



DOCTRINE_02

The Agentic NLE

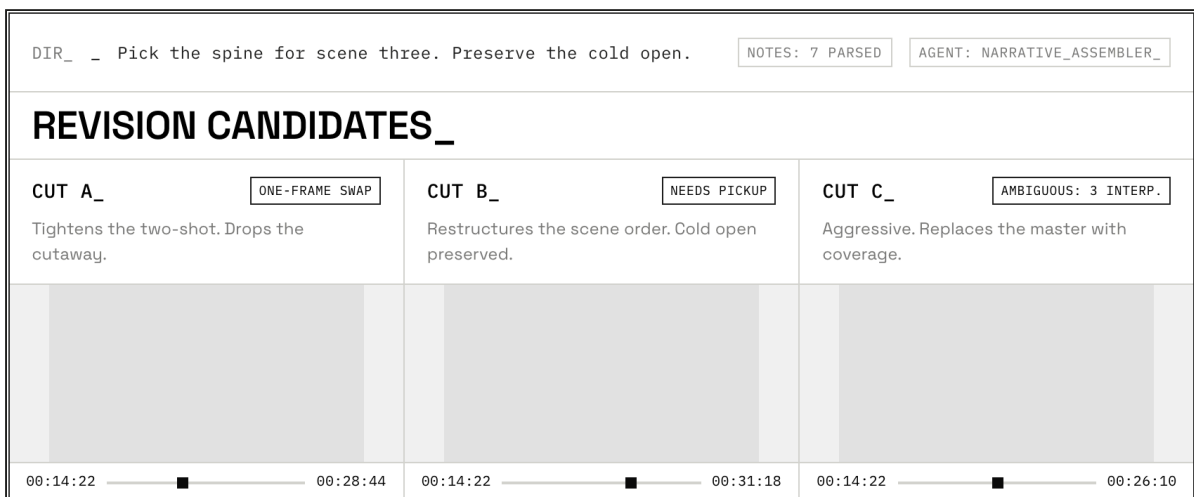
The future of editing isn't AI replacing editors. It's AI accelerating them. So why does the industry keep building AI tools for tape-era NLEs instead of NLEs built for AI?

An editor opens the workspace mid-project. Yesterday's first cut went out to the client at 6pm. The notes came back at 11pm.

The review-agent read them, scored each for feasibility, and sent three interpretations of the structural changes through the assembler overnight. By morning they sit side by side in an A/B/C panel. Reaction shots already slotted into the silences as J-cuts and L-cuts.

The timeline isn't a strip of horizontal tracks. It's a grid of scene blocks, because that's how this editor has been working for months, and the surface has shaped itself to match. Scopes sit tucked away. The performance grid is up.

The editor scrubs the second pass, swaps a take in scene three for the backup the assembler flagged, marks one shot to feel colder, drags one frame six frames earlier on the master. The system makes a note. It won't make that mistake again.



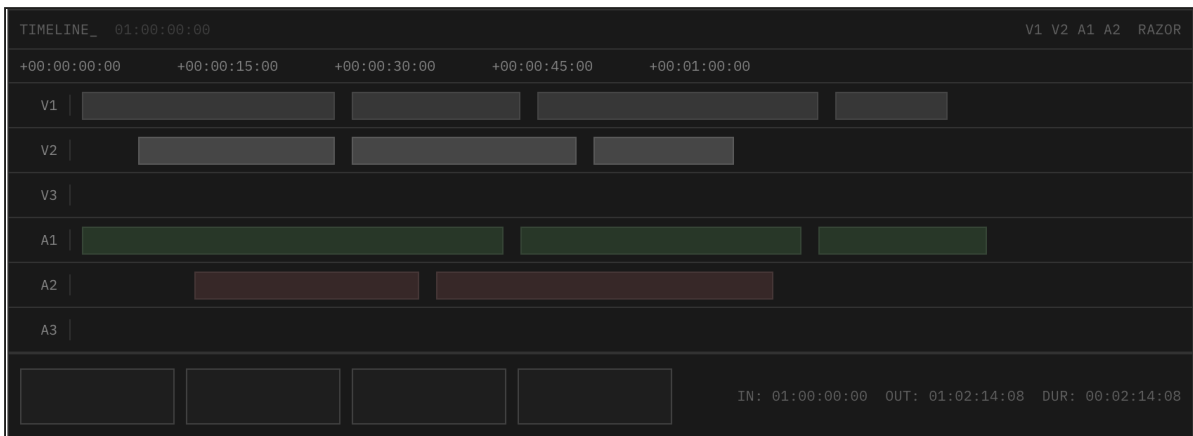
02 — A/B/C REVISION PANEL · POST-CLIENT FEEDBACK

The 1991 in the 2026 software

Look at any timeline today and you're looking at 1991.

The track grid is a fossil. V1, V2, A1, A2, video one, video two, audio one, audio two, were the names of physical tape decks stacked on top of each other in an edit suite. The labels survived into the software. So did the geometry. So did the address.

Hours, minutes, seconds, frames is the way you found your place on a reel that physically rotated past a playhead. Even the razor, the primary cutting tool, is a mechanical metaphor for a knife you used to slice celluloid. The source and record monitors are two boxes that represent two real machines, one that played the take, one that recorded the cut.



07 – FRAME-LEVEL ESCAPE HATCH · 1991 TRACK GRID

None of this is wrong. It was correct for the medium it was designed in. Tape was linear, physical, expensive to re-record, and an editor's job was to choose the right pieces and place them in time. The track grid mapped that work cleanly onto a screen.

Then the medium changed. Footage became file-based. Cuts became free. The job changed too. Editors moved from selection-and-placement into orchestration of dialogue, pacing, performance, colour, sound, story. The medium got bigger. The interface stayed the same.

Every AI feature shipped into editing software since then is a renovation of one room in this house. Auto silence-cut tools automate the razor: same mechanical move, faster. Magic-mask roto automates matte tracking: same keyer, but the keyer no longer needs you to draw the spline. Auto-transcribe-edit treats spoken text as a clip you can slice and rearrange, but the clip still lands on the same track at the same address. Generative b-roll fills the media bin with synthetic footage you didn't shoot, but the bin is the same bin and the clips drop onto the same V2.

Smart reframe rebuilds the crop tool. Scene detect rebuilds the catalogue. Auto-colour-match rebuilds the matching shot. Each one is a useful renovation. None of them touch the floor plan.

The ceiling shows up the moment you try to give the software a goal. Tighten this dialogue. Build to a fast crescendo. Make scene one feel nostalgic and warm and scene two feel sterile and cold. These aren't track-level instructions. They're directorial instructions. They reach across colour and pacing and audio and structure at once, which is what direction does.

To run them, the AI has to translate the instruction down into a hundred track-level operations: shorten this clip on V1, slot a reaction on V2, swing a curve on the colour node, lift the sub on the music stem. The instruction lands. The orchestration is gone. What an editor reads off the timeline afterwards is not the directorial choice. It's the renovation list.

This is the cap. Goal-oriented AI does not fit on a track grid because the track grid is not the language of cinema. It is the language of one craft inside cinema, plumbed through a substrate from a

previous era. As long as the substrate stays the same, the AI features layered on top of it can only get more precise about their renovations. None of them will move the house.

Storyline over Timeline

Start the house over. The primary unit isn't a clip on a track. It's a scene block.

A scene block is a region of the story, not a region of the timeline. It holds whatever clips, audio, masks, and decisions belong to that beat: interview setup, action sequence A, the reaction, the resolution. The block knows what it is. The editor names it and the system tags it. The micro-cuts inside a block are handled by an agent that understands pacing and intent. The editor stays at the block layer for most of the work and drops below it only when something specific needs hand-work.

Above the blocks sits the workspace itself, and the workspace is the interface. Not a fixed surface with panels in fixed places. A mutable surface that reshapes around what the editor is doing. A/B panels when two cuts need comparing. A storyline view when story is the unit. A performance grid when one actor's beats are the unit. A colour bay when the grade is up. The surface knows what the editor is working on and rearranges to put the right tools in front.

And the workspace evolves with the editor. The frames that get hand-corrected, the takes that get swapped, the masks that get redrawn over and over become signal. The system learns this editor likes reactions held a beat longer. This editor prefers a slightly cooler key on dialogue scenes. This editor wants the assembler to cut to the wide before the reaction, not after.

The next session opens sharper than the last. The longer the editor uses it, the more the tool becomes theirs. Same engine, completely different feel per user. The way the best agentic tools already work in other crafts.

The track grid does not vanish. It survives as an escape hatch. When something needs to be done frame-accurate, by hand, with no abstraction in the way, the editor breaks through. The virtual razor is still there. So is the timecode. So are the tracks. The difference is they're no longer the surface you start at. They're the surface you fall through to when the abstraction can't reach what you want to change.

This is the swap. Operator-level controls hide under director-level interaction. Cinema's working vocabulary moves above the tape-era vocabulary. The tape-era vocabulary stays available, but it stops being the front door.

DIR_ - Tighten the interview. Hold on reaction two beats longer.

LAYOUT: ARI_DEFAULT_V14

LEARNED: REACTIONS_HELD_+8F

STORYLINE_

SCENES: 12

RUNTIME: 04:32:18

AGENTS: 2 ACTIVE

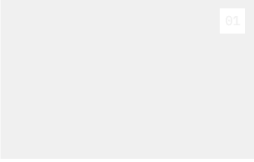
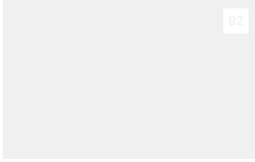
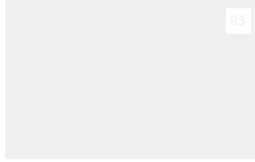
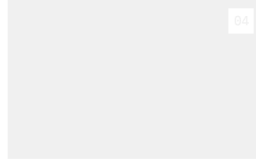
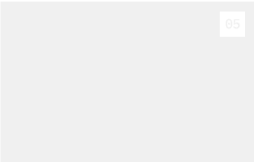
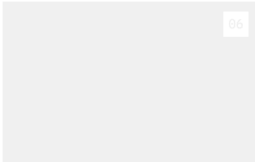
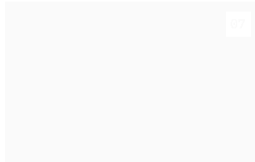
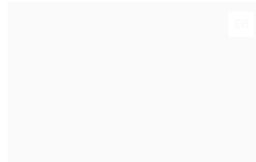
01	INTERVIEW_SETUP_	00:42:15	CAM_A
02	ACTION_A_	01:18:04	CAM_B + CAM_C
03	REACTION_	00:14:22	CAM_A CLOSE
04	RESOLUTION_	00:56:08	CAM_A + CAM_B
05	CREDITS_	00:08:00	ROLL

01 - STORYLINE VIEW · DEFAULT STATE

DIR_ - Find a take where the subject lands the line dry. ACTOR: M. CARTER TAKES: 14

PERFORMANCE GRID_

SCENE: 02 ACTION_A CIRCLE: 3 TAKES

 <p>TK_01 94</p> <p>wry → grounded 50MM / TUNGSTEN</p>	 <p>TK_02 87</p> <p>guarded → open 50MM / DAYLIGHT</p>	 <p>TK_03 42</p> <p>overcooked → false 85MM / TUNGSTEN</p>	 <p>TK_04 91</p> <p>dry → precise 50MM / TUNGSTEN</p>
 <p>TK_05 38</p> <p>rushed → uneven 35MM / DAYLIGHT</p>	 <p>TK_06 89</p> <p>subtle → true 50MM / TUNGSTEN</p>	 <p>TK_07 55</p> <p>hesitant → soft 85MM / MIXED</p>	 <p>TK_08 61</p> <p>earnest → flat 35MM / DAYLIGHT</p>


03 - PERFORMANCE GRID · TAKE SELECTION

DIR_ - Make scene one warmer. Protect the highlights in the window. SCOPE: VECTORSCOPE + WAVEFORM BLOCK: 01 INTERVIEW_SETUP


COLOUR BAY_

NODE: 03 / 12 GAMUT: REC.709


VECTORSCOPE




WAVEFORM




RGB PARADE



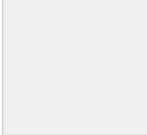
PRIMARY WHEELS_

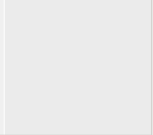


CURVES_



PREVIEW_





SOURCE GRADED

04 - COLOUR BAY · GRADING SUITE

The crew

Inside this workspace, the work gets done by a crew. The first member is the DIT-assistant, the agent that meets the footage at ingest. Point it at a folder of raw takes and it builds proxies, syncs multicam, aligns sound, transcribes the lot.

It also tags. Not just speech to text but camera move, lighting condition, lens, depth, an emotional read on the performance. It scores usability and bubbles the circle takes to the top, marks the soft-focus ones for review, files the obvious throwaways out of sight. The editor never opens a clip bin to hunt for a take. The crew has already laid them out in the order the script implies.

The Narrative Assembler builds the assembly. Feed it a script, a storyboard, or a beat sheet, and it walks the tagged footage block by block, choosing takes that match intent and sequencing them at a pace the story asks for.

It understands cinematic rhythm: when to hold a wide, when to cut to the reaction, when to use a J-cut to let dialogue overlap the next beat, when to let an L-cut carry the previous shot's audio into the new image. Ask for three versions of a scene at different paces and it returns three, each editable in the same surface, not a separate timeline. Compare, pick the spine, work from there.

The VFX and gap-filling agent does the work that used to require a node tree and a tracking pass. The editor highlights a distracting logo on a shirt and asks for it gone. The agent tracks, masks, paints, holds the integrity of the shot across the cut.

When the script calls for a wide of a neon-lit street and the rushes don't include one, the agent runs a generative video model in the background and proposes options the editor can audition and place. The b-roll is style-matched to the surrounding footage on grain, colour, and lens behaviour before it lands in the timeline. The intervention is reversible at any point.

The Audio and Foley agent does what a sound designer would, in the rhythm the picture asks for. It separates dialogue, music, and effects into clean stems on its own. It generates the room tone that should be there but wasn't captured. It builds footsteps that match the surface someone is walking on, fabric movement that matches what they're wearing, wind that matches the wide shot of the cliff.

It scores under the cut, swelling and falling on the beats the editor has chosen, not on the bar lines of a library track. The mix is rough by default but the structure is there from the assembly onwards.

The Finishing and Color agent grades for story, not for match. A directorial note becomes a per-block grade: make scene one feel nostalgic and warm, scene two sterile and cold, and the system propagates the choice across every shot in those blocks.

Because the system holds depth information from the rushes, it can relight as well as grade: shift a key light in post, recover a face that was underexposed, hold continuity across takes shot hours apart in different daylight. The grade is reversible to source. Nothing is baked.

DIR_ _ Render proxies for the offline cut. Tag audio stems.					CREW: 5 AGENTS	QUEUE: 3 JOBS
CREW DISPATCH_					EST. COMPLETION: 00:14:00	
■ DIT-ASSISTANT_ running Ingesting rushes from CARD_B. Generating 1:4 proxies.	NARRATIVE_ASSEMBLER_ ■ VFX_ idle Awaiting editor selection on A/B/C revision panel.	■ AUDIO_FOLEY_ running Generating ambience beds. Spotting SFX for impact cuts.	■ FINISHING_ idle Colour and finishing hold until picture lock signal.			
00:04:12 PROXY_04_DONE 00:03:58 PROXY_03_DONE 00:03:41 PROXY_02_DONE	00:00:00 IDLE 23:42:18 CUT_C_RENDERED 23:38:05 CUT_B_RENDERED	00:01:33 REVIEW_READY 23:55:12 TRACK_SOLVED 23:51:08 PLATE_EXTRACTED	00:02:18 BED_02_RENDERED 00:01:44 IMPACT_07_SPOTTED 00:00:12 DIALOGUE_CLEANED	00:00:00 IDLE 00:00:00 IDLE 22:14:33 LUT_APPLIED		

05 – CREW DISPATCH PANEL · AGENT STATES

Five specialists. None of them the editor. All of them reporting to the editor. The work that used to take five different application suites and five different file handoffs runs as a single crew on a single surface, and the editor stays in the chair the whole time.

What the editor does with the time the crew gives back is the point. The numbers tell the story before the argument does. Roughly two-thirds of an editor’s working day today goes to sync, ingest, conform, the find-the-take grind, the masking that has to be done by hand because no software understands what the shot is about. Take that work away and the day rebalances. The recovered hours land on choice. Which version, which performance, which rhythm, which tone. The shift is from operator to director.





The iteration loop is where this lands hardest. Review a cut, name what's wrong: the pacing in the middle drags, the colour is too saturated, the reaction at minute three needs to land harder. The system recalculates. A new pass plays back inside the same surface. The cycle is seconds, not days, and the editor is making directorial choices in every cycle, not technical ones.

Client feedback is where this lands second-hardest. The notes come in. The review-agent reads them alongside the cut and scores each one for feasibility. This one is a one-frame swap. This one needs a take that wasn't shot. This one is ambiguous and could go three ways.

By the time the editor opens the workspace, the obvious changes are already proposed. The ambiguous ones come back with three candidate interpretations sitting in an A/B/C panel, each one a working pass the editor can scrub and choose between.

The editor reads the notes once, picks the spine, refines what needs refining, and the revision goes back the same morning. The revision cycle that used to take days runs in hours.

DIR_ _ Approve all one-frame swaps. Flag ambiguous notes for review. PARSED: 7 / 7 CONFIDENCE: 94%

REVIEW INBOX_

CLIENT: FRAME_DRIFT TURNAROUND: OVERNIGHT

ORIGINAL NOTE	FEASIBILITY	PROPOSED MOVE	ACTION
The interview feels too long. Can we lose the second question?	ONE-FRAME SWAP	Trim O1_INTERVIEW_SETUP by 00:08:14. Remove Q2 beat.	APPROVE REFINE
Action scene needs more urgency. It drags.	NEEDS PICKUP	Replace TK_03 with TK_07 in O2_ACTION_A. Tighten by 6F.	APPROVE REFINE
Not sure about the reaction shot. Feels staged.	AMBIGUOUS: 3 INTERP.	Alt A: Replace with TK_04. Alt B: Extend by 2F. Alt C: Remove entirely.	APPROVE REFINE
Make the ending hit harder. More final.	NEEDS PICKUP	Add 2F hold on final frame. Extend music tail 00:04:00.	APPROVE REFINE
The colour on scene one feels cold.	ONE-FRAME SWAP	Lift warmth +12 units in O1_INTERVIEW_SETUP grade node.	APPROVE REFINE
Missing a close-up of the hands during the exchange.	NEEDS PICKUP	Insert CU_TK_02 at 02:14:08. Requires plate extension.	APPROVE REFINE
Overall length is fine but pacing in the middle lags.	AMBIGUOUS: 3 INTERP.	Alt A: Compress midsection 15%. Alt B: Remove B-roll bridge. Alt C: Add music sting.	APPROVE REFINE

06 – REVIEW-AGENT INBOX · PARSED CLIENT NOTES

The track grid does not disappear and the editor does not become a manager. The escape hatch matters. When a specific frame needs to be six frames earlier, when a specific mask needs to be hand-drawn, when the system's choice is wrong, the editor breaks through. The work happens at the level the work needs to happen at. The difference is that the editor chooses that level. The substrate doesn't impose it from below.

When the editor does break through, the system watches. The frames that get manually adjusted, the masks that get manually redrawn, the takes that get manually swapped become training data for the next pass. Not training data in a model-fine-tuning sense. Training data in a stylistic sense. This editor likes to hold reactions a beat longer than the assembler proposes. This editor

prefers a slightly cooler colour pass than the assembler suggested. The system learns the editor, not the medium. Over time, the assembler stops needing to be corrected on the same thing.

The first time an editor works inside a system like this, the most surprising part is what they stop doing. The hour spent finding the right take of the right line, gone. The hour spent matting out the wrong logo, gone. The two hours spent syncing audio to picture across four cameras, gone. What replaces them is not less work. It's more of the work the editor wanted to be doing the whole time.

Read this way, what looks like automation is the opposite. It's the recovery of judgement time. The editor is doing more of the work that is editing: taste, rhythm, intention, story. And less of the work that was a side effect of how the substrate was built. The role doesn't shrink. It concentrates.

Call it what it is. The editor becomes an edit orchestrator. Crew below, surface around, judgment in the chair.

So who's building it?

The components exist. Goal-oriented language models exist. Specialist agents exist. Generative video, generative audio, depth-aware grading, scene-level tagging. All of it ships in something today, somewhere in the stack. What doesn't exist yet is the integrator that thinks like an editor and gets out of the editor's way.

Someone is going to ship it. The question is who, and the question is when, and the question is whether editors get to shape what gets built or whether it gets built around them.

A good NLE is the most patient instrument an editor ever holds. It waits. The next good NLE will do more than wait. It will collaborate, on terms an editor can recognise as theirs. That tool isn't here yet.

So. Who's building it?