# ChatGPT

# ReelForge v1 Backend Technical Manifesto

## Microkernel-Like Architecture and Orchestrator Core

ReelForge's backend is designed with a **microkernel architecture** – a minimal central "kernel" orchestrating a suite of modular AI components [1]. The **core system** is a Python FastAPI server that handles API requests and coordinates tasks, while specialized **plug-in modules** (skills for scriptwriting, voice synthesis, video editing, etc.) operate independently. This separation of core and extensions promotes flexibility and maintainability, critical for a lean 2-engineer team [1].

At the heart is a **custom Orchestrator AI agent** that functions as the "brain" of the microkernel. It sequences the workflow, invokes the necessary modules, and handles decision logic. Each module (plugin) is a loosely-coupled service – often implemented as a Celery task – that the orchestrator can call as needed. The orchestrator (or **"agent"**) tracks the overall job state in a database (PostgreSQL) and uses inter-process communication (via Redis and Celery message queues) to dispatch work to modules asynchronously.

This architecture ensures **extensibility**: new features (say a different TTS engine or a new video effect module) can be added as new plug-ins without altering the core. Modules communicate with the core via well-defined interfaces (function calls or queued tasks returning results). The core itself only handles essential logic – task scheduling, data passing, and result aggregation – keeping it simple and robust [1].

We leverage **Celery** for workflow orchestration. Celery is a proven distributed task queue in Python that lets us offload heavy or long-running tasks to background workers, keeping the FastAPI web layer snappy [2]. With Celery, we can define a pipeline of tasks (using chords or chains) corresponding to our video generation stages, and Celery will manage their execution and retries automatically [3]. Redis (via a managed service like Upstash) serves as the Celery broker and result backend. This setup prevents blocking the API – the user's request to generate an ad can return immediately with a job ID, while the actual work happens in background workers. Celery's robust scheduling and retry features ensure that if a module fails (say the video render crashes), it can retry or fallback gracefully, which is essential for an autonomous system.

**Tech Stack Summary:** The backend stack is **Python** (fast to iterate, rich AI ecosystem). We use **FastAPI** for the HTTP API (to interface with the Next.js frontend and any internal tools). **PostgreSQL** (via a managed service like Supabase) stores persistent data: user accounts, brand profiles, job statuses, ad metadata, and performance metrics. **Redis** (e.g. Upstash free tier) is used for ephemeral data: task queues, caching, and short-term state. All components run in Docker containers orchestrated via Docker Compose for easy development and deployment. This architecture is lightweight yet scalable – we avoid monolithic design and instead have a nimble core orchestrating modular services, much like an OS kernel spawning processes.

# Real Video Rendering with FFmpeg (No Templates)

ReelForge generates **fully rendered real video files** by programmatically stitching media content with FFmpeg – rather than relying on pre-made templates or simple animations. This means the system has fine-grained control over every frame, enabling dynamic composition of voiceovers, B-roll footage, music, subtitles, and effects into a final MP4 output. **FFmpeg** is invoked via subprocess commands or Python bindings to handle all video editing tasks. We favor direct FFmpeg CLI usage for reliability and performance, given its powerful filter capabilities (video/audio transitions, overlays, scaling, etc.) [4] .

Using FFmpeg ensures **maximum flexibility**: *any* visual effect or transition that professional editing software can do, we can achieve through FFmpeg's filters [4] . For example, adding a crossfade between two video clips or a swipe transition is done with FFmpeg's built-in `xfade` filter – no expensive rendering engine or manual editing needed [4] . We will output videos in common formats (e.g. H.264 MP4) suitable for social platforms.

In practice, the Editing module will construct a complex FFmpeg filter graph to assemble the video timeline. It will concatenate multiple video clips (B-roll segments), overlay the AI-generated voiceover and music tracks, apply transitions between segments, and burn-in subtitles or graphics – all in one rendering pass. By using FFmpeg's **filter_complex**, we minimize intermediate files and keep the pipeline efficient. This approach yields professional-quality results while keeping costs near zero, leveraging open-source technology.

Critically, this design avoids the constraints of template-based video tools. **There is no rigid template** – each ad can have a unique structure, pacing, and style decided by the AI. The system isn't limited to swapping text in a preset animation; it truly *edits* the video from raw components. This makes ReelForge's output far more flexible and tailored. For instance, competitor tools often produce formulaic videos from a handful of templates, but with our FFmpeg approach we can dynamically change scene order, timing, transitions, and overlays on a per-ad basis.

# AI Content Generation Pipeline (Modular Components)

The end-to-end process of generating a short-form video ad is broken into modular pipeline stages. Each stage is handled by a dedicated component (often an AI model or specialized service) and orchestrated by the core agent. Below we detail each component's pipeline, including tools/libraries, data flow, and how variants or retries are handled:

## 1. Concept & Trend Extraction

**Purpose:** Analyze the brand/product and current digital trends to derive a creative concept or hook for the ad. This sets the high-level direction (theme, angle, style) for the ad so that it's both on-brand and culturally relevant (e.g. aligns with a trending meme or popular format).

**Process:** The agent takes as input the brand brief (product details, value propositions, target audience, any brand preferences) and optionally market research data. It then uses an LLM (e.g. GPT-4 or GPT-3.5 via OpenAI API) to brainstorm ad concepts. We craft a prompt that might include: a summary of the product, the brand's tone (from the brand profile), and a request for *several creative video ad ideas*. The LLM could be

asked for multiple concepts or just one strong concept with rationale. If multiple concepts are returned, the agent will choose one based on heuristic scoring – e.g., it may prefer a concept that leverages a current trend or one that matches the brand's past successful themes.

To incorporate **trends**, the system can tap into external data sources. For MVP, this may be as simple as feeding the LLM some context like "Here are a few current TikTok trends: X, Y, Z" or using a static list of recent viral challenges. In future, we might integrate with social media APIs or Google Trends to dynamically pull what's hot. The agent could maintain a small knowledge base of recent ad trends (updated periodically) and embed that info in the prompt. For example: *"Product: Fitness Drink. Trend: Many TikTok ads use fast jump-cut transitions and meme captions. Task: Propose a concept combining the product's benefit with a TikTok-style challenge."*

**Output:** A structured **Ad Concept** object, e.g.:

```
{
  "concept": "Highlight refreshing energy boost via a morning routine trend",
  "hook": "Use the popular #5AMClub TikTok theme to show the drink powering a
sunrise workout",
  "tone": "Upbeat, motivational",
  "key_points": ["Natural ingredients", "No sugar crash",
"limited-time discount mention"]
}
```

This concept metadata will accompany the job through the pipeline. It informs later steps (script writing will incorporate these key points and tone; visuals will reflect the morning workout theme, etc.).

**Tools & Libraries:** OpenAI GPT-4 (for rich idea generation) or a local model (if cost must be minimized, perhaps a fine-tuned Llama-2 for brainstorming). We may also use simple web scraping or API calls to get trending keywords (if implemented). In later versions, a vector database (Pinecone) can store known successful concepts and their performance, allowing the agent to retrieve inspiration from what worked for similar brands.

**Retries/Variants:** If the concept returned is weak or too generic (agent can detect if the LLM output is low-confidence or repetitive), the agent can **prompt GPT again with adjusted instructions**. For example, if the first concept is bland, it might prompt "Give a more daring/creative idea". The system could generate 2–3 concepts in parallel (as separate tasks) and automatically pick the best – using a heuristic or even a secondary AI evaluation. A heuristic might be checking for uniqueness (no concept should be just "make a slideshow of product images") or alignment to a trend (maybe using a list of trend keywords and seeing if any appear). Because this step is relatively fast (just text), trying a few variations costs little.

## 2. Script Generation

**Purpose:** Create a compelling short script (voiceover narrative and any on-screen text cues) based on the chosen concept. This is essentially the copywriting step – producing the words that will be spoken or displayed in the ad.

**Process:** The agent now prompts an LLM to write the **ad script**, providing it with: the selected concept (from step 1), key product/brand info (unique selling points, call-to-action), and format guidelines (e.g. *"15-second script, 3-5 short sentences, upbeat tone, end with brand name and slogan"*). We use GPT-4 for highest quality language (or GPT-3.5 if cost-saving, as it can usually produce decent ad copy). The prompt might explicitly request a breakdown by scene or line, for easier downstream processing. For example:

*"Write a script for a 15-second video ad. The ad concept: {concept_hook}. The brand: {brand_name}, which is {brand_description}. Key points to mention: {list}. Tone: {tone}. Provide the script as 3-4 numbered lines (each line ~1 sentence) for separate scenes."*

The LLM's output is parsed into a structured format, e.g. an array of scenes. Each scene entry includes the **voiceover text** and possibly a suggestion for visuals. We can encourage the model to include a visual cue per line (e.g. "Scene 1: {VO line} [Visual: …]"), which we'll use to guide B-roll selection. Even if the LLM doesn't provide explicit visual hints, the voiceover text itself will be used to infer visuals.

**Data Schema:** We store the script as an ordered list of segments, e.g.:

```
[
  {
    "scene": 1,
    "voiceover": "Mornings feel impossible? Boost your energy with GlowDrink.",
    "visual_hint": "Tired person at dawn turns upbeat after drink"
  },
  {
    "scene": 2,
    "voiceover": "Packed with natural caffeine – no crashes, just pure
refreshment.",
    "visual_hint": "Close-up of drink and natural ingredients"
  },
  {
    "scene": 3,
    "voiceover": "Join the 5AMClub with GlowDrink. Sunrise workouts just got
easier!",
    "visual_hint": "Person jogging at sunrise, smiling"
  },
  {
    "scene": 4,
    "voiceover": "Feel the glow. GlowDrink – available now. #GlowUp",
    "visual_hint": "Product shot with branding and hashtag"
  }
]
```

Each segment will carry through the pipeline, triggering corresponding voiceover and visuals.

**Tools:** Primarily the LLM (OpenAI API). We might use **HuggingFace Transformers** as a backup (for instance, a fine-tuned GPT-2 or Llama for ad copy), but given the nuance needed, a top-tier model like GPT-4 is

preferred initially. Optionally, a **grammar/style checker** (like LanguageTool or a mini GPT prompt) can post-process the script to ensure correctness and that it fits time constraints (e.g., ~2.5 words spoken per second means ~40 words for a 15s ad; we can quickly word-count and trim if needed).

**Orchestration:** The script generation is typically done synchronously by the orchestrator (it's relatively quick). If multiple variants are needed, the orchestrator can spawn parallel tasks to multiple LLM calls (Celery tasks for each variant).

**Retries & Quality Control:** The agent will validate the script against basic criteria: - **Length check:** If the total word count is too high for the intended video duration, the agent can prompt the LLM to shorten the script (e.g. "Make it snappier and under 30 words total") and regenerate. - **Brand alignment:** Ensure the brand name or key message appears. If the script forgot a required element (say the CTA or a feature), the agent can append an instruction ("Include the discount offer in the last line") and regenerate that part. - **Multiple Variants:** We can leverage the LLM to get a few different takes (humorous vs. inspirational, etc.) and either present them to the user or automatically pick one. The selection could be random for A/B testing or based on a prediction of performance. For example, Pencil's platform generates multiple variants and scores them for predicted outcomes; we could similarly use an internal scoring (perhaps using our future performance data) to choose the most promising script. Initially, though, we might simply pick the first or let the user choose.

In case of **LLM failure** (e.g. API error or content filter issue), the agent will automatically retry with a simplified prompt or fallback to a smaller model. Ensuring the script is ready is a gating item – subsequent steps depend on having voiceover lines, so the system will make a couple of attempts if needed, rather than proceeding with an empty script.

## 3. Voiceover Generation (TTS)

**Purpose:** Convert the script lines into spoken audio. A convincing voiceover is crucial for an engaging ad, so we use text-to-speech to narrate the script in a human-like voice.

**Process:** For each script segment (scene/line), a TTS engine generates an audio clip. We prefer an open-source TTS like **Coqui TTS** (which supports many voices and can run locally) to avoid API costs. Coqui has pretrained models for English that produce natural-sounding speech. We'll choose a default voice appropriate for the brand/tone (e.g. a friendly female voice for a cheerful tone, or a deep male voice for a serious tone). The system can allow **voice selection** per brand – stored in the brand profile – or choose automatically based on context (like Creatify offers 140+ voices to match different needs [5] , we can start with a few). Initially, one high-quality general voice is used for simplicity.

The TTS module (running as a Celery task or invoked function) takes a text string and returns or saves an audio file (e.g. WAV/MP3). Each line's audio is generated separately so that we have fine control in editing (we can insert pauses between lines or align them with different video clips easily).

After generation, we measure each clip's duration (via FFmpeg's `ffprobe` or the TTS output length) and attach that metadata to the segment. This helps in video editing to know how long each scene should be on screen to accommodate the VO.

**Tools: Coqui TTS** library for local inference (fast on CPU for short sentences). Alternatively, for a quick start or more voice variety, we might use **Microsoft Azure Cognitive Services TTS** or **Google Cloud TTS** on a free tier – they have natural voices and multiple accents. But since cost-efficiency is key, Coqui (which is offline and free) is ideal if performance is acceptable. We also consider using **WhisperX** (an alignment tool) after generating the audio to get precise word-level timestamps, which can later help in subtitle syncing or cutting scenes. WhisperX can align the generated speech with the original text [5], though this might be optional if our TTS can output timings.

**Data Output:** Each scene now has an `audio_file_path` and a duration. For example:

```
{
  "scene": 2,
  "voiceover": "Packed with natural caffeine – no crashes, just pure
refreshment.",
  "audio_file": "voiceovers/123_ad_scene2.wav",
  "duration_sec": 3.5
}
```

These files are stored in a temporary storage (could be local disk or a cloud storage bucket). They will later be mixed into the video.

**Retries/Variants:** - If the TTS generation fails for a specific line (rare, but perhaps due to unusual proper nouns or length), the agent can retry with a different service or slightly modified text (e.g. if an acronym isn't spoken properly, we can insert spaces or phonetic spelling). For MVP, we likely use one TTS engine and count on it working; fallback could be an online API if our offline model struggles. - **Voice variation:** The agent could generate multiple voiceover variants if unsure which voice is best (e.g. a calm voice vs. an energetic voice) and then use performance data later to decide which type to prefer. Initially, we will simply pick one to avoid complexity, but the architecture supports switching voices easily (it's just a parameter for TTS model or API). - All audio clips are normalized for volume and quality. If needed, we'll run a normalization filter (ffmpeg loudnorm) to ensure consistent volume across clips and with the background music.

**Communication:** The orchestrator can fire off TTS tasks for all lines in parallel, since they are independent. This can speed up processing (generate all segments concurrently on separate CPU cores). Celery helps manage this fan-out. Once all voice clips are ready, the orchestrator moves to the next stage.

### 4. B-Roll Retrieval (Stock Video Selection)

**Purpose:** Find relevant visual footage (video clips or images) to pair with each segment of the script. These are the "B-roll" background visuals that will play while the voiceover is heard, illustrating the story.

**Process:** For each script segment, the system searches a stock media library for a short clip that matches the content or mood. We will integrate a free stock footage API such as **Pexels** or **Pixabay** (both offer free videos/images) to retrieve clips by keyword. The query is derived from the segment's voiceover text and/or the `visual_hint` from the script. For example, if the line talks about "morning workout", the query might

be "morning exercise" or "sunrise jogging". The agent might also include the product type or concept keyword to find something contextually appropriate (e.g. "energy drink workout").

We prioritize finding footage that is **royalty-free** and high-quality. Pexels API returns a list of videos with metadata (tags, length, download URL). The retrieval module will: 1. Send the search query to the API (with our API key, staying within free tier limits). 2. Get top N results (say 5). 3. Filter or rank them. We can use heuristics like: - Duration: prefer clips longer than the voiceover duration (so we can trim down, rather than run out of footage). - Content: If we have **OpenAI CLIP** or a similar image-text model locally, we can embed the candidate thumbnails and the query text to pick the closest semantic match [6] . For MVP, a simpler approach is to trust the API's relevance ordering or look at the title/tags from the API response for a match with our keywords. - Orientation: if we plan output in vertical format, maybe prefer vertical videos (Pexels does provide width/height; we can filter for aspect ratio). Alternatively, we'll handle cropping later, but it's good to consider.

1. Select the best match and download it (perhaps in 720p to save bandwidth, which is sufficient for preview/testing; we can choose resolution based on need).
2. Store the video file locally (or in a temp directory).

**Data:** We enrich our segment data with the chosen B-roll info:

```
{
  "scene": 2,
  "voiceover": "...",
  "audio_file": "...",
  "duration_sec": 3.5,
  "video_file": "broll/energy_drink_natural.mp4",
  "video_duration_sec": 10.0,
  "selected_start_sec": 2.0
}
```

We may also choose a specific sub-section of the clip to use. For instance, if the stock video is 10 seconds but we only need 3.5 seconds, we could decide to use seconds 2.0 to 5.5 of it (perhaps skipping an irrelevant part at the beginning). We could automate this by, say, starting at a point where the action is good or just at 0 for simplicity. Advanced use could involve analyzing frames (e.g. avoid a segment where nothing is happening by using CLIP to detect where content changes).

**Tools:** - **Stock APIs**: Pexels (Free with API key, millions of videos/images). Pixabay similar. We favor these to avoid cost. - **CLIP (Contrastive Language-Image Pretraining)**: A model from OpenAI (or LAION's CLIP) that can encode images and text to a vector space. We can use CLIP to improve matching: embed the voiceover text and candidate video frame to compute a similarity score. This can be done via HuggingFace's `openai/ clip-vit-base-patch32` model. Running CLIP on a handful of images is fast (even on CPU) since the data is small. This helps ensure the chosen footage truly aligns with the script semantics (e.g. if our query "morning exercise" returns some irrelevant result first, CLIP might catch that a later result is more semantically on-point). - Possibly **FFmpeg** for processing the video (we will definitely use FFmpeg in the editing stage to trim and scale clips, but at retrieval stage, we mostly just download and store).

**Retries & Variants:** - If no adequate video is found (e.g. an obscure query), the agent can **fallback strategies**: broaden the search query (remove specific brand terms, or use a more general keyword). It can also try an alternative source (if Pexels fails, maybe Pixabay). In worst case, for that scene we could use a stock image (and apply a Ken Burns zoom effect to simulate video) or a solid color background with text – ensuring the pipeline doesn't break. - The agent might retrieve multiple candidates for each scene and **choose the best** via CLIP or a simple score. We can even log the top 3 so that if the user or future versions want to try a different visual, they exist. - **Deduplication:** Ensure we don't pick the same stock clip for two different scenes if it would look repetitive (unless that's intentional). The retrieval can track used URLs and skip them for subsequent scenes.

- **Parallelization:** We can fire multiple search tasks in parallel (one per scene) via Celery. Stock API calls are IO-bound, so parallelism will improve speed, especially if the ad has 4-5 scenes. Each task downloads its clip. The orchestrator waits until all scenes have a video before moving on.

**Meta-data storage:** We maintain a library mapping of what content was used (e.g., store Pexels ID and attribution info if needed for later licensing checks). Pexels videos are free to use without attribution, but we might still log the source for our records or if we ever need to filter out certain creators.

## 5. Music & SFX Matching

**Purpose:** Add a background music track and any sound effects to enhance the ad's emotional impact and polish. Music sets the mood and pacing, while sound effects (SFX) can emphasize certain moments or transitions.

**Process (Music):** After script and concept, we determine the **mood** or energy level of the ad. This could be inferred from the concept's tone (e.g. "Upbeat, motivational" indicates a fast-paced, happy soundtrack). The agent can also analyze the script sentiment or ask GPT: *"What mood does this script convey? (e.g. upbeat, dramatic, calm)"*. Based on this, we select a music track from a small curated library of royalty-free music. Initially, to keep things simple and cost-free, we might bundle a handful of Creative Commons or public-domain music tracks covering a range of moods (e.g. an upbeat corporate track, a high-energy electronic track, a gentle inspirational track). These can be sourced from free libraries (like Free Music Archive or Pixabay Music).

Each track in our library would be tagged with mood descriptors (uplifting, energetic, calm, etc.). The agent picks the track whose tag matches the desired mood of the ad. For example: - **Upbeat/Energetic** – use "Upbeat_Uplifting_Corporate.mp3". - **Cinematic/Dramatic** – use "Epic_Orchestral.mp3". - **Calm/Trustworthy** – use "Soft_Ambient_Background.mp3". - **Trendy/Youthful** – maybe a track with a pop or EDM vibe.

If the brand has specific preferences (stored in their profile, like "Rock music" or "No vocals"), we incorporate that in the selection logic. The track is then loaded and will be mixed into the final video at a low volume under the voiceover.

One challenge is ensuring the music length >= video length. If the chosen track is shorter, we can loop it or repeat sections (many stock music tracks are loopable). If it's longer, we'll cut it to fit the video duration (with a nice fade-out at the end using ffmpeg audio filters).

**Process (SFX):** Sound effects are optional in v1, but we plan for them. SFX can accent transitions or important points: - A "whoosh" or "swipe" sound during a scene transition (to make cuts feel smoother). - A "ding" or sparkle when showing the product or a price (to draw attention). - Crowd cheer or drum roll for a big reveal, etc.

We'll maintain a small library of common SFX (these can be obtained from freesound.org or similar, and are small files). The agent will decide where to place SFX by analyzing the script and concept: - Transition sounds: by default, add a subtle whoosh at each scene change if it fits the style. - Emphasis: if the script has a moment like "50% off now!", we might put a *ding* right after that line. - For a playful ad, maybe a little "pop" sound when text appears.

For MVP, we might hardcode one or two uses: e.g., always use a soft whoosh during scene transitions to enhance them. It's a simple addition that increases perceived quality.

**Tools:** - No AI needed for music selection beyond simple rule-based matching (though we could use a model like **MusicLM** in future or an audio similarity model to dynamically pick music, but that's overkill initially). - We could use **HuggingFace's CLAP (Contrastive Language-Audio Pretraining)** in future to match audio to a mood description, but given our curated library approach, not required now. - **ffmpeg** will be used to adjust volumes and mix multiple audio tracks (voice, music, SFX) during rendering. - For SFX timing, we may just schedule them manually in the timeline (e.g., whoosh at scene boundary exactly when we cut or transition at time t, which we know from accumulating previous scene durations).

**Data:** The system will produce a `music_track` selection (file path) and possibly a list of SFX events:

```
{
  "music_track": "music/upbeat_corporate.mp3",
  "sfx_events": [
    {"file": "sfx/whoosh.wav", "time_sec": 3.5},
    {"file": "sfx/ding.wav", "time_sec": 12.0}
  ]
}
```

Where time_sec for SFX is when to play it in the final timeline (relative to video start). These will be used in editing.

**Variants:** - We may have multiple suitable music tracks per mood. The agent could randomly pick one or choose based on past feedback (if we learn that for a fitness product, a high-tempo EDM track performed better than a corporate pop track, next time we'll pick accordingly). - The architecture allows experimenting: e.g., generate two versions of the ad with different music to see which the client prefers or which performs better (A/B test). Initially, we stick to one to reduce complexity.

**Retries:** If for some reason the chosen music track doesn't feel right (maybe the agent can detect if the beats per minute don't match the cut pace), it could try an alternate track. However, such nuance may be done via manual tuning at first (we'll likely test a few and set a good default). The system design keeps the music separate, so swapping a track is trivial without rerunning other steps.

## 6. Editing & Video Rendering

**Purpose:** Compile all the assets (voiceover audio segments, B-roll video clips, background music, SFX, and text overlays) into the final video file. This is the most processing-intensive step, essentially performing the role of a video editor autonomously.

**Process:** The editing module receives everything needed to produce the timeline: - Ordered list of video clips with intended durations, - Corresponding voiceover audio files for each clip, - Background music track, - Subtitles (if any) and graphics (logo, etc.), - Transition instructions (e.g., use crossfade between clips).

The module uses **FFmpeg** (invoked via a Celery task) to stitch this together. We can break this down into sub-steps:

**a. Prepare Clips:** Each B-roll clip may need trimming to match its scene duration. For scene *i* with duration $d\_i$ (which is usually the voiceover length), if the raw clip is longer than $d\_i$, we trim it to $d\_i$ + *overlap* (a slight overlap if we plan crossfades). If it's significantly longer, we might choose the best segment of it (for now, likely we'll take the start). If a clip is shorter than the voiceover (shouldn't happen often if we filtered by duration, but if so) we have options: slow it down slightly (not ideal), loop it, or freeze the last frame. A simple solution: allow a little freeze-frame at the end of a too-short clip to cover the remaining audio. These adjustments will be done via FFmpeg trim filters.

We also ensure **aspect ratio and resolution** consistency. If our target output is (say) 1080x1920 portrait (to cover TikTok/Reels use-case), and a stock clip is 1920x1080 landscape, we have to rotate or crop. Most likely we'll do a center-crop to vertical: FFmpeg can crop 1080x1920 out of 1920x1080 by cutting off the sides (ensuring the main action is centered – often acceptable). Alternatively, we can scale the clip to fill height and blur the sides, but that's more complex for v1. We choose one approach and apply it to all clips for consistency. This can be done in the FFmpeg filter graph per clip (e.g., `crop=1080:1920` or using the `scale` and `pad` filters accordingly).

**b. Assemble Timeline:** We then concatenate the clips in order, adding transitions between them. We will utilize FFmpeg's `xfade` **filter** for a smooth crossfade (or other effect) between each pair of clips [4] . For instance, a 0.5-second crossfade where clip A fades out while clip B fades in gives a professional touch. Implementing a crossfade in FFmpeg is done by overlaying the end of clip A with the start of clip B and applying the filter [4] . We have to slightly extend clips beyond their voiceover to have material to crossfade; e.g., we might let each clip run 0.25s past the end of its voiceover to use for transition overlap.

If crossfades prove too tricky in the first implementation, a simpler fallback is **cut transitions** (direct cuts). The architecture allows turning transitions on or off via a setting. But we aim for at least basic fades since they elevate quality. (There's even an `ffmpeg-concat` tool with preset transitions, but we can manage via filter_complex ourselves.)

**c. Overlay Audio:** We mix the audio tracks: - **Voiceover**: We have multiple voiceover files, one per scene. We'll concatenate these in order to form one continuous voiceover track (or instruct FFmpeg to place each at the correct timeline position). Since there might be a fraction of a second gap we want between lines (for natural pacing or to accommodate a transition), we can insert a short silence between audio clips if needed. However, a benefit of crossfading video is we can actually let the audio flow continuously scene to scene with no silence, which sounds more fluid. - **Music**: The music track will be added to the timeline starting at

time 0 (unless we wanted to delay it). It will run under the entire video. We set its volume to a low level (e.g. -15 dB relative to full scale) so it doesn't overpower the voice. FFmpeg's `amix` or `sidechaincompress` filter can duck the music when voice plays, but a fixed lower volume is easier and usually fine. - **SFX**: For each SFX event, we will overlay its audio at the specified time. This is done by treating it as another input in FFmpeg and mixing it in at that timestamp (using the `adelay` filter to start it at the right time and then mix).

The result is a single mixed audio track for the final video.

**d. Subtitles & Text:** We want to add subtitles for the spoken lines (and possibly additional on-screen text like the brand name or a final call-to-action). We have the script text, so simplest is to use FFmpeg's **drawtext filter** or the **subtitles filter**: - The **subtitles filter** allows us to provide an .SRT or .ASS file with timestamps and it will burn those captions onto the video frames. We can generate an .SRT from our script easily: each voiceover line's start time and end time (we accumulate the start times as we build the timeline, since we know durations). For example:

```
1
00:00:00,000 --> 00:00:03,500
Mornings feel impossible? Boost your energy with GlowDrink.

2
00:00:03,500 --> 00:00:07,000
Packed with natural caffeine – no crashes, just pure refreshment.
```

and so on. Using the `subtitles` filter requires FFmpeg built with libass, and it will use default styling unless we supply .ASS for custom font/style. We can start with default (white text, bottom center). - The **drawtext filter** gives more control in-filter (we can set font, size, color, position, and even add background box or shadow for readability). We would need to add one drawtext filter per subtitle line with an enable condition matching the time range of that line. It's doable but a bit fiddly to compute exact timing in seconds. The advantage is we can use brand fonts or colors if needed by pointing drawtext to a font file and setting color codes.

For MVP, the simplest path: generate an .SRT and let FFmpeg burn it. It's one command and handles line wrapping automatically. We just ensure the text is not too long per line (the LLM script is already concise). We also might set a slightly larger font size for mobile viewing.

Additionally, **graphics overlays**: if the brand has a logo image and wants it visible, we can use FFmpeg's `overlay` filter to put a PNG logo at, say, top-left or bottom-right throughout, or just at the end for a final branding. The brand profile might specify this ("watermark logo in corner" or "end card with logo"). For now, we can plan an end card: the last scene could actually be a generated image (like a simple background with the brand logo and a call-to-action text). We can create this end slate image on the fly (for example, use PIL/ Pillow library to generate an image with brand's color and some text) and then feed it as the last "video clip" to FFmpeg. But if time is short, we may just overlay the logo on the last video scene.

**e. Rendering:** We invoke ffmpeg with all the filter directives compiled. This might look like:

```
ffmpeg -i clip1.mp4 -i clip2.mp4 -i clip3.mp4 -i voice.wav -i music.mp3 -i
whoosh.wav -filter_complex "[0:v]trim=0:3.5,setpts=PTS-STARTPTS[v0]; \
 [1:v]trim=0:4.0,setpts=PTS-STARTPTS[v1]; \
 [v0][v1]xfade=transition=fade:duration=0.5:offset=3.5[v01]; \
 [v01][2:v]xfade=transition=fade:duration=0.5:offset=7.5[outv]; \
 [3:a]adelay=0|0[vo]; [4:a]volume=0.2[bgm]; [5:a]adelay=3500|3500[sfx]; \
 [vo][bgm][sfx]amix=inputs=3:dropout_transition=0[outa]; \

[outv]subtitles=script.srt:force_style='Fontsize=24,PrimaryColour=&HFFFFFF&'[outv2]"
\
 -map "[outv2]" -map "[outa]" -c:v libx264 -c:a aac output.mp4
```

(The above is illustrative: it trims clips, crossfades them, mixes 3 audio sources, and burns subtitles – actual filter graph might differ, but that's the idea.)

FFmpeg then outputs `output.mp4`. The task will run on our server; for a ~15s video with a couple clips, this should only take a few seconds on a modern CPU. Even if it takes, say, 10-20 seconds, that's acceptable in an asynchronous job scenario.

**Output:** The final video file (e.g., `ad_12345.mp4`) is stored, and the path or URL is saved to the database. The API can then return this to the frontend or allow the user to download it.

**Orchestration:** We treat editing as a single Celery task that depends on all previous tasks' outputs (Celery Canvas makes it easy: we can have a chord that waits for all voiceovers and downloads to finish, then calls the editing task). The orchestrator provides the editing task with a manifest of all asset file paths and needed instructions (durations, transitions, etc.). Because video rendering can be CPU-intensive, we might dedicate a separate worker process or queue for it to avoid blocking other smaller tasks.

**Error Handling & Retries:** This is a complex step, but errors typically come from: - Missing files or mis-synced inputs. (We ensure in code that all files exist before calling FFmpeg.) - FFmpeg command syntax errors. (We will iterate during development to get it right; once working, it should be stable for similar inputs.) - Memory/CPU issues. (15s videos won't overflow memory; we keep resolution reasonable, e.g. 720p or 1080p.) - If FFmpeg does fail (non-zero exit), the task will report an error. We log the ffmpeg output for debugging. The agent can catch this and attempt a simplified render: e.g., if crossfade filter caused an issue, try a straight cut version (we could maintain a fallback rendering routine that just concatenates without fancy filters). - In worst case, the agent returns an error status for that job to the user, but our goal is to minimize that through testing common scenarios.

**Quality considerations:** We will experiment with bitrates to balance file size and quality (targeting social media, a 1080x1920 15s video at ~5Mbps is fine, resulting in a few MB file). We'll also ensure the audio mix is clear (voice not drowned by music).

This autonomous editing approach is a standout feature – it effectively replicates what a human editor might do in Adobe Premiere, but via code. Notably, it's more flexible than template-based generation and can incorporate any new effect or overlay by adjusting the FFmpeg instructions. It is "programmatic editing".

## 7. Subtitles & Captions

**Purpose:** Generate and include subtitles (closed captions) for the video to improve accessibility and engagement (many viewers watch ads on mute, so subtitles are crucial). Additionally, having the text on-screen reinforces the message.

We already covered implementation in the editing step, but let's separate the pipeline logic: - **Generation:** The subtitle text is basically the script. Since we have exact voiceover text, there's no transcription needed (we bypass the need for speech-to-text). However, if we want **word-level timestamps** for fancier caption styles (like highlighting each word as spoken), we could use **WhisperX**. WhisperX is an alignment tool that, given the audio and the transcript, returns timestamps for each word [5] . This could enable, for example, karaoke-style captions or popping each word individually. For v1, we likely won't go that far, but it's an available technique. More straightforward: we use known segment timings for line-level captions.

- **Formatting:** We want the subtitles to look good on short videos. Likely position them center-bottom with a semi-transparent background or shadow. If brand style is important, we might use the brand's font and a brand color for the text outline or highlight. These style choices can be configured. Initially, white text with black outline or shadow (a common readable choice over any video) is fine.

- **Multiple languages (future):** The architecture could allow translating the script via an LLM or translation API and generating subtitles in multiple languages. Not needed for MVP, but we mention it to highlight extensibility.

- **Closed vs Open captions:** We are burning in (open captions) because that guarantees they show up on platforms (and it allows stylization). We could also provide a separate .srt file as output if the user or platform prefers uploading it (some ad platforms accept SRT for auto-caption). It's trivial for us to save the SRT we generated for burning and include it.

- **Agent adjustments:** The agent ensures that each subtitle line isn't too long to read quickly. If a single voiceover sentence is long, we can split it into two subtitle lines or even break a scene into two if needed. Ideally, our script generation already gave short lines. We target no more than ~8-10 words per subtitle line for readability.

**Tools:** For alignment (if used): **OpenAI Whisper + WhisperX**. WhisperX can take our generated voice audio and align each word [5] , but since TTS gives perfectly known text, the only reason to run it would be to get precise per-phoneme timing which might be overkill. We likely skip WhisperX in MVP, but keep it in mind if we notice any drift in timing or for future word-level effects.

**Orchestration:** Subtitle creation happens just before rendering. The orchestrator will create the SRT file from the script and collected timings (the sum of durations up to each scene gives the start time). It then passes that file to FFmpeg. If FFmpeg's subtitle rendering fails (perhaps issues with font), we can adjust and retry (or try drawtext). But these are deterministic once set.

**Quality:** Subtitles will be reviewed for correctness (since it's AI-generated text, but we assume it's correct as it's our own script). If the brand or user edits the script text after generation (maybe we'll allow them to tweak wording before final render in the UI), the updated text is what gets captioned.

In summary, subtitles are an automated byproduct of our script – easy to generate and highly beneficial for ad performance.

## 8. Visual Effects & Transitions

**Purpose:** Enhance the video with visual effects (FX) and smooth transitions to increase professionalism and viewer retention. This includes things like transition animations, filters, overlays, and any dynamic visual element beyond raw cuts.

**Transitions:** We have planned crossfades as the default transition between scenes (a safe, neutral effect). But the system is capable of more: FFmpeg's `xfade` filter supports dozens of transition types (wipe, slide, pixelize, etc.) [7] [8] . We could allow the agent to choose a style based on the concept or trend – e.g., a flashy TikTok-style ad might benefit from quick **flash cuts** or whip-pan transitions, whereas a calm story might just dissolve. For MVP, we stick to one (crossfade or straight cut) globally, but we note the potential. The architecture could easily swap in a different `transition=<type>` in the FFmpeg command based on an input parameter.

**Visual FX:** - **Color Filters:** Possibly apply color grading or filters to match brand aesthetics (e.g., make all scenes slightly warmer tone if brand wants a sunny vibe, or add high contrast if trying to be edgy). This could be done via FFmpeg filters (`hue`, `eq` filters for brightness/contrast). - **Speed Changes:** Speed ramping (slow-mo or fast-mo) could highlight moments. For instance, a dramatic effect: slow motion during a key scene. We could achieve by adjusting video playback rate in FFmpeg. Unless conceptually needed, we skip it in v1. - **Zoom/Pan (Ken Burns):** If any scene uses a still image (in case video footage wasn't available and we fall back to an image), we can animate it by slowly zooming in or panning to avoid a static look. FFmpeg has a `zoompan` filter to create that effect. This will be our backup if we ever have to use a static image. - **Overlay Graphics:** We discussed logo overlay. Another common overlay is a **call-to-action text or price tag** appearing on screen. For example, "$19.99 for a limited time!" might pop up visually. We can either bake that into subtitles (if it's spoken) or as a separate text overlay graphic (if it's not spoken, e.g., sometimes you show text that isn't spoken word-for-word). The agent can decide if any extra on-screen text is needed. Perhaps based on brand preference, always show the website URL or slogan at the end. - **Animation/Stickers:** Not likely in v1, but conceptually we could add simple animations (like a slight bounce on the final logo, or an SVG graphic flying in). This would require more complex rendering logic (possibly using something like Motion Canvas or generating an overlay video). We probably won't do custom animations initially.

**Integration:** All these effects are implemented via FFmpeg filter graph modifications. The microkernel design shines here: adding a new effect doesn't require changing how other modules work. For example, if we later incorporate a **stable diffusion video filter** to stylize footage into a cartoon look, that could be an optional plugin at this stage. Or integrating an API like RunwayML for certain transitions. We simply would drop in that step if needed.

**Agent's role in FX:** The AI agent can choose to toggle certain effects based on: - Brand style: e.g., **Brand personalization** might specify "use our brand colors". This could mean adding a colored border or background shape in scenes, or colorizing subtitles. The agent will ensure that happens in the editing instructions. - Trend: If a TikTok trend involves e.g. a quick glitch effect between cuts, the agent could mimic that by instructing the editing module to use a particular filter (there are glitch filters in ffmpeg or we could

simulate by a rapid frame cut). - Content considerations: The agent won't add effects that clash (e.g., don't put a goofy transition on a serious PSA ad).

For MVP, the *presence* of transitions and subtitles is already a big plus over raw cuts. We will implement crossfades and subtitle overlays for sure, maybe a logo overlay. These are baseline FX.

**Testing FX:** We will test with sample outputs to ensure the transitions align with audio (we don't want voiceover from scene 1 still going while scene 2 visuals are fully on – unless that's intended overlap). Ideally, the crossfade happens just as the first voice line ends, so it's seamless. This timing will be fine-tuned by adjusting offsets.

In conclusion, our pipeline's editing stage not only "stitches" but also "polishes" the video with transitions and overlays. This modular approach to FX means we can continuously improve visual style without touching upstream logic – for instance, swapping out the transition effect library or adding an end-screen template is isolated to the editing component.

## AI Agent Framework and Autonomy

Rather than a rigid pipeline, ReelForge's backend is driven by an **AI Agent Orchestrator** – a hybrid of rule-based workflow and intelligent decision-making. This agent is the executive that decides *which* modules to invoke, *when* to invoke them, and *how* to adjust if something goes wrong or suboptimal. The design is informed by emerging patterns in autonomous AI systems, where a central agent uses modular tools to achieve high-level goals [9] .

**Hybrid Approach:** The agent combines **symbolic logic** (hard-coded rules/heuristics we program) with **AI reasoning**. For straightforward decisions (e.g., if no video is found, use an image), we use deterministic logic. For more complex decisions (e.g., evaluating the creativity of a script or the suitability of a concept), we can employ an LLM as part of the decision loop. This is similar to how AutoGPT or other agentic frameworks let an AI reflect and choose tools, but we tailor it to our use-case. The agent essentially treats each pipeline component as a **modular skill** it can deploy [10] [9] . Skills are like "concept brainstorming", "copywriting", "image search", etc., each implemented by our modules. The agent's job is to orchestrate them in the right sequence and conditions.

**Memory and State:** The agent maintains context about: - **Brand**: stored in a **Vector Database** (like Pinecone or Weaviate) as well as relational DB. For each brand, we store embeddings of their background info (product descriptions, prior ads, style guidelines). This forms a long-term memory. When the agent is working on a new ad for that brand, it can query this memory: e.g., retrieve similar past campaigns or brand facts to include. Pinecone allows semantic similarity search in milliseconds, so the agent might do a Pinecone query like "brand mission statement" or "past ad for similar product" to ground the new generation. - **History**: The agent records previous outputs and their performance. This is key for learning (more on feedback loop below). This could be as simple as "Ad ID 47: concept was X, CTR=2%" stored as a data point. The agent can embed the textual aspects (concept, script) and store with a vector labeled with performance. Later, when making a new ad, it might ask: *"What past ads for this brand (or similar brands) had highest conversion?"* and retrieve those to inform the current strategy. This moves towards **Case-Based Reasoning**, using Pinecone to find analogous situations. - **Conversation/Session**: If the user interacts (say chooses one of multiple concepts or provides feedback), the agent keeps that in short-term memory for the current job.

The agent's memory approach is influenced by techniques in chatbot development – e.g., using vector stores for long-term context [11] – but applied to marketing content. Each brand could effectively have a dedicated knowledge base the agent consults.

**Decision Logic:** The agent operates in a loop of Plan -> Execute -> Evaluate: 1. **Plan:** Given a request (e.g., "Generate a video ad for Brand X's new product Y"), the agent plans the tasks needed. Typically, that's the standard pipeline, but it could diverge. For example, if the brand just wants to repurpose an existing video with minor edits, the agent might skip concept and script generation and instead use a different pipeline (e.g., just do editing and captioning on provided footage). The agent decides this based on inputs (if user provided a script or assets, it modifies the plan).

1. **Execute (with Monitoring):** The agent invokes each module (tool) in turn, monitoring outputs. After each module, it can perform checks:
2. Did the module output meet expectations (e.g., script length check, TTS audio clarity, video relevance)?
3. It keeps **confidence scores** where possible. For instance, when multiple scripts or visuals are candidates, it might assign a confidence or use the model's probabilities.

4. If a confidence is low (like the concept generator wasn't confident, or CLIP similarity for a chosen video was mediocre), the agent might trigger a *retry or alternative tool*. E.g., low visual match confidence could trigger a secondary search query or using a different approach (like generating an image with DALL-E if stock fails).

5. **Evaluate/Reflect:** After a full pipeline run, the agent can evaluate the result. This could involve an automated critique: use GPT-4 to analyze the final script or concept for quality (LLM as a critic). Or use a heuristic like Pencil's approach of predictive scoring – Pencil's system predicts how well an ad might perform using a model trained on $1B ad spend. We may not have that data initially, but we could have a lightweight version (perhaps use OpenAI's **Moderation or classification models** to avoid any compliance issues and ensure positivity, etc.). The agent could also ensure no brand safety issues – e.g., check the script doesn't accidentally violate policy (like making unsubstantiated health claims, which an AI might do if not guided; we'd include guardrails in prompts, but this evaluation is a backstop).

If the evaluation identifies a problem or improvement, the agent can loop back. For instance, if the script is good but perhaps too long, the agent might trim it and regenerate voiceover (without having to restart everything manually – it's autonomous enough to iterate locally).

**Autonomous Improvement:** The agent is designed to **learn from feedback**. This is where the feedback loop (point 5) ties in: - We will feed **real performance data** (click-through rates, conversion rates, etc.) back into the system. The agent's logic includes using these outcomes to adjust future decisions. For example, if for Brand X, ads that used "humor" concept have consistently outperformed "serious" ones (based on past A/B tests or metrics), the agent will favor humor in the next concept. It might store a flag in brand memory like "humor_variant_success: true". - Over time, the agent can even do **prompt tuning**: subtly refine the prompts given to GPT for script generation by including additional context from what has worked. E.g., "Last time emphasizing 'free shipping' improved results, make sure to mention it." It can insert such hints automatically if the data shows that pattern.

We can implement a simple form of learning: a rules engine that updates weights for certain choices. For instance: - Maintain a score for each voice type per brand (if one variant ad used a female voice and got higher conversions, increment that voice's score; agent will pick it more often). - Maintain scores for concept styles (meme-style vs testimonial vs straightforward pitch). - These scores update with each campaign outcome. Next generation, the agent probabilistically biases towards higher-scoring traits. This is essentially a **multi-armed bandit** approach for creative choices.

Additionally, we plan for **automated A/B testing**: The agent can create two variants of an ad (maybe different final scene or different opening line) and suggest the user run both. With integrated performance tracking, it would then identify the winner and note that insight. This goes beyond what many competitors do – they might provide multiple creatives, but not automatically learn from the results in a closed loop. Our architecture's ability to loop in performance data means the agent can *truly optimize actual outcomes, not just predicted scores* [12] . This is a key differentiator: while Pencil uses a large dataset to predict winners, we aim to dynamically learn and **optimize based on real metrics**, making the system smarter over time for each specific user.

**Frameworks & Implementation:** We draw inspiration from **LangChain's agent tools** and Microsoft's **Semantic Kernel** for modular skill orchestration [13] . However, we keep our implementation lightweight to avoid bloat. Essentially, we have a Python class `ReelForgeAgent` with methods to call each module and some decision-making methods. The vector DB (Pinecone) is accessed through a thin wrapper: e.g., `agent.fetch_brand_memory(brand_id, query) -> documents`. We'll likely maintain some prompt templates for self-reflection, like *"Critique this script for effectiveness and brand alignment"* which we can feed to GPT-4 as needed.

**Pinecone/Weaviate usage:** Each time an ad is created, we can embed key aspects (the script text, etc.) and upsert to Pinecone with metadata including performance. When the agent starts a new job, it can do vector similarity queries like: - `similar_ads = pinecone.query(vector=embed(new_concept), filter={"brand": brand_id}, top_k=3)` This might return past ads that have similar messaging or were done for similar campaigns. If those have performance data attached, the agent can see "these similar ones had poor conversion, so maybe this concept needs tweaking" or vice versa. This is advanced behavior and would evolve as we gather data, but the infrastructure will be there from v1 to collect and store embeddings.

**Safety & Guardrails:** The agent also includes a **Governance** aspect. It needs to ensure outputs are within acceptable content guidelines (no profanity, no discriminatory or false claims, etc.). We will implement a content filter (OpenAI's moderation API or a local word blacklist) to check the script and maybe even analyze the final video (e.g., ensure the images chosen aren't offensive or against policy). The agent will adjust if any issue: e.g., if moderation flags a word, the agent can automatically swap it or re-generate that part. This is important for brand trust.

In summary, the AI Agent framework turns what could be a static pipeline into a **flexible, adaptive system** that not only automates the steps but intelligently navigates them, improving with experience. It treats each component as a tool in a toolbox, assembling them in various ways to meet the end goal of a high-performing ad [9] . This gives ReelForge a significant edge in **adaptability** – as new AI models or data come in, the agent can integrate them as new skills without re-architecting the whole system.

# Attribution Tracking and Performance Feedback Loop

A core principle of ReelForge is **continuous improvement via real-world feedback**. To enable this, the system must track each generated ad's performance in the wild and feed that data back into the AI agent. We establish an attribution framework that uniquely identifies every video and ties it to analytics data.

**Ad ID Injections:** Every generated video ad gets a unique ID in our database. We embed this ID in the campaign's tracking parameters (UTM codes) or within the ad metadata: - For ads that will be used on platforms like Facebook, Google, TikTok, etc., the typical approach is through **UTM parameters** on the destination URL. For example, if the ad's call-to-action is a link to the brand's site, we append `?utm_source=ReelForge&utm_campaign=Ad_<ADID>` (and other utm tags like medium, etc.). This `ADID` (or a human-friendly variant) ensures that when a user clicks and lands on the site, analytics platforms record which creative drove it. - Even for platforms that are not link-out (like Instagram Story ads often still have a swipe link, so UTMs apply, or for TikTok maybe linking to a profile or so – in those cases, we might rely on platform analytics more directly rather than UTMs). - If an ad is used in a context without links (say an organic social post generated by our tool), tracking is trickier. In those cases, we rely on the platform's reporting of views/engagement by post ID. Since our system created the content, we know the content ID on the platform if we publish via API. But MVP might not publish; more likely the user downloads the video and posts it themselves. In that scenario, we may have to rely on the user to tell us how it did, or implement computer vision to identify our video (not ideal). So initial focus: track performance for ads used in paid campaigns or link-equipped contexts.

**Integration with Ad Platforms:** We plan to connect to APIs of major ad platforms: - **Facebook Ads API (Meta Marketing API):** Once the user uses our video in a Facebook ad campaign, if they connect their ad account to us (which Pencil and AdCreative also do), we can fetch metrics for that creative. Specifically, the API can provide stats at the ad level (impressions, clicks, spend, conversions, etc.). We would need the ad's ID from Facebook – we can either require the user to input that or if we integrate deeply, we create the ad via API which returns the ID. Pencil directly integrates to launch ads and even swap creatives. For v1, maybe we won't auto-launch, but we want at least read access to performance. - **Google Ads API:** Similarly, fetch metrics for video ads (YouTube or GDN). Google's API is complex, possibly we might integrate with Google Analytics or GA4 as a simpler proxy (e.g., see traffic by UTM campaign). - **TikTok Ads API:** TikTok provides an API for creative and campaign metrics as well. Many SMBs may not use TikTok API though, but we should design for it. - **Shopify / E-commerce integration:** Because many DTC brands use Shopify, we can ingest conversion events (sales) from their store. For example, using Shopify's API or webhooks, we can get order data including referrer or UTM. If an order has our `utm_campaign=Ad_123`, we attribute that sale to our Ad ID 123. This closes the loop all the way to revenue, not just clicks. We'll set up a small ETL pipeline to gather these events. Possibly a service like Segment or directly hitting each API with scheduled tasks.

**Data Pipeline:** We'll create a scheduled Celery task (or use Celery Beat) to periodically fetch new performance data. For instance, every night, fetch the previous day's metrics for all active ads from connected APIs. Or if using webhooks (Shopify allows webhooks on order creation), we can get near-real-time conversion info.

We'll have an **Analytics Database** (could just be some tables in PostgreSQL): - Table `ad_performance` with columns: ad_id, date, impressions, clicks, spend, conversions, revenue, etc., updated daily. - Table `ad_events` or `conversions` for individual conversion events if we want granular (each sale with value,

ad_id). - Potentially a `campaign` table if multiple ads per campaign, but since we focus on creative, ad-level is fine.

**Real-time aspect:** If an ad is performing extremely poorly, theoretically the agent could intervene (e.g., alert the user or automatically generate a new variant). That's more phase 2, but our system can certainly be extended to do that. For now, real-time ETL ensures that whenever the user checks our dashboard or whenever the agent is tasked with a new generation, the latest results are there.

**Closing the Loop:** The agent consumes this data in its learning module. For each ad ID, after some run time, we'll have metrics like CTR (click-through rate), CVR (conversion rate), etc. We can compute a **performance score** for each ad. For example, a simple scoring function:

```
score = 0.5*(CTR / avg_CTR) + 0.5*(CVR / avg_CVR)
```

(for that brand or industry baseline) – to gauge if it's above average. Or more straightforward: maybe label ads as "Hit" if conversion rate > X or ROAS > 1, etc. These labels and scores get associated with the ad's features (concept, script text, visuals, etc.). The agent then uses this as a training signal: - It might do a **vector store update**: when storing the embedded script in Pinecone, also store the performance score in metadata. Then we can query, weighting similar scripts by their score. - It might update brand-specific preferences: e.g., in Brand X's record, store "best_hook_style: motivational" if we notice motivational ones consistently got better results than humorous.

**Feedback Ingestion Example:** Suppose ReelForge generated 5 ads for Brand X over 2 months. Some had a young casual tone, others a formal tone (maybe we tried different angles). The performance data shows the casual-tone ones had 2x higher engagement. The agent's logic will then favor casual tone in subsequent generations for Brand X. Concretely, we might have a parameter in brand profile `preferred_tone` that is updated automatically based on what worked, which then is used in prompt injection for script generation ("Use a casual, friendly tone because that resonates with our audience").

Another example: If one particular **scene structure** worked (say ads that immediately showed the product in first 2 seconds had better retention), the agent can learn that pattern and adjust future outputs to do the same (this insight could come from analyzing video metrics like 3-second view rate if available).

**Cross-user learning:** In early stage, data per brand might be limited. We can also glean general insights by aggregating across brands (respecting privacy and without sharing proprietary info). For instance, maybe across all e-commerce apparel brands, those that used a human model in the video got higher CTR than those that just showed product. Our system could identify that trend by tagging whether a human was present in footage (we could know from video search keywords or using a vision model). The agent could then bias towards including humans in future apparel brand videos. This is analogous to how Pencil trained on a large dataset to predict winners, but we will be doing it in a more heuristic/on-the-fly manner at first, then possibly train our own predictive model when data grows.

**Interface & Dashboard:** We will expose these analytics to the user in a clear way (likely a simple dashboard showing each video and key metrics – though building a full frontend for that can be phase 2; initially, maybe a manual report or just the fact that our agent uses it). The important part for the manifesto is the backend's capability: we store and utilize these metrics internally.

**Validation of Loop:** This closed-loop approach is a **major advantage** over static creative generators. Competitors like Pencil emphasize predictive scores before launch, but they also integrate with FB to swap underperformers. ReelForge will actually learn and **adapt content generation itself** from actual outcomes, which is even more powerful. We position this as **"outcome-driven optimization"** [12] – the system gets smarter with every ad served and every click recorded, tuning its creative strategies in the direction of what the market responds to.

In sum, the attribution and feedback subsystem ensures that ReelForge is not generating ads in a vacuum – it closes the feedback loop from creation to deployment to result, fueling an ever-improving creative agent.

## Brand Personalization Engine

One of ReelForge's pillars is **brand-specific personalization**. Every brand has unique values, style, and visual identity, and our system is built to reflect that in the generated ads. To do this, we maintain per-brand memory and customization settings, and weave those into each step of the pipeline.

**Brand Profiles:** In our PostgreSQL DB, we have a `brands` table capturing key info: - **Basic info:** name, industry, target customer profiles. - **Brand voice/tone:** e.g., playful, professional, edgy, compassionate. (This can be provided by the user or inferred from their website copy via an initial analysis – e.g., run their homepage text through an LLM to summarize tone.) - **Key messaging points:** common slogans, value propositions, do's and don'ts (e.g., "never mention competitor names", "highlight our 30-day guarantee"). - **Style guidelines:** - Visual: brand colors (in hex), preferred fonts (we might allow them to upload a TTF or choose from Google Fonts), logo image, any specific design motifs. - Audio: perhaps a preferred voice style or background music genre if any. - Pacing/editing: if the brand has known preferences (e.g., "our ads should be fast-paced and energetic" vs "steady and informative"). - **Brand assets:** any media they uploaded – logos, product images, previous ads. We can use these in videos (e.g., overlay logo, or include product photos as cutaways if no stock footage matches a product).

**Vector Memory:** We use Pinecone (or Weaviate) to store richer brand data in vector form: - We embed the brand's own marketing content: like their About Us page text, product descriptions, prior ad copy. These become retrievable chunks. - When generating scripts, we can do a similarity search in this vector DB to pull in any relevant snippets. For example, if the script is about a specific product feature, the agent might fetch how the brand's website describes that feature and then include wording from it to stay on-brand. - Also, any *historical interactions* or preferences can be stored. If the brand's team liked a certain phrasing previously, we keep that as part of memory.

**Custom Prompt Injection:** At runtime, when the agent calls the script generator or concept generator, it injects brand-specific context. For instance: - Prepend to the prompt: "Brand voice: {brand_voice_description}. Brand tagline: '{tagline}'. The script should adhere to this style and language." - If the brand has example ads or copy, we could even feed a few examples in-context (few-shot learning). E.g., include a past successful tagline as an example. - If the brand has a particular formatting (maybe they like to always start ads with a question, etc.), we incorporate that instruction. - This ensures the output isn't generic but feels tailor-made for the brand. (For example, if a brand avoids slang, the prompt injection will explicitly say "Do not use slang or internet abbreviations.")

**Template Variations:** We plan to support different **ad templates** or frameworks at a conceptual level. For example: - **Problem-Agitate-Solution**: a classic formula. - **Testimonial**: where the script is like a customer talking. - **Demonstration**: straightforward show and tell of product features. - **Storytelling**: a mini narrative.

The agent can decide which template to use for a given ad. It might base this on: - **Trend & Concept**: If a certain style is trending (e.g., "TikTok storytime format"), it picks that. - **Brand**: If the brand is very formal, maybe a testimonial or direct format is better than a meme-y one. - **Past performance**: If we tried two types before and one did better, lean on that.

These templates are implemented as prompt skeletons or different pipeline flows. For now, we might explicitly have a few prompt templates for script generation and just choose between them. Over time, we could allow the agent to even blend them or create new ones.

**Example:** If the brand sells an AI SaaS B2B product, a humorous meme-style ad might not suit. The brand profile might indicate "professional, explanatory tone". The agent thus would choose a straightforward template ("Highlight problem and solution with product screenshots and calm voiceover") rather than a goofy skit. Conversely, a fun DTC brand for teens could do a TikTok dance meme, which the agent could plan for by using a template that includes an on-screen person (maybe in future using an avatar, etc.). Our competitor Creatify, for example, offers AI avatars that can act as spokespersons for a brand [14] [15] . In our system, we could incorporate such a capability later as a plugin (e.g., if brand wants an avatar, that could be an alternative to stock footage for scenes where someone speaks to camera).

**Optional Fine-tuning on Brand Assets:** For ultimate personalization, we consider fine-tuning models: - **Fine-tune LLM on brand tone/data:** We could take the brand's existing copy (website, brochures, emails) and fine-tune a smaller language model (like a 7B parameter model) to speak in that style. That model could then be used for script generation for that brand exclusively. This requires some data and training time, but for a paying client it might be worth it. As v1, we likely won't train new models due to resource constraints, but we might use prompt-based fine-tuning (like GPT-3's few-shot or OpenAI's embedding-based customization). - **Fine-tune TTS voice:** If a brand has a signature voice (say the founder's voice or a specific actor's voice) and provides samples, we could fine-tune a TTS model to clone that voice. Tools like Coqui or ElevenLabs allow voice cloning from a few minutes of audio. This could be huge for brand consistency (their ads always have the same narrator voice). It's an advanced feature – probably not day-1, but our architecture (modular TTS) can accommodate swapping the model with a fine-tuned one for a specific brand. - **Custom Vision Models:** If a brand has a unique visual style (e.g., a fashion brand with a particular aesthetic), one could train a GAN or diffusion model on their catalog to generate on-brand imagery. That's far beyond MVP and not necessary as long as stock footage suffices, but conceptually possible in our plug-in system (an image generation module conditioned on brand style).

For now, brand customization in v1 will rely on **data-driven prompt engineering and style settings rather than model re-training**. That covers a lot of ground with minimal effort: - Use their colors & logo in the final video (visual consistency). - Use their preferred language and terminology in the script (verbal consistency). - Maintain consistent tone in voiceover (auditory consistency, by picking the same voice each time unless changed intentionally).

**Brand Library**: We'll also provide a "Brand Library" interface (likely in the frontend) for the user to upload and input these preferences. Our backend just needs to store and utilize them. This is akin to what Pencil offers with its brand library feature [16] , which they emphasize is critical to on-brand outputs.

**Memory Example:** Suppose brand "HealthyCo" always capitalizes the name of their product "VitaDrink" and uses the phrase "Feel the vitality". We store that. The script generator prompt injection will ensure the LLM knows to use the word "vitality" and capital "VitaDrink". If the raw LLM output said "vitadrink" in lowercase or missed their slogan, the agent will catch it and fix it (either by instructing the LLM, or a simple post-processing).

Another example: brand style says no exclamation points (some brands find them off-tone). Our agent can quickly scan and remove or reduce them in the final script text. These are small but important details for brand trust.

**Personalized Analytics:** The performance feedback will also be filtered by brand. The agent learns what works **for that brand**, not just generally. A playful ad might work for one brand but flop for another; our feedback loop is keyed by brand so we adjust appropriately per brand persona.

In summary, the **Brand Personalization Engine** ensures *everything* the system creates feels like it was made by someone on the brand's internal team who deeply understands their identity. Technically, this is achieved by injecting brand-specific data at every possible point: prompts, visual choices, and output styling. This level of personalization is something generic competitors can miss if they focus on one-size-fits-all generation. By having a memory and profile per brand, ReelForge can produce content that not only performs but also that brand managers feel comfortable using (it "sounds like us" – which is often a make-or-break factor for adopting AI-generated content).

## Deployment Strategy and Cost Efficiency

We design v1 deployment to be **lean and cost-effective**, leveraging free tiers and lightweight infrastructure suitable for a two-person startup. The system will initially run on Docker containers via **Docker Compose**, which is simple and sufficient for an MVP.

**Environment Setup:** - We create a Docker Compose file defining services: - `web` : our FastAPI app (this could also run the Celery worker in the same container for simplicity, or we might split into `api` and `worker` if scaling or resource isolation is needed). - `redis` : using the official Redis image, acting as the message broker for Celery and a cache. We can point it to a cloud Redis (like Upstash) to avoid local state, but during dev or for local deployment Compose is fine. - `postgres` : a PostgreSQL container for local dev. In production, we'll use a managed DB instead of a container, but Compose can have a placeholder. - Possibly a `vector_db` if self-hosting (Weaviate or a simple service). However, Pinecone doesn't have a self-host option (it's SaaS). For MVP, Pinecone's free tier (up to a decent number of vectors) is probably easiest: we just hit their API from the `web` container. Alternatively, Weaviate has a docker image but it might require more RAM. Given budget, we might stick with Pinecone free SaaS for now (very generous for small scale). - All these are configured with environment variables for secrets (OpenAI API keys, Pinecone keys, etc.), which we keep out of the repo for security, using Docker secrets or an `.env` file not committed.

**Cloud Deployment:** We can host this on a single VM or a platform like Heroku/AWS Lightsail/DigitalOcean: - A single **$40-$80/month** droplet or EC2 instance can easily run our containers (especially if using CPU for all tasks). We can start with something like 4 vCPUs, 8GB RAM. This should handle a modest load (each video generation might use 1 CPU core for ffmpeg and some for AI tasks). - If we have credits on AWS or GCP, we could use those. AWS for example: use ECS with Fargate or a small EC2 with Docker Compose in a screen session. But Compose on a single node is fine to start. - Using **Supabase** (hosted Postgres) on free tier for the DB: typically free tier can handle up to some amount of rows/storage which likely suffices at first (say 500MB or 1GB of data, which for text and logs is plenty). - Using **Upstash Redis** free: they allow a certain number of connections and data (like 10k commands/day free, should be okay initially). - Domain and SSL: maybe we use Vercel/Netlify for frontend, but our backend can use a free SSL from Let's Encrypt on our server, or if on Heroku, they manage SSL. These are minor or free.

**Cost Breakdown Aim:** Keep monthly infra < $200: - VM/Compute: ~$50 (assuming some credit or minimal usage). - DB: ~$0 (free tier) initially, maybe $15 if upgrade for more storage later. - Redis: ~$0 (free upstash) or $10 for a basic instance if needed. - External API costs: - OpenAI API: depends on usage. GPT-3.5 is $0.002/1K tokens. Each script gen maybe 100 tokens => $0.0002 each. If we generate 1000 scripts, that's $0.20 – negligible. GPT-4 is pricier ($0.03/1K prompt, $0.06/1K completion). If we use it sparingly (maybe for concept only, or final quality improvement), monthly could be <$50. We can mostly use GPT-3.5 to control cost and fall back to GPT-4 for tricky cases or maybe have the user on a higher plan unlock GPT-4. - Pinecone: free tier covers plenty of vector operations for dev (up to 1M vector dimensions stored and moderate queries). If we exceed that, the next tier is ~$0.07 per 1000 queries etc., which at low usage is still cheap (and might be covered by some credits). - Stock media APIs: Pexels is free and unlimited with key as long as usage is reasonable. No cost. - Coqui TTS: open source, runs on our machine, no cost. It might use CPU heavily during generation, but not a monetary cost. - If we use any other Hugging Face models or similar, we'll run them locally (so just CPU time). - Monitoring/CI: we might use free GitHub Actions for CI (they give a certain amount of minutes free, should be fine for our pushes) and something like UptimeRobot for ping (free for a few monitors). Possibly Sentry for error tracking (also has a small free tier). - So indeed, under $200 is achievable, likely under $100 until usage grows.

**Scaling Plan:** As usage grows or if we need better resilience: - We'd move to **Kubernetes** (perhaps on a managed K8s like GKE/AKS or DigitalOcean Kubernetes) once we need to scale horizontally. But K8s adds ops overhead for a small team – we'll likely postpone until necessary. We explicitly choose Docker Compose in MVP to avoid devops complexity. - We can gradually modularize into separate services if needed: e.g., a dedicated video rendering service (maybe on a machine with better CPU/GPU), a separate service for the AI agent if needed. But initially, one process can handle the pipeline sequentially (and Celery can spawn threads or processes as needed). - The microservice/microkernel separation is logical in code, but physically we might still run them in one container process for simplicity (just calling functions). Only heavy tasks go to Celery background threads.

**GPU usage:** Currently, everything can run on CPU, albeit slower. If we later incorporate heavier models (like a large local LLM or real-time stable diffusion video, etc.), we might need a GPU. We could then use a cloud GPU instance on demand (e.g., spin up only when needed, or use a service like RunPod for specific tasks). But none of our MVP components strictly require a GPU: - TTS Coqui can run CPU (fast enough for short sentences). - CLIP is lightweight on CPU for a few images. - FFmpeg uses CPU (unless we enable GPU acceleration, which isn't necessary for short videos). - GPT API is on OpenAI's servers. - So we intentionally avoid needing our own GPU to keep costs down.

**Managed Services Choices:** We lean on managed solutions where they save us time/cost: - **Supabase** for Postgres gives us a free DB plus nice web UI for DB inspection and authentication if needed. - **Upstash** for Redis eliminates running a Redis server 24/7 ourselves (which might be overkill for small usage). - Using these also means less memory/CPU used on our box and easy scalability (Supabase can scale the db as we pay, etc.).

**Development vs Production parity:** We'll use Docker so that the environment is consistent. Locally, we develop maybe without containers for fast iteration, but we ensure everything runs in containers for production.

**CI/CD:** We will utilize **GitHub Actions** to automate tests and deploy: - On push to main, run unit tests (if any) and maybe a quick integration test (could simulate one full pipeline on dummy data). - If tests pass, build the Docker image and push to a registry (GitHub Container Registry or Docker Hub). - Then deploy: since we may just have one server, this could be a simple SSH and pull or using Docker Compose on the server. We can automate with something like: - Use an Action to SSH into the server and run `docker compose pull && docker compose up -d` to deploy the new version. This can be done securely via GitHub Actions with an SSH key. - If using a platform like Heroku, we could just push code (but Heroku might not handle heavy ffmpeg well, plus limited free tier now). - If using AWS, maybe use ECS and have the Action update the task definition with the new image. - We'll choose a simple path first – maybe a single EC2 or Lightsail box.

**Monitoring & Logging:** - **Logging:** All services will log to stdout, which Docker captures. We can aggregate logs easily. For more robust solution, maybe set up a lightweight ELK stack or use a hosted log service (Datadog or Logflare have free plans). But to avoid cost, maybe just use Papertrail free plan (50MB/month) or log to files on disk and rotate. - **Error Tracking:** Use Sentry's free developer plan to catch exceptions in the FastAPI and Celery processes. This way we get alerts if something crashes for a user. - **Uptime monitoring:** UptimeRobot can ping our health endpoint every 5 minutes free. We'll implement a `/health` route that checks that the DB and Redis are reachable and returns OK. That gives us basic uptime info and alerts if the server is down. - **Performance monitoring:** Not crucial at first, but we could monitor resource usage with something like **Netdata** (open-source, easy to run as a container, with dashboards for CPU/mem/disk/network). This helps ensure e.g. ffmpeg isn't eating all memory or that we aren't swapping. We can set up Netdata on the same box for internal use – it's free.

**Security:** - Secrets (API keys for OpenAI, etc.) will be stored in environment variables, not in code. We'll use GitHub Actions Secrets store for CI, and on the server, we'll have them in a .env file that Docker Compose reads (excluded from repo). - The FastAPI endpoints should be secured – since the Next.js frontend and our backend are separate, we should have an authentication mechanism (maybe a simple API key or token for the MVP). Possibly we can integrate with NextAuth on the frontend or use Supabase Auth. But given MVP and maybe single-user testing, we might not implement full OAuth, maybe just a basic auth or a secret key required in headers. We do want to ensure others can't hit our API and abuse it (especially since it could use our OpenAI credit). - The system handles user data (brand assets, etc.) so we'll enforce HTTPS, and store data securely. We'll also be mindful of not exposing any PII in logs. - Container security: We base our images on official Python slim images, etc., to reduce vulnerabilities. We can run a scanner (like Trivy) occasionally.

**Cost Estimation:** Tallying up expected monthly: - Cloud VM: $50 (assuming we have to pay; if we have AWS credits we might be at $0 for a while). - DB: free (Supabase free tier). - Redis: free (Upstash). - Pinecone: free (Starter plan). - Domain: $0 (maybe we use a .vercel.app or something for MVP or a $10 domain/year negligible). - APIs: maybe $20-$50 if usage picks up (OpenAI mostly). - So perhaps ~$50 without credits,

which is well under $200. We budget up to $200 to account for scaling or if we decide to use some paid services.

As we onboard more users, costs will rise (more API calls, maybe need bigger DB instance, etc.), but presumably by then either revenue or funding could cover it. The architecture is designed to allow swapping free options with paid scalable ones easily: - If Outgrowing Supabase free, upgrade to their paid or switch to AWS RDS. - If Upstash free is limiting (they have a cap on commands), either upgrade (their paid is still cheap) or run our own Redis on the box if it's not too loaded. - If Pinecone free runs out (I think free is 1 index, 1 pod, which might be fine for quite some time since our vector count per brand won't be huge at first), either pay for Pinecone or host Weaviate on our box (if we have memory headroom). - We prefer to buy managed services when necessary because with 2 engineers, focus should be on product, not maintaining infra.

**DevOps Work**: One of us can set up the Docker and CI in a couple of days; after that, deploying new changes is trivial, which enables rapid iteration.

**Later Scaling (K8s)**: If we ever have many concurrent users and heavy jobs: - We'd containerize each worker type (maybe a separate deployment for the ffmpeg rendering if it needs more CPU, one for the web API, etc.) and use Kubernetes to schedule across maybe a cluster of VMs. - Use HPA (Horizontal Pod Autoscaler) to spin up more workers on demand (like if many videos are queued). - But K8s overhead might be too much for 2 devs at start, so likely we'll go as far as we can with simpler methods (maybe just manually scaling by increasing VM size or adding another worker process).

**Backup/Recovery:** Use managed DB's backup (Supabase does nightly backups on paid, but on free we might rely on us making dumps occasionally). Since a lot of our generated content is ephemeral (videos can be regenerated), losing some data isn't catastrophic except brand profiles and performance data – we'll ensure those are backed up. Could set up a cron to dump DB to cloud storage weekly.

All in all, this deployment setup aligns with a garage startup ethos: maximize free resources, avoid unnecessary complexity, and ensure we can quickly ship improvements. We consciously avoid heavy platforms or enterprise tools until usage demands (no need for a full Kubernetes or expensive SaaS monitoring right away – we'll use just what's needed). This keeps costs low and focus on development. As usage grows, we have a clear path to scale up using the microservices-friendly architecture (the microkernel design means we *could* containerize each module and scale independently if needed).

# Competitive Landscape and ReelForge's Edge

The automated ad creation space is heating up, with notable players like **Pencil**, **AdCreative.ai**, and **Creatify** offering their takes. It's vital to understand what they do and how ReelForge's architecture gives us an edge in **flexibility, feedback integration, and performance optimization**.

### Pencil (AI-Powered Ad Creative Platform):

**What Pencil Does:** Pencil (founded 2018, now part of The Brandtech Group) generates ad variants (static and short video) using generative AI and heavily focuses on *predicting* performance. Their platform ingests a brand's assets (images, logos, past copy) and creates multiple ad versions, each scored by a predictive

model for conversion or ROAS. Notably, Pencil's predictions are powered by a large dataset (~5 years, $1B+ in ad spend across 4000 brands) used to train their AI. They integrate with Facebook/Instagram to directly launch ads and automatically swap out underperformers, closing the loop from creation to deployment.

They position as helping make ads that "actually work" by combining human creativity with AI speed [17]. Unique features include an **"Explainable" score** for each ad (to instill confidence in the predicted performance) and a **brand library** to ensure outputs meet brand guidelines [16]. Pencil often provides an editor for final tweaks and highlights that they augment creatives, not replace them [18].

**Pencil's Strengths:** - Data-driven predictions (their USP): They claim ~84% accuracy in predicting winning ads [19], which impresses data-driven marketers. - Integrated workflow: connecting to ad accounts for launching and getting results, plus an insights dashboard for creative analytics. - Focus on video storytelling: They do generate short video ads with decent templates and allow basic testing/feedback on those [20]. - Affordable entry for SMBs ($14/mo starter plans) which lowers adoption friction.

**Where ReelForge Wins:** - **True Learning vs. Prediction:** Pencil's core strategy is predicting performance upfront using a static trained model [21]. That's powerful, but it's inherently based on past average data and needs validation by real outcomes. ReelForge emphasizes *actual outcome-driven optimization* – our agent learns from each specific brand's real metrics and adapts accordingly [12]. Instead of relying solely on a global model's prediction, we tailor to the brand's audience behavior. This can outperform a generic predictor, making ReelForge "even smarter in practice than prediction-based systems" [12]. Essentially, we use feedback loops to personalize performance optimization, whereas Pencil primarily uses a one-size-fits-all predictive model. - **Flexibility & Creativity:** Pencil's outputs, while varied, can be somewhat *formulaic* – often recombining existing assets in templates [22]. ReelForge's architecture (GPT-powered scripting + dynamic stock footage) can produce more diverse and fresh creative angles [22]. Our agent can truly brainstorm new concepts (even trending meme angles) beyond the formula Pencil might stick to. This agility in trying unconventional ideas could find wins where Pencil's data-trained model doesn't venture. - **Multi-Platform Adaptability:** Pencil is deeply tied to Meta (Facebook/IG) environment – great for those channels, but not proven on newer platforms like TikTok or YouTube Shorts [23]. ReelForge from the outset is platform-agnostic. Our video outputs can be easily formatted to any aspect ratio and our agent can optimize for various channels. For example, we can generate a 9:16 TikTok style or a 1:1 for Instagram feed as needed. If we offer equal ease across platforms (Reels, TikTok, Shorts, etc.), we exploit a gap where Pencil might be less focused [23]. - **Autonomy & Speed of Iteration:** With only two engineers, our system is engineered for rapid iteration. We can quickly integrate a new model or trend as a module. Pencil is a bigger team now under a holding company, targeting enterprise – potentially slower to experiment with edgy trends. We have a nimble "microkernel" that can plug-n-play new capabilities (say a new diffusion model or the latest GPT) faster in response to market changes.

In short, Pencil is a strong incumbent with predictive analytics and integration, but ReelForge counters with *adaptive learning*, *creative flexibility*, and *multi-platform agility*. We also target outcome-optimization, positioning as the smarter system that actually improves as it gets data, rather than relying on static predictions.

## AdCreative.ai (Automated Ad & Creative Generation):

**What AdCreative Does:** AdCreative.ai (Paris-based, recently acquired by Appier in 2025) focuses on generating **ad graphics and copy** (and has added some video capabilities). They cater to producing ads at

scale, boasting over 3 million users and 1 billion creatives generated – indicating a very broad adoption especially among small businesses. The platform creates ads (images, short videos, text) aimed at high conversion, and provides a **"Creative Insight Score"** that predicts how well an ad might perform (similar in spirit to Pencil's score). Users can connect ad accounts (Facebook, Google, etc.) to have performance data fed back into AdCreative's dashboard, helping identify which creatives work best.

AdCreative's strengths are **scale and versatility**: they support multi-format (static, dynamic, video) and have a massive stock image library (100+ million images via iStock/Unsplash) [24]. They emphasize saving time (generating dozens of variants in minutes) and *empowering small budgets with AI-optimized creatives* [25]. They also integrate analytics and recommendations (like telling users which ad to use based on predictive insights) [26] [27]. Their recent acquisition by Appier signals a push into a larger marketing suite.

**AdCreative's Strengths:** - **High Throughput Automation:** Great for pumping out many ad variations quickly (useful for A/B or for populating many sizes/formats). They tout generation 70% faster than manual design [25]. - **Multi-Format Support:** Unlike Pencil, they weren't limited to video – they do banners, carousels, etc., and now even some video support [28]. This one-stop solution appeals to users needing ads across platforms. - **Analytics & Recommendations:** They have real-time analytics, presumably pulling in performance and suggesting optimizations [26] [27]. The system might say "Ad A has 50% higher CTR than average, use similar elements" – it guides users. - **Templates and Industry Tailoring:** They offer templates and insights tailored to specific industries/sectors [28], which helps non-experts get relevant designs. - **Scale of user base:** A large community means they're battle-tested (but also indicates their output might be more generalized to cater to so many).

**Where ReelForge Wins:** - **Quality over Quantity (Smart Targeting):** AdCreative is about scale and speed – cranking out lots of decent creatives [29]. ReelForge focuses on generating *the right creative*. Our AI agent doesn't just churn variations blindly; it uses the brand's data and feedback to produce fewer, **but more on-target** ads. We emphasize that our AI "intelligently creates winning videos *tailored* from sales insights" rather than flooding you with variants [29]. This can be a differentiator: instead of giving 10 average options, we aim to give 2-3 excellent, data-informed options. - **Short-Form Video Mastery:** AdCreative started in images and only recently added short video. ReelForge is video-first. We can likely produce *higher-quality videos* with real editing (ffmpeg) versus AdCreative's possibly template-based videos. Also, our pipeline involving actual footage, dynamic voiceovers, etc., yields more engaging content than perhaps AdCreative's simpler animated templates. Creatify's blog noted AdCreative primarily does images and simple dynamic ads, not high-quality video like Creatify does [30] [31]. So there's room for ReelForge to claim the high ground in video ad quality/flexibility. - **Feedback Loop Integration:** AdCreative does pull performance data for insights, but do they *close the loop* into generation? It appears they use it to inform the user (analytics dashboards) rather than automatically tweaking the generation. ReelForge's architecture is built to incorporate performance signals directly into the creative process (the agent learning). This means over time, a ReelForge user's outputs become increasingly fine-tuned to their audience, whereas AdCreative provides analytics but may rely on the user to make the creative adjustments. We effectively offer a *self-optimizing* creative agent, which is next-level compared to a tool that simply generates and reports. - **Flexibility & Extensibility:** AdCreative has a lot of integrations and features, but that can also mean complexity and perhaps rigidity. Our microkernel approach means we can adapt to new trends quickly. For example, if a new social platform emerges or a new type of ad format, we can plug that into our pipeline quickly. AdCreative being a bigger platform might not pivot as fast. We also aren't limited by a fixed library of templates; our generative approach can create novel combinations (like using any stock video available, not just pre-approved ones). - **Focus on Effectiveness:** We position ourselves as not just generating *ads* but

generating *ads that win*. We tie everything to actual metrics (like ROAS, conversions). AdCreative's marketing does mention conversion focus [26] [27] , so we both value it, but we take it a step further by automating the optimization via the agent.

Additionally, after the Appier acquisition, AdCreative might be moving upmarket or integrating into a bigger suite (Appier does AI ad optimization as well). There might be a gap for scrappy small teams who want a dedicated creative tool that's more personalized and not part of a large enterprise stack. ReelForge can fill that by catering intimately to each brand.

**Summary:** AdCreative is a volume-oriented, broad tool. ReelForge counters with **depth and personalization** – we don't just pump out ads, we learn what your brand's customers want and hone in. We provide fewer but *better* choices and reduce the noise. Our continuous improvement loop means the longer you use ReelForge, the better it gets for you, which is a powerful value proposition over a tool that doesn't learn your nuances.

## Creatify (AI Video Ads with Avatars and Generative Video):

**What Creatify Does:** Creatify (launched ~2023) is focused on **short video ads at scale**, often touting the ability to create hundreds of video ad variations from a single product URL [32] . Their standout feature is the use of **AI avatars** – digital actors that can serve as spokespersons in ads [33] . A user can pick from a library of avatars or even create a custom one (a "digital twin" of themselves from a few photos) [34] . These avatars can speak with 140+ realistic TTS voices in many languages [5] , allowing personalized spokesperson ads without hiring actors.

Creatify's workflow: the user enters a product page URL, Creatify scrapes the page for info and media, and automatically generates a video ad (script, visuals, voiceover) [35] . It basically automates the entire creative process from minimal input. They emphasize **UGC-style ads** (user-generated content style) – meaning their ads often look like selfie videos or influencer promotions made by their avatars [36] . They position it as a fast alternative to hiring creators for UGC ads, claiming high-converting results quickly [36] .

They also incorporate **generative video styles** – "AI shorts" feature that can produce videos in comic, 3D, or other styles using generative models [37] . They are experimenting with trendy visual styles for TikTok (e.g., turning images into animated videos).

**Creatify's Strengths: - Human-like AI Avatars:** This is a unique offering. Multi-cultural avatars with lip-sync tech give a "real spokesperson" vibe to ads without filming [33] . It's quite attractive for brands that benefit from a human presence (many do, as people connect with faces). - **One-click ad generation from URL:** Extremely simple UX for SMBs – no need to input much, just URL and choose styles, making ad creation accessible to non-creative folks. - **Mass Variation & Testing:** They allow creating hundreds of variations (different avatars, different hooks) easily [32] , encouraging rich A/B testing. They even have a **batch mode** to spit out many videos at once [38] . - **Focus on Video (especially mobile vertical):** They're all-in on video ads for social, which is the key growth area. They tailor templates to trending platform styles [39] . - **Traction with SMBs:** 10k+ paying customers by early 2025 and growing [40] , indicates product-market fit among small businesses wanting affordable video ads. Their $49/mo pricing is in the sweet spot for many SMBs [41] . - **Avatar novelty:** Being able to see yourself or a custom persona in an ad is a selling point (e.g., a realtor could have their AI twin speak in an ad, which Creatify specifically markets to local professionals) [42] .

**Where ReelForge Wins: - Dynamic Footage vs. Avatar-Only:** Creatify's avatar approach is powerful, but not every ad is best served by a talking head. Some products benefit from dynamic B-roll, real environment shots, etc. ReelForge's use of *actual stock footage* for B-roll can create more **visually engaging** ads in cases where showing context or multiple scenes is better than a single person talking. We can, for example, show the product in use, show happy customers, scenic shots, etc., which an avatar-only ad might not cover. Essentially, we have more flexibility in visual style – not just a spokesperson format. - **No Uncanny Valley/ Trust Issues:** While avatars are improving, some viewers might still sense they're AI and find it gimmicky or less trustworthy. ReelForge uses real footage and real human voiceovers (even if TTS, they sound natural) without leaning on an entirely virtual person. This avoids potential uncanny valley effects. We can also incorporate *real* human stock videos, which might appear more authentic than a synthesized avatar. For brands that prefer a polished montage style or product-focus rather than an influencer-style, our approach is more suitable. - **Brand Customization:** Creatify can make a custom avatar of the business owner, which is cool, but beyond that, how much brand-specific tailoring do they do? They largely auto-generate from the URL, which might lead to quite templated outputs (lots of ads with the same structure, just different text and avatars). ReelForge, with its brand memory and custom prompts, can produce more **brand-distinct** content. We actively incorporate brand style (colors, fonts, messaging). Creatify's output style might be somewhat homogeneous given it's template-driven (lots of those avatar-on-plain-background scenes). Our microkernel allows any style the brand wants – e.g., if brand wants an animated explainer, we could integrate Doodly or if they want live-action vibe, we use stock videos, etc. - **Feedback-Driven Improvement:** Creatify generates many variants but does it automatically learn which variant was best? It's not clear they close the loop; they likely rely on the user to pick the winner or run ads to see results. ReelForge's feedback loop would automatically glean what's working (e.g., maybe one avatar or one message consistently wins) and then emphasize that. This saves the user from guesswork in subsequent campaigns. - **Agent Orchestration & Extensibility:** Our agent could actually incorporate an avatar module too if we wanted (e.g., integrate Synthesia or D-ID API for talking avatars). The architecture is not mutually exclusive with that – it's a plugin. So anything standout from Creatify, we can add as a component. For instance, if down the line we see a need for spokesperson videos, we can plug in an "Avatar Generation" module: feed script to an avatar API, get video, and treat it as another B-roll option. Because of our modular design, we can adopt competitor features more easily than they can adopt our entire orchestration approach. - **Better Multi-Modal Combos:** Creatify's strength is avatars + scraped images. But maybe they don't combine multiple footage types (e.g., mix avatar speaking + cut to product B-roll + back to avatar). ReelForge can do that kind of multi-scene storytelling – we can have a human narrator voice (TTS) and show various visuals supporting the narrative. It's more like a mini-commercial. Creatify often ends up with one avatar doing the whole ad, which might get monotonic. Our ads can have *variety within the 15s* (various shots, etc., akin to professional ad structure), which could be more engaging.

Think of it this way: Creatify is automating the "influencer makes an ad" format. ReelForge is automating the "creative agency produces a polished spot" format. Both have markets; we lean toward the latter scenario which appeals to brands wanting more refined ads or those in niches where a talking avatar might not be enough (e.g., outdoor tourism might need scenic shots, etc.).

**On Feedback Integration:** none of the competitors (Pencil, AdCreative, Creatify) explicitly advertise an agent that self-optimizes content using performance feedback – they all provide data and the *means* to test, but the creative generation itself doesn't automatically evolve per brand. ReelForge aims to be the first where the creative gets *measurably better each cycle* because of learning, essentially acting as a creative director that gets to know the brand's audience intimately.

**Competitor Summary Table (for clarity):**

- *Pencil:* Best for predictive scoring and quick video creatives with data backing. Weak on multi-channel, very Meta-focused. ReelForge advantage: learning from actual outcomes, adapting content per brand, more flexibility in creative direction.
- *AdCreative.ai:* Best for massive scale, many formats, and quick wins for generic use. Weak on deep customization and video sophistication. ReelForge advantage: generating higher-quality, highly tailored videos with a focus on actual performance improvement, not just volume.
- *Creatify:* Best for instant video ads with AI avatars, great for UGC-style. Weak on varied formats (mostly spokesperson videos) and possibly authenticity concerns. ReelForge advantage: can produce more diverse video styles (not just one format), incorporate brand's unique style, and no uncanny valley issues; plus, we learn and optimize rather than just output en masse.

By understanding their approaches, we built ReelForge's system to *beat them on flexibility and intelligence*. Our microkernel agent can theoretically emulate their features (predictive scoring, high-volume output, avatar inclusion) by adding modules, but their systems would find it hard to emulate ours (an agent that truly learns and is modular by design) without a complete overhaul. **This architectural edge is our long-term moat** – we can keep integrating the latest AI advancements (be it a new model or data source) faster than more monolithic competitors, and our system becomes smarter with usage, creating a compounding advantage.

# Step-by-Step MVP Build Plan (First 30–60 Days)

Building ReelForge v1 is an ambitious but achievable project for a two-engineer team. We'll follow a phased approach over the next ~8 weeks to incrementally develop, integrate, and test each component. Rapid iteration and learning are key – we'll prioritize getting a basic end-to-end pipeline working early, then refine each part. Below is the breakdown with estimated timeline and key tasks, as well as suggested resources to guide implementation:

## Week 1–2: Core Setup and Architecture

- **Project Scaffolding:** Set up the Git repository and basic file structure. Create the FastAPI app with a simple health endpoint to verify deployment. Define a Pydantic data model for an Ad Job (to track inputs and status).
- **Task Queue Initialization:** Integrate Celery with Redis. Write a trivial example task (e.g., add two numbers) to ensure Celery worker communication works. This step validates our async orchestration setup.
- **Database Schema:** Initialize PostgreSQL (local or Supabase). Define schemas/tables:
- `brands` (id, name, settings JSON for style),
- `ads` (id, brand_id, status, creation_time, etc.),
- maybe `ad_components` (ad_id, component outputs like script text for debugging).
- **Basic Orchestrator Logic:** Implement a stub function for the agent that, given a brand_id and maybe a prompt, will sequentially call placeholder functions for each pipeline stage. At this point, these functions can just log something or return dummy data (e.g., script_gen returns a hardcoded script, TTS returns a path to a sample audio file, etc.). The goal is to have the control flow in place.
- **API Endpoint:** Create an endpoint `/generate_ad` (POST) that triggers the orchestrator (possibly via a Celery task) and returns a job ID. Also an endpoint like `/ad_status/{id}` to fetch status/

progress. Initially, these can work synchronously for simplicity (e.g., just call orchestrator directly and block, returning the final result for quick testing).

- **Docker Compose config:** Write a docker-compose.yml with services for web, redis, and db. Ensure `docker compose up` brings up the API and it can handle a request.
- **Testing & Iteration:** By end of week 2, we should be able to hit an API endpoint locally that runs through the stub pipeline and returns (even if the "video" is just a dummy file). This confirms our skeleton.

*Learning resources:* - FastAPI official tutorial (for setting up Pydantic models, dependency injection). - Celery documentation on first steps and using `celery.delay` or Celery chords for workflow (the Medium article "Mastering Celery" gives good examples of chains/groups) [43] [44] . - Supabase quickstart for setting up a Postgres DB and connecting (if we go that route). - Docker docs on Compose to ensure all services can communicate (especially FastAPI to Redis, etc.).

## Week 3–4: Implementing Pipeline MVP (Script -> Voice -> Video)

- **Script Generation (basic):** Integrate with OpenAI API. Set up an account and API key (possibly use OpenAI's free trial credits). Write a function `generate_script(brand, concept)` that calls `openai.ChatCompletion` with a prompt. For now, we can hardcode the prompt format. We'll test it with a simple example: e.g., brand = "Acme Coffee", concept = "highlight morning routine", and see if GPT-3.5 returns a reasonable short script. We'll iterate prompt wording until we get the structure we want (like numbered list of lines).
- **Testing:** Use the OpenAI Playground or quick Curl to refine the prompt before coding. Then integrate and actually log the output.
- Use `temperature` parameter to control variability.
- **Voiceover with Coqui TTS:** Install Coqui TTS and get a simple model running (they have CLI or Python interface). Try out generating audio from a test sentence and saving to file. Select a model like `tts_models/en/ljspeech/vits` (LJ Speech female voice – widely used). Integrate this in a function `synthesize_voice(text)`. We may need to install a few dependencies (PyTorch, etc.). Keep an eye on performance – generating one sentence is usually a few seconds on CPU.
- **Alternate quick win:** If Coqui setup is time-consuming, use gTTS (Google Text-to-Speech library) for a placeholder – it's just an API call and outputs an MP3 with a robotic voice. Not production quality, but fine to demonstrate pipeline. We can swap to Coqui in a later iteration when ready.
- **Video Assembly with ffmpeg (basic):**
- Initially, don't worry about dynamic B-roll. Use a placeholder video or image to simulate visuals. For example, use ffmpeg to create a blank background video of a certain length or take a static image and pan it.
- Combine the voiceover audio with this visual. For MVP, maybe just show a blank video with the voiceover overlaid (which proves we can produce an mp4 with audio).
- Write a function `render_video(segments, music=None)` that:
  - Takes the list of scenes with their audio files and durations.
  - For each, if we don't yet have real clips, perhaps use a stock filler (like generate a colored screen or a generic stock clip). Even a single static image stretched to the duration can suffice for a test.
  - Use ffmpeg (via subprocess) to concatenate these segments and overlay the audio. For a first attempt, we might not do transitions – just butt join.

- E.g., create a text file `inputs.txt` for ffmpeg concat demuxer listing segment videos, then run `ffmpeg -f concat -safe 0 -i inputs.txt -i voiceover_all.wav -c:v libx264 -c:a aac -shortest output.mp4`.
- If comfortable, directly use filter_complex as earlier, but that might take more time. The demuxer approach is simpler for initial output.

- Ensure ffmpeg is accessible (install it in container or host). Test by manually running a known command outside code, then integrate.
- **Stock Video Integration (simple):** By end of week 4, try calling the Pexels API for one scene:
- Sign up for Pexels API (free). Hardcode a query like "morning coffee" and see if we get a URL. Download it with `requests`.
- This is mostly to ensure we know how to use their API. We may not fully integrate multi-scene search yet, but get the basics.
- Also, we can try using an image search if video is too large to download often. But likely we can handle small video clips.
- **Run End-to-End:** Put it all together:
- Take a sample brand (we can hardcode brand info in absence of UI input),
- Run concept generation (if implemented; or skip concept and directly do script),
- Generate script via GPT,
- TTS each line,
- Compose a dummy visual for each line (colored screens or repeat one stock video multiple times),
- Render with ffmpeg,
- Return the output video path.
- Evaluate the result manually – we should get an MP4 we can play: does it have voice audio? Are scenes roughly aligned with lines (if we did multiple)?
- This is a critical milestone: **First Working Ad** (maybe rudimentary). It demonstrates the pipeline feasibility.

*Learning resources:* - OpenAI API docs (chat completion format) and perhaps LangChain quickstart if needed (though probably not needed for just one call). - Coqui TTS documentation on usage (their GitHub readme has examples of synthesize). - FFmpeg docs or StackOverflow for concatenating videos and adding audio. E.g., FFmpeg Wiki on Concatenation [45] or using the `-filter_complex "[0:a] [1:a]concat=n=...:v=0:a=1[outa]"` for audio. - Pexels API docs (simple REST calls).

## Week 5: Refinement & Feature Completion

- **Refine Script & Concept**: Add the **Concept Extraction** step before script:
- Could use GPT to output a one-line concept and trending angle. If short on time, we might skip a dedicated trend fetch and just incorporate a generic instruction in script prompt like "optionally include a trending meme if appropriate". But ideally, implement a `generate_concept(brand)` that maybe calls GPT with a prompt like "Brainstorm a creative video ad concept for {brand}. Provide a theme and hook that's popular with {brand's audience or trending online}." Use the result to steer script.
- **Integrate B-Roll selection fully**:
- For each script line, form a search query. A simple approach: take the main noun and verb. Or just use the whole sentence or the `visual_hint` if our script included that.
- Call Pexels API, get result, download that clip to a temp file.
- Trim the clip to needed length (we can use ffmpeg `-t duration` to cut it).

- Do this for each scene. If any fail, fall back to a default video or image.
- This step will involve dealing with video files, which can be large. We might use lower resolution to save time (Pexels provides different video files, maybe choose one with width ~720 for speed).
- **Music & SFX integration**:
- Add a selection of 2-3 royalty-free music files to our project (perhaps find some CC0 music from Free Music Archive). Or use Pexels Music API similarly.
- Implement a simple logic: if concept/tone contains words like "happy, upbeat" choose an upbeat track, if "relax" choose a chill track, etc.
- Mix the chosen music into the final ffmpeg command: e.g., use `-i music.mp3` and then in filter_complex `[music]volume=0.2[m]; [voice][m]amix=...`.
- Add a basic SFX on transitions: find a free "whoosh" sound, and overlay it at each scene boundary (we know times when one scene ends and next begins, so schedule it accordingly).
- This might be tricky to coordinate precisely – if short on time, maybe just put one whoosh at mid-video as a placeholder for demonstration.
- **Subtitles overlay**:
- Create an .SRT from our script lines and the known timings (the sum durations as we go).
- Use ffmpeg's subtitles filter to burn it in. Ensure we have libass in our ffmpeg build (if using a stock ffmpeg binary it usually does).
- As an alternative, use drawtext for one line to test – but SRT is easier if it works.
- Check the output to see if text appears correctly.
- **Orchestrator & Error Handling**:
- By now, the pipeline is complex. We wrap each stage in try/except so that if something fails, we mark the job failed but gracefully (maybe store error message).
- Possibly implement Celery chord: e.g., parallelize TTS for lines and maybe parallelize Pexels downloads (though might not be needed if sequence is fine).
- Ensure the main thread polls or waits for completion properly. Possibly have the API request trigger a Celery workflow and immediately return job_id.
- **Basic UI for Testing** (optional but helpful):
- Perhaps whip up a minimal Next.js or simple HTML form to hit our API. Or use a tool like Postman. But a small webpage where we input brand and see the video output in browser would be cool.
- However, given time, focusing on the backend is priority; we can test by manually checking the output files for now.
- **Intermediate Demo:** At the end of week 5, we'd aim to produce a *real* sample video: e.g. a 15-second ad for a fictitious brand using all modules (script by GPT, voiceover by TTS, several stock video clips stitched, background music, subtitles, transitions). It may not be perfect, but should showcase the full capability.
- We'll pick a scenario and run it. Evaluate where the quality is lacking (maybe voice volume vs music, or a weird stock footage mismatch) and note those.

*Learning resources:* - CLIP usage if needed: maybe skip if Pexels results are okay. But if needed, look at HuggingFace transformers pipeline for CLIP feature extraction. - ffmpeg advanced use: Transitions via `xfade` (the Superuser link we found shows example) [46] . - YouTube or blog examples for adding subtitles via ffmpeg (to confirm syntax). - Celery chords documentation: using `group` and `chain` for parallel tasks and then joining. - Possibly LangChain documentation for how to implement a multi-step agent (not strictly needed, but reading how LangChain tools work could inspire how to structure our agent logic in code).

## Week 6: Testing, Tuning, and Polish

- **Testing with Multiple Brands/Scenarios:** Try the pipeline with a few different inputs:
- A product (e.g., "Organic Tea") vs. a service brand (maybe "ACME Insurance") to see if the content adapts.
- Adjust brand voice settings (simulate one brand as fun, another as formal) and see if we can manually tweak the prompt injection to reflect that.
- Fix any glaring issues discovered: e.g., if GPT sometimes gives too long a script, add a token limit or instruction ("keep it under 40 words").
- If voiceover syncing is off, adjust by adding tiny silence padding or aligning via Whisper if needed.
- **Performance & Reliability:** See how long a full generation takes. Perhaps ~30s to 1min is okay for a 15s video on CPU. If something is too slow (TTS can be a bit slow per line), consider optimization:
- Use multiprocessing for TTS lines.
- If GPT call is slow, it's network-bound but usually a couple seconds.
- FFmpeg for 15s video should be a few seconds.
- So maybe total ~10-20s. That's acceptable. If more, we might consider using faster models (maybe faster TTS or using multi-thread).
- **Costs Check:** estimate how many OpenAI tokens we use per run (so we can plan costs). E.g., prompt ~100 tokens, output ~60 tokens -> $0.0003 per run with gpt-3.5. Negligible. So we're fine.
- **Implement Feedback Data Hooks:** We likely won't have time to integrate real FB/Google APIs now. But we can design the tables for metrics and maybe simulate:
- Add a field in `ads` for `utm_id` or similar that we would inject.
- Maybe have a dummy endpoint `/report_performance` that we can manually call with ad_id and some metrics to simulate receiving data.
- The agent isn't immediately using this feedback in generation for v1, but we lay groundwork. Possibly just print a statement like "If CTR < 1%, consider generating alternative concept" for demonstration.
- The key is to show in the manifesto that we have thought of it, but coding it fully might be beyond 60 days. Minimal: ensure we store the UTMs and prepared to collect data.
- **Brand Memory Integration (light):** If time permits, plug Pinecone:
- Take some brand-related text (like brand description) and upsert to Pinecone index.
- When generating script, do a Pinecone query for e.g. "unique selling point" and see if it finds something to include.
- This might be overkill for MVP, but we could do a simple example where brand's tagline is stored and we prompt the LLM with it. Actually easier: just store in DB and fetch normally. Pinecone is more for unstructured doc Q&A which might not be needed for short scripts.
- Possibly skip this for now due to time, or just mention that the brand description is passed into the prompt (which covers the same idea).
- **Documentation & Resources:** As we code, keep notes on how each piece was implemented (this will help writing our manifesto or technical docs). Also note any open-source licenses for assets if we included them (music, etc.) to ensure compliance.
- **Prepare Demo for Investors/Users:** By end of week 6, we want a stable version to demo. That means:
- The API doesn't crash on normal inputs.
- The output video is reasonably good (maybe not perfect polish, but showcases all features).
- We have one or two example outputs saved to show (in case live generation is risky).

- Possibly integrate front-end: even a command-line or basic web form to input brand name, product, and hit generate (we can do it via a simple HTML + JS calling our API).
- **Team Review:** The two of us review which parts of the system are most fragile or need improvement if we had more time. We'll prioritize those in the next steps or mark as future work. For instance, maybe the transitions are rough or the voice sounds monotonic – we note these.

**Learning Resources per Component Recap:**

- **LangChain & Agents:** "LangChain: Agents" documentation for understanding tool use (we didn't necessarily use LangChain, but reading it gave ideas on orchestrator design) [13] .
- **OpenAI Prompting:** OpenAI Cookbook on best practices for prompt design (to improve our script output if needed).
- **ffmpeg:** Official ffmpeg Wiki and OTTVerse tutorial for transitions [47] .
- **Celery Workflow:** Toptal's article on Celery for complex workflows which might give tips on chords and error handling [48] .
- **Coqui TTS:** Coqui forum or docs for optimizing voice quality and speed.

Throughout these 6 weeks, we allocate time for **learning and troubleshooting**. Integrating each AI component often involves some fiddling (install issues, API quirks). We should buffer some extra time for that (which we did by spreading pipeline over multiple weeks).

By the end of 60 days, we aim to have an MVP that can: - Accept a brand/product input (through an API or simple UI), - Automatically generate a fully edited 10-30s video ad (with voiceover, images/videos, music, subtitles), - Do it in one go with no human intervention in the loop, - And have logging/monitoring in place to track each step for debugging.

This MVP might not have all the bells and whistles (e.g., fully implemented continuous learning on its own yet, or a fancy frontend), but it will **demonstrate end-to-end functionality** and prove our architectural choices. From there, subsequent weeks can be spent on adding the more advanced features (feedback loop integration, brand fine-tuning, multi-language, etc.) and hardening the system for production use.

## Additional Insights & Best Practices

In building and deploying ReelForge v1, there are a few extra considerations to ensure a robust, secure, and maintainable product, especially for a small startup team:

- **Content Safety & Brand Reputation:** We must prevent the AI from producing anything that could harm a brand's image. This means implementing basic **content moderation** on AI outputs. For example, after script generation, run the text through OpenAI's Moderation API or a simple profanity filter to catch disallowed content. The agent should have guardrails (in prompts we already instruct the AI to avoid certain claims or tones). But as a fail-safe, if the script contains something like hate speech, or even just off-brand tone (too snarky when brand is formal), the system should flag it. Possibly, we set up a manual review trigger in such cases (at least in MVP, developers review any flagged output before it goes to end user). As we iterate, we can tighten the AI's prompt or fine-tune to the brand voice to avoid this.
- **Licensing and Attribution:** Since we're using stock footage and music, we must ensure they are properly licensed for our use. Pexels content is free for commercial use with no attribution required

– still, a good practice is to keep records of source to prove legality if needed. If we ever use a clip that needs attribution or purchase, we must incorporate that. For MVP, we stick to free assets (Pexels/Pixabay CC0 content). If a client supplies assets (like their own photos), we need to handle those files securely and separate from others.

- **Scalability of Video Storage:** Videos can be large, and generating many will require storage. For MVP we can store on disk, but as soon as it's viable, consider uploading final videos to a cloud storage (like S3 or Cloudinary) and just storing a URL in our DB. That way the frontend can stream from CDN and we don't overload our server bandwidth. This also allows easy purge/expire of older videos to save space. Under $200 budget, maybe we rely on free tier of Cloudinary (which offers some free video storage/transform). Or since output videos are relatively short, storing a few on disk is fine early on, but it's a point to monitor.
- **Privacy & Security:** If we're handling any personal data (like if a brand's info includes personal names or if we ingest their sales data for feedback), comply with basic privacy principles. Likely not huge for MVP, but worth noting as we ingest Shopify or ad account data, ensure data is only used internally for model improvement and not exposed. Secure those API credentials (OAuth tokens for FB/Google, etc.) carefully in our DB (encrypted if possible).
- **Continuous Integration & Delivery:** Set up a GitHub Actions pipeline early (as planned in week 1-2) so every code push runs tests. Even if tests are few initially, having the pipeline catches things like lint issues or dependency conflicts before deploy. We should also enforce code style (maybe use black or flake8 via a pre-commit hook) to keep the codebase clean as both engineers contribute.
- **Continuous** Learning **vs. Model Drift:** One caution: when our agent learns from its own output feedback, we must ensure it's not reinforcing bad habits or noise (like overfitting to small sample of results). This is a common issue in feedback loops (e.g., if it generates 5 ads, one does best by chance, and then agent always does that style, possibly missing better options). To mitigate:
- Use A/B tests and ensure enough data before concluding. The agent might require a certain confidence or number of samples to say "Style A > Style B" rather than one instance.
- Keep some diversity or exploration in generation to avoid converging too early on one style (multi-armed bandit approach with some epsilon for exploration).
- Continuously validate that performance actually improves; if not, adjust the learning strategy.
- (This is more future concern when we implement full closed-loop optimization.)
- **Monitoring and Alerts:** Since we're in production with minimal team, set up alerts for critical failures:
- If Celery queue builds up (jobs failing and retrying or stuck), we should know. Could integrate with a tool like Sentry (for exception alerts) or even a custom alert if queue length > X.
- UptimeRobot will alert if API is down. Also, have an alert if disk space is running low (videos could fill disk if not managed).
- Use a simple analytics like Prometheus or even just log events to a file – e.g., how long each stage takes – so we can identify bottlenecks and anomalies.
- **Graceful Degradation:** If a module fails, the system should ideally still produce an output, even if suboptimal:
- If script generation fails (OpenAI outage), maybe use a fallback model (like a smaller local model or even template a generic ad "Check out our product!" to not fully break).
- If TTS fails, maybe fall back to a default voice or a simpler engine (e.g., gTTS).
- If video retrieval fails for a scene, use a placeholder background with text overlay of the script (so the ad still has something visual).
- This way, the platform can deliver something rather than nothing, and mark it needs review. As a startup, maintaining trust means not delivering blank or error outputs if we can avoid it.

- **Modular Code Structure:** Keep each component's code encapsulated: e.g., a `scripts.py` for script gen logic, `media.py` for video editing logic, etc. This makes it easier to swap out implementations or upgrade components (like replacing our custom ffmpeg code with a higher-level library or a cloud video rendering service if we ever chose).
- **Logging for ML Components:** It's wise to log the inputs/outputs of AI components for debugging. E.g., log the final prompt given to GPT and the script it returned (sanitize sensitive info). Log which stock video IDs were picked for which line. This helps later if an output was weird – we can trace back why (maybe the prompt was missing context, etc.). These logs can be stored in our DB (maybe a JSON field in `ads` record) or just text logs.
- **CI/CD for Machine Learning:** If we start fine-tuning models or using large models, incorporate that into our pipeline carefully. For example, if we fine-tune a small transformer on brand data, store the model artifacts and version them. Possibly integrate with something like DVC or just a naming convention. While not needed day 1, planning for experiment tracking is good (so we know which model version created which output – especially if we do brand-specific fine-tunes).
- **Community and Support:** Given we might rely on open-source components (like Coqui, etc.), engage with those communities for help and keep an eye on updates. Also monitor OpenAI policy changes or pricing changes which could affect us.
- **Ethical Considerations:** Using AI avatars or voices (we're not doing deepfakes of real people in MVP, but if we ever allow custom voices or avatars, ensure we have consent and it's used ethically – e.g., if creating a "digital twin" of a person, that person should be the user or have given rights).
- **Prometheus-lite Monitoring:** If we want metrics (CPU, memory), a simple approach is using *Docker stats* or *cAdvisor* and perhaps pushing metrics to a free monitoring service. Since this is mentioned, a lightweight might be using **UptimeRobot** plus maybe a self-hosted Grafana if we can spare RAM. Alternatively, just ensure to measure the duration of jobs and number of jobs – if we find a backlog growing, that signals need to scale.
- **GitHub Actions & Testing:** Write a few unit tests for critical functions (like does `generate_script` return at least one line? Does `synthesize_voice` actually produce a file?). In CI, we can't run ffmpeg with heavy media (no GPU, etc.), but maybe we mark heavy tests as integration to run manually. We can mock external API calls in tests to not consume tokens.
- **User Feedback Mechanism:** Since we claim to learn from feedback, it's good to allow the user to rate or choose preferences too. Maybe in the UI, after seeing the video, they can press a thumbs-up or down. We can log that as an explicit feedback to supplement performance metrics. This helps especially if an ad hasn't gotten real-world data yet but the user knows something is off (they can signal and we adjust next output accordingly).
- **Continuous Deployment:** With a small team, automating deployment is key. We set up auto-deploy from GitHub Actions to our server, maybe triggered on push to main (after tests pass). This ensures we can ship fixes or improvements quickly (which is our advantage over bigger competitors).
- **Documentation:** Write internal docs (even a README) on how to run the system, how each component works, and troubleshooting tips. This is important if we bring on new team members or if one of us is unavailable, the other can manage. Also document how to add a new component (so in future if we integrate, say, a new model for concept generation, we follow a pattern).
- **Plan for Future Features:** Keep the code extensible. E.g., we know "avatar video" is a possible future addition – maybe structure the code so that the `retrieve_broll` function could under the hood either fetch a stock video or generate an avatar video depending on input parameters. By abstracting "get visual for scene" as a single function, we can later plug different sources. Same for voices (if we add more voices or use ElevenLabs for premium quality, we can branch on config).

- **Analytics for Business:** Implement a way to track usage of the system itself: how many ads generated, average time, which features used. This will help with everything from debugging to pitching investors (we can show engagement metrics). Even if it's just logging or sending events to a simple Google Analytics or a Mixpanel, it's useful.

Finally, keep everything **lean and focused** – as per user's direction, "no fluff, no bloat." This means if a feature is not immediately adding value to hitting MVP goals, push it to backlog. For instance, multi-language support is great (and our architecture could handle it by swapping TTS model and translating script with GPT), but if our initial clients are all English-speaking DTC brands, we don't need it in v1. Similarly, fine-tuning custom models per brand sounds cool but might be overkill until we see a need. We should concentrate on a few key differentiators: - High-quality video output (so invest effort in that polish like transitions, resolution, sync). - Feedback loop integration (even if minimal at first, just showing we capture data). - Modular design (so we can quickly adapt to any early user feedback, like "can it make it more X?", we can tweak that module).

By adhering to these practices, we aim to deliver a stable v1 that impresses and forms a solid foundation for future enhancements. We'll monitor, learn from actual usage, and iterate rapidly – that's our advantage as a small focused team.

---

1  Microkernel Architecture Pattern: Understanding Software Architecture Patterns[3] | by Algorithm Alchemist | Medium
https://nerdnodes2023.medium.com/microkernel-architecture-pattern-understanding-software-architecture-patterns-3-1a22f0640118

2  3  Demystifying Python Celery: key components and result storage
https://www.vintasoftware.com/blog/celery-overview-archtecture-and-how-it-works

4  7  8  47  CrossFade, Dissolve, and other Effects using FFmpeg's xfade Filter - OTTVerse
https://ottverse.com/crossfade-between-videos-ffmpeg-xfade-filter/

5  AI Voice Actor For Video Ads With 140+ Realistic TTS Voices - Creatify
https://creatify.ai/features/text-to-speech

6  12  15  18  21  22  23  29  34  35  36  40  41  42  reelforge-competition-deep-research.pdf
file://file-6MGG1v2c3TUfuTWemmHHgU

9  10  11  How AI Agents Are Redefining Enterprise AI Architectures
https://superbo.ai/beyond-llms-how-ai-agents-are-redefining-enterprise-ai-architectures/

13  Top 7 Free AI Agent Frameworks - Botpress
https://botpress.com/blog/ai-agent-frameworks

14  24  30  31  32  33  37  38  39  Should Performance Marketers Use Creatify AI or AdCreative AI? - Craft short video ads for your products from any URL
https://creatify.ai/blog/adcreative-ai-comparison

16  How it works - Learn how Pencil works in minutes
https://www.trypencil.com/how-it-works

17  19  Pencil Review: Fire Your Editor With This Facebook Ads AI Software
https://aazarshad.com/resources/pencil-ai-review/

[20] [25] [26] [27] [28] Why Adcreative.ai Outperforms Canva, Adobe Express, and Pencil in Ad Optimization

https://www.adcreative.ai/post/adcreative-ai-vs-canva-adobe-express-and-pencil-the-clear-choice-for-high-performing-ads

[43] Mastering Task Orchestration with Celery: Exploring Groups, Chains ...

https://mazaherian.medium.com/mastering-task-orchestration-with-celery-exploring-groups-chains-and-chords-991f9e407a4f

[44] Dagster & Celery: Mastering Data Orchestration

https://www.getorchestra.io/guides/dagster-celery-mastering-data-orchestration-dagster-celery

[45] How to combine videos using FFmpeg concat? - Gumlet

https://www.gumlet.com/learn/ffmpeg-concat/

[46] crossfade between 2 videos using ffmpeg - Super User

https://superuser.com/questions/778762/crossfade-between-2-videos-using-ffmpeg

[48] Orchestrating a Background Job Workflow in Celery for Python - Toptal

https://www.toptal.com/python/orchestrating-celery-python-background-jobs