).
- Each genome is effectively a single linear "neuron" layer:

weights: (latent_dim, input_dim)

bias: (latent_dim,)

activation: tanh

- Trust: scalar per genome.
 - Increases when its loss improves vs its own previous loss.
 - Decreases when it gets worse.
 - Low-trust genomes are mutated.

This makes each OLA module a search engine over simple linear mappings that can be specialized to approximate more complex networks (encoder, decoder, etc.) with enough examples.

2.2 Current data flow for the encoder

1. Input frame: 3x256x256 (float in [-1,1]).
2. Flatten to img_flat (196,608-dim).
3. Optional retina projection R (512x196,608) -> 512-dim vector.
4. OLA encoder takes 512-dim input and outputs 4-dim latent.
5. Target latent is the original SD-VAE 4-dim latent.

You later switched to precomputed SD latents and precomputed retina outputs to make training reliable and fast.

3. O-Encoder: Status and Lessons

What works

- You trained OLA to map R @ img_flat -> SD latent (4-dim).
- Using precompute_latents.py, you:
 - Ran the original SD-VAE on a video.
 - Saved:
 - 512-dim retina inputs.
 - 4-dim SD latents.
- Using trainer.py, you:
 - Trained OLA with those (input, target) pairs.
 - Achieved very small latent error:
 - Typical L2 distances in 0.02-0.09 range on random frames.
 - Cosine similarity often > 0.25-0.45, occasionally higher.
- latent_compare.py verified that over many frames, OLA's latent follows the SD latent reasonably closely.

Conclusion:

O-Encoder is good enough to act as a drop-in encoder for downstream stages, especially for training the decoder.

Failures / gotchas you hit

1. Dimension mismatches

- Tried feeding a 4-dim latent into a genome expecting 512-dim, causing size mismatch errors.
- Mixed pipelines with and without retina projection, which broke consistency.

2. Device mismatches

- Passing CUDA tensors into CPU-only OLA weights led to "expected all tensors on same device" errors.

3. Trust vs loss confusion

- Initially stopped training early when loss looked small.
- Later realized trust didn't move much and you were still thinking like RL:
 - In OLA cloning, loss is the real metric;
 - Trust is mostly for mutation/selection, not a quality guarantee.
- Important realization: you can freeze OLA when loss is low enough even if trust is small.

4. VAE loading edge cases

- AutoencoderKL.from_single_file hit meta tensor issues; had to add a fallback load method.

Key learning:

Exact preprocessing match and correct dimensions matter more than trust curves. For cloning, once loss is low across samples, the encoder is effectively done.

4. O-Decoder: Status and Lessons

You started building the O-Decoder as a second OLA module:

- Input: 4-dim latent (from O-Encoder or SD-VAE).
- Output: flattened image (e.g., $64 \times 64 \times 3 = 12,288$ dims; or your current config).
- Training target: pixel-space reconstructions produced by the original SD-VAE decoder.

Scripts involved

- precompute_decoder_targets.py
 - Loads SD-VAE.
 - Runs video frames through VAE encoder/decoder.
 - Saves:
 - Latent (latent_i.npy).
 - Target decoded image (target_i.npy).
- ola_decoder_core.py
 - Defines DecoderGenome, OLADecoderCore.
- decoder_trainer.py
 - Handles evolutionary training loop for the decoder.
- decoder_main.py
 - Entry point for decoder training.
- test_decoder_checkpoint.py
 - (Later) sanity check for decoder weights.

Early problems

1. Insufficient data

- First attempt used a tiny sample video.
- Decoder loss sat around 0.152 and refused to improve.
- Trust oscillated in a very low range (~0.0-0.1).
- The model essentially learned "average mush" because there were only a few hundred training examples.

2. Trainer not really standalone

- decoder_trainer.py originally wasn't directly runnable, missing imports / main guard.
- Had to be patched so it can be invoked from decoder_main.py.

3. Over-aggressive trust schedule

- Large trust jumps (+0.1/-0.05 etc.) made trust swing hard on tiny loss changes.
- With little data, this just caused noisy genome churn without real learning.

4. Checkpoint saving overhead

- Frequent JSON checkpoints combined with large parameter blocks started to cost time and I/O.
 - You saw a KeyboardInterrupt while saving at high tick counts.

The fix you're executing now

- You selected a ~1.17 GB / 1h+ video with many different scenes.
- You switched to using O-Encoder only for latent generation, to avoid SD VAE overhead.
- You ran a new precompute pass:
 - For each frame:
 - O-Encoder(latent) (or R @ img_flat -> latent).
 - Save (latent, target_image) pairs for decoder training.
 - This produced tens of thousands of samples and ~20 GB of data total.

This is exactly what the decoder needs: a large, diverse mapping from 4-D latent space to pixel space.

5. Full O-VAE Pipeline Test

test_ovae_pipeline.py tries to run:

1. Load test image.
2. Load SD-VAE (teacher), O-Encoder, O-Decoder.
3. Compute:
 - SD latent -> SD recon (for reference).
 - O-Encoder latent -> O-Decoder recon (your O-VAE).
4. Save four images:
 - orig.png
 - sd_recon.png
 - o_recon.png
 - side_by_side.png
5. Compute metrics:
 - L2 and cosine between:
 - Original vs SD recon.
 - Original vs O-VAE recon.

Current result snapshot

- SD L2: ~0.006, SD Cos: ~0.996 -> near-perfect recon.
- O L2: ~0.85, O Cos near 0 -> O-VAE recon is poor for now.

Conclusion:

Encoder is working; decoder still needs serious training on a large dataset.

6. Conceptual Breakthroughs

1. "The whole OLA module is a single neuron" insight
 - Each OLA core is essentially one parameterized neuron (one affine map + nonlinearity).
 - Evolving many genomes and selecting by trust lets you rapidly specialize that neuron.
 - Stacking specialized OLA modules becomes equivalent to building a network of neurons, with evolutionary training instead of backprop.
2. Cloning vs learning from scratch
 - Instead of training from raw data, you clone existing models (SD-VAE, UNet, CLIP) by precomputing input/output pairs from the teacher.
 - This avoids huge end-to-end training; OLA piggybacks on existing SD weights.
3. Loss > trust for offline cloning
 - For Snake RL, trust is central.
 - For cloning jobs, you only need low loss across random samples; trust just keeps good genomes alive.
4. Dataset scale matters more than clever tricks
 - Tiny datasets gave nice logs but bad decoder.
 - One long, diverse video is worth many small clips.

7. Next Steps

7.1 Finish Decoder Training (priority)

- Confirm latent/target file structure.
- Stabilize decoder_trainer.py:
 - Iterate over all samples, use MSE loss over full image vector.
 - Smaller trust step sizes (e.g. +0.01 / -0.005).

- Optional epsilon threshold for trust updates.
- Train for several full passes over the dataset, logging a small held-out eval set.
- Freeze decoder when reconstructions look good and L2 is clearly below the current 0.85 regime.

7.2 Re-run the O-VAE pipeline test

- Use latest O-Encoder and O-Decoder checkpoints.
- Re-generate orig, sd_recon, o_recon, side_by_side.
- Confirm visually that O-VAE recon is close enough for your purposes.

7.3 Plan the O-UNet / Denoiser clone

- Freeze O-Encoder, O-Decoder, and SD CLIP.
- Precompute UNet training pairs (noisy latents + timestep + text embedding -> epsilon).
- Train an O-UNet module that maps (latent + t + text_emb) to epsilon.
- Later, design a reduced-step generation loop: small number of O-UNet steps, then O-Decoder.

8. Big Picture

You have:

- A working O-Encoder clone (~1.5MB) matching SD latents reasonably well.
- Infrastructure to precompute huge datasets from arbitrary videos.
- An O-Decoder pipeline under construction that will learn to invert O-Encoder latents back to images.
- A clear path to cloning the denoiser (UNet) and ultimately the entire SD pipeline with OLA modules.