**⊛ ChatGPT**

# Building a 5–60 min Mythological AI Video with Python & Open-Source Models (2026)

## 1. Understanding the Goal and Constraints

The linked video (`Neelkanth – … | 4K AI Movie | DRIFT MAXX STUDIOS`) is a high-resolution mythological short film. Recreating something similar with little money means avoiding commercial platforms that bill by the second and instead relying on **open-source text-to-video and image-to-video models**, free voice-synthesis models and **your own compute**. Long-form AI video is still produced as **dozens to hundreds of short shots (2–15 s each)** that are later stitched together. Most models top out at 10 s clips (often at 720p or 1024×576) and require a decent GPU; CPU-only workflows are extremely slow.

Commercial tools (Runway Gen-4/Gen-4.5, Pika Labs, Luma Dream Machine, Kling, etc.) provide convenient pipelines and higher resolutions but become costly when producing minutes of video. For instance, Runway's standard plan provides 625 credits per month, which equals ~52 s of Gen-4 video [1] ; Luma's Ray2 model charges 320 credits for a 10-s, 720 p clip [2] ; Pika's basic plan has 80 credits, and a 5-s Turbo clip costs 10 credits [3] . Generating a 60-minute film would burn through thousands of credits, so using open-source models on your own hardware is the only viable way to achieve "minimal cost."

## 2. Key Open-Source Video Models and Frameworks

| Model / Framework | Key capabilities (clip length, tasks, resolution) | Open-source & cost |
| --- | --- | --- |
| **CogVideoX (Zhipu AI)** | Text-to-video and image-to-video diffusion models. The **CogVideoX-1.5** release (Nov 2024) supports **10-second** videos with higher resolutions and includes an **I2V variant** that accepts an image as a background [4] . Earlier versions (5B, 2B) support three tasks—text-to-video, video continuation and image-to-video [5] . The model compresses video using a 3D VAE and generates at **16 fps** with default resolution **768 × 1360** [6] . Inference improvements allow the 2B model to run on **GTX 1080 Ti**, while the 5B model runs on **RTX 3060** [7] . | **Apache 2.0 licence** (2B and newer). Free to run locally; you need a GPU (~8–16 GB VRAM). Fine-tuning is possible via `cogvideox-factory` and can be done on a single RTX 4090 [8] . |

| Model / Framework | Key capabilities (clip length, tasks, resolution) | Open-source & cost |
|---|---|---|
| **Stable Video Diffusion (SVD / SVD-XT)** | A diffusion-based **image-to-video** model released by Stability AI. SVD generates **2–4 second** clips from a conditioning image; the SVD-XT variant produces **25 frames** instead of 14 [9]. The default resolution is **576 × 1024** [10], and the pipeline allows custom frame rates between 3–30 fps [11]. Memory-saving techniques (model offloading, feed-forward chunking, decode chunk size) reduce VRAM requirements to **< 8 GB** [12]. | **Open-source via Hugging Face**. Free to run locally (PyTorch). Requires GPU for reasonable speed; CPU-only is very slow. |
| **AnimateDiff** | A framework for turning **any text-to-image model** into an animated version by injecting a **motion module** into the frozen base model [13]. It animates images or prompts without per-model tuning, using LoRA or DreamBooth-personalised models [13]. Popular in ComfyUI workflows for stylised motion. | MIT-licensed; free to use. Clip lengths depend on workflow; typical settings produce 8–16 frames ($\leq$ 4 s). |
| **Open-Sora** | Open-source initiative mirroring OpenAI's Sora. Version 1.1 (Apr 2024) supports **2 s–15 s** videos, resolutions from **144p to 720p** and multiple tasks (text-to-video, image-to-video, video-to-video) [14]. It can generate **any aspect ratio** and offers a full video-processing pipeline [14]. Version 2.0 (Mar 2025) scales to **11 B** parameters and claims on-par performance with larger commercial models while keeping training cost to $\approx$ **USD 200 k** [15]. | Apache-licensed. Free to run if you can download the checkpoints. Clip generation limited to 2–15 s and up to 720p; longer durations require chaining. |
| **ComfyUI** | A **node-based GUI/CLI** for generative AI workflows. It lets you build complex graphs connecting image, video, 3D and audio nodes. ComfyUI is **100 % open source** and runs locally; you can create, save and reuse workflows, preview results in real time and extend functionality via custom nodes [16]. | Free; runs on your own GPU. Serves as the orchestration layer—great for building agentic pipelines and mixing models (CogVideoX nodes, SVD nodes, AnimateDiff, etc.). |

| Model / Framework | Key capabilities (clip length, tasks, resolution) | Open-source & cost |
|---|---|---|
| **Real-ESRGAN** | A super-resolution model for enhancing low-resolution images and videos. It extends ESRGAN into a practical restoration tool trained on synthetic data [17]. It supports portable executables and GPU/CPU inference [18] [19]. | BSD Licence. Free to use. Useful for upscaling AI-generated clips to higher resolutions. |
| **TTS models (XTTS-v2, Mozilla TTS, ChatTTS, MeloTTS, Coqui TTS, Mimic 3, Bark)** | A 2025 survey lists these as the top open-source text-to-speech models. XTTS-v2 offers high-quality multilingual synthesis; Mozilla TTS is customisable; ChatTTS is optimised for low-latency dialogue; MeloTTS runs on low-resource devices; Coqui TTS provides a toolkit with many voices; Mimic 3 is privacy-friendly for offline use; and Bark produces expressive speech with non-speech sounds [20]. | All open-source (various licences). They can be used locally for narration and character voices; some require GPUs for faster synthesis. |

### Why Use Open-Source Models?

- **Cost control:** You only pay for hardware/electricity rather than per-second credit fees.
- **Flexibility:** Models can be fine-tuned on your own dataset and integrated into custom pipelines.
- **Privacy:** Content never leaves your machine, important for mythological or culturally sensitive stories.
- **Limitations:** Clip duration is limited; you must generate long videos by stitching many short clips; requires GPU(s) and technical expertise.

## 3. Commercial AI Video Platforms (for comparison)

| Platform | Pricing (2026) & clip limits | Notes |
|---|---|---|
| **Runway (Gen-4 / Gen-4.5)** | **Free plan:** 125 credits one-time; equates to ~25 s of Gen-4 Turbo [21]. **Standard:** \$12/month includes 625 credits (~52 s of Gen-4 or 25 s of Gen-4.5) [22]. **Pro:** \$28/month with 2,250 credits (~187 s of Gen-4) [23]. | Credits recharge monthly; upscaling and 4K cost extra. Convenient UI but expensive for long videos. |
| **Luma Dream Machine** | **Free plan:** 8 draft-mode videos (low resolution) [24]. **Lite:** \$7.99/month with 3,200 credits; **Plus:** \$23.99/month with 10,000 credits; **Unlimited:** \$75.99/month with 10,000 credits and relaxed mode [25]. Cost per 10-second clip: Ray2 720p = **320 credits**; 1080p = **340 credits** [2]. | Produces cinematic videos but with strict credit usage; editing pipeline built-in. |

| Platform | Pricing (2026) & clip limits | Notes |
|---|---|---|
| **Pika Labs** | **Basic:** \$8/month with 80 credits. Each 5-s Turbo clip uses **10 credits**; Pro model "Pikatwists" uses **60–80 credits** [26] . **Standard:** \$28/month with 700 credits [27] . | Generates 5–10 s clips at 480p–1080p. Good for stylised scenes; still credit-based. |
| **Kling AI** | **Free plan:** 166 monthly credits, 720 p, max length 10 s [28] . **Basic:** \$6.99/month with 660 credits and 1080 p up to 30 s [29] . **Pro:** \$25.99/ month with 3,000 credits, 4 K up to 60 s [30] . Credit consumption roughly **1 credit per second** (10–25 credits depending on resolution) [31] . | Newer platform oriented toward dynamic motion; cheaper than Runway but still limited to short clips and subscription. |

## 4. Designing an Agentic Pipeline in Python

1. **Pre-production & Storyboard**
   *Break the story into scenes and shots.* A 5-minute film might need ~60–80 shots (4–5 s each); a 60-minute film may require hundreds. Determine the *camera angle, character actions, and mood* for each shot.
   *Generate keyframes.* Use text-to-image models (e.g., Stable Diffusion XL or personalised LoRA models) to create high-quality keyframe images for each shot. These images lock character appearance and composition, which is crucial because text-to-video models can drift in style.

2. **Video Generation**
   For each shot, choose a suitable model:

3. **Text-to-video** (CogVideoX or Open-Sora) when the scene doesn't need an exact starting frame. CogVideoX generates 10 s clips at up to 1360 × 768 [6] ; Open-Sora 1.1 supports 2–15 s at up to 720p [14] . You can set seeds and prompts to maintain style.
4. **Image-to-video** (Stable Video Diffusion or CogVideoX-I2V) for scenes where you want continuity with a pre-designed image. SVD's 2–4 s clips at 576 × 1024 [9] are ideal for dialogue shots; CogVideoX-I2V supports arbitrary resolutions [4] .
5. **AnimateDiff** for stylised motion or when using custom LoRA/SD models; it can animate static images into short motion sequences [13] . This is effective for background loops or mythical creatures.

**Tips for quality & efficiency**: * Keep clip durations short (5–10 s) to maintain coherence.
* Use consistent seeds and negative prompts to stabilise characters.
* Render at moderate resolution (e.g., 576 × 1024) then upscale with Real-ESRGAN. Memory-saving options in SVD can reduce VRAM to < 8 GB [12] , making generation feasible on consumer GPUs.

1. **Post-processing & Assembly**
2. **Upscaling & enhancement:** Use **Real-ESRGAN** to upscale clips to 1080p or 4 K; the model is designed for general video restoration [17] .

3. **Frame interpolation:** For smoother motion, apply a frame interpolation model (e.g., FILM or RIFE) to increase frame rates to 24 fps or 30 fps.
4. **Audio:** Generate narration and dialogues with open-source TTS models like **XTTS-v2** or **Bark**. The Northflank survey notes that XTTS-v2 provides high-quality multilingual voices, while Bark is expressive and can produce non-speech sounds [20] . For music and ambience, use royalty-free tracks or generative audio models (e.g., Stable Audio Open).

5. **Editing:** Use a non-linear editor (e.g., DaVinci Resolve or Blender's VSE) or automate merging with FFmpeg. This stage involves arranging shots, adding sound, colour grading, subtitles and transitions.

6. **Orchestrating the Workflow**

7. **Option A – Python scripts:** Build a custom package with classes for `Shot` (prompt, duration, seed, keyframe), `VideoGenerator` (wraps CogVideoX, SVD, AnimateDiff, etc.), and `Assembler` (calls FFmpeg). Use the diffusers library to load checkpoints (e.g., `StableVideoDiffusionPipeline`) and schedule tasks. 3D arrays or data classes help manage hundreds of shots.
8. **Option B – ComfyUI:** Use ComfyUI's node graph to visually connect models. ComfyUI is open-source, provides live previews and reusable workflows, and can run locally [32] . Nodes for CogVideoX, SVD, AnimateDiff and Real-ESRGAN exist; you can design a graph that takes a story file, generates keyframes, animates them and saves clips. An agentic controller (Python or JavaScript) can programmatically modify node parameters and trigger batch jobs.

9. **Automation:** To handle 60 minutes of content, write scripts that loop through scenes, monitor GPU memory and reuse models. Use caching to avoid regenerating identical backgrounds. Parallelisation across multiple GPUs or cloud instances can speed up processing; consider renting GPU servers when needed.

10. **Cost & Hardware Considerations**

11. **Compute:** A consumer GPU with 8–24 GB VRAM (e.g., RTX 3060–4090) is recommended. CogVideoX-2B can run on a **GTX 1080 Ti** [33] ; SVD can fit into **< 8 GB** with memory-saving settings [12] . CPU-only rendering is possible but **extremely slow**.

12. **Time:** Generating a 10-s clip may take 1–10 minutes depending on model and GPU. A 60-minute film could require **hundreds of hours** of compute unless you parallelise.
13. **Storage:** Keep intermediate frames and videos organised (use unique IDs and directories). High-resolution clips and audio will consume tens to hundreds of gigabytes.
14. **Licensing:** Check each model's licence. CogVideoX is Apache-2.0 for the 2B model [34] ; Stability AI's SVD has its own non-commercial terms (check `LICENSE_ML` on the repository). Real-ESRGAN uses the BSD-3-Clause license [35] . Commercial use may require additional agreements.

## 5. Putting It All Together – An Example Workflow

1. **Prepare the script**: Write or convert the story into a shot list (scene number, description, length, mood, characters).

2. **Generate keyframes**: For each shot, use a stable diffusion model (e.g., SDXL, LoRA of mythological style) to generate 1–2 high-quality images. Save images to `/keyframes` directory.
3. **Choose generation mode**:
4. Shots with complex motion (e.g., battles, flying) → use CogVideoX text-to-video.
5. Shots requiring consistency with a keyframe → use SVD or CogVideoX-I2V.
6. Stylised abstract sequences → use AnimateDiff with a LoRA.
7. **Render clips**: Loop through shots and call the appropriate model pipeline. Save each clip to `/clips`. Keep prompts, seeds and model names for reproducibility.
8. **Upscale**: Run Real-ESRGAN on clips to reach the target resolution (e.g., 1080p or 4 K).
9. **Post-process**: Optionally apply frame interpolation and colour correction.
10. **Audio & voice**: Use an open-source TTS model (e.g., XTTS-v2) to generate narration and dialogues. Add background music.
11. **Assemble**: Use FFmpeg or an editor to stitch clips chronologically, add audio, transitions and subtitles.
12. **Review & iterate**: Watch the film, note scenes requiring re-generation, adjust prompts or seeds and re-render.

## 6. Conclusion

Creating a mythological AI film similar to the provided YouTube example is achievable with **open-source models and Python** as long as you have access to a capable GPU and are prepared to invest time. The **pipeline relies on generating hundreds of short clips** (2–10 s), upscaling them, adding audio, and stitching them together. Open-source models like **CogVideoX**, **Stable Video Diffusion**, **AnimateDiff**, and **Open-Sora** provide the core generation capability, while **ComfyUI** or custom Python scripts orchestrate the process. **Commercial platforms** (Runway, Luma, Pika, Kling) offer convenience but charge by the second and are impractical for 5–60 minute films. By leveraging free models, local compute, and careful workflow design, you can produce high-quality AI videos at minimal cost while retaining creative control.

---

[1] [21] [22] [23] AI Image and Video Pricing from $12/month | Runway AI
https://runwayml.com/pricing

[2] [24] [25] Dream Machine Pricing: Plans and Credits | Luma AI
https://lumalabs.ai/pricing

[3] [26] [27] Pika
https://pika.art/pricing

[4] [5] [7] [8] [33] [34] GitHub - zai-org/CogVideo: text and image to video generation: CogVideoX (2024) and CogVideo (ICLR 2023)
https://github.com/zai-org/CogVideo

[6] [2408.06072] CogVideoX: Text-to-Video Diffusion Models with An Expert Transformer
https://arxiv.org/abs/2408.06072

[9] [10] [12] Stable Video Diffusion
https://huggingface.co/docs/diffusers/en/using-diffusers/svd

[11] Stable Video — Stability AI
https://stability.ai/stable-video

[13]  AnimateDiff

https://animatediff.github.io/

[14]  [15]  GitHub - hpcaitech/Open-Sora: Open-Sora: Democratizing Efficient Video Production for All

https://github.com/hpcaitech/Open-Sora

[16]  [32]  ComfyUI | Generate video, images, 3D, audio with AI

https://www.comfy.org/

[17]  [18]  [19]  [35]  GitHub - xinntao/Real-ESRGAN: Real-ESRGAN aims at developing Practical Algorithms for General Image/Video Restoration.

https://github.com/xinntao/Real-ESRGAN

[20]  Best open source text-to-speech models and how to run them | Blog — Northflank

https://northflank.com/blog/best-open-source-text-to-speech-models-and-how-to-run-them

[28]  [29]  [30]  [31]  Kling AI Pricing Breakdown & Comparison: Get Your Best

https://www.litmedia.ai/resource-video-tips/kling-ai-pricing/