

AI STRATEGY GUIDE

# The Multimodal AI Marketing Engine: Text, Image, and Voice Working Together

How to connect text, image, and voice AI into a single production system that multiplies output without multiplying headcount

**Carlos Rivera**  
Founder, NetWebMedia

14 pages  
[netwebmedia.com](https://netwebmedia.com)

# Contents

---

## **The Multimodal Opportunity: Why Single-Modal AI Teams Are Leaving ROI on the Table**

The compounding output gap between teams using isolated AI tools versus teams running integrated multimodal systems.

---

## **The Modality Map: What Each AI Type Does Best in Marketing**

A decision framework for matching AI modality to marketing task — with the specific use cases where each type delivers its highest ROI.

---

## **Orchestration Architecture: Connecting Text, Image, and Voice AI into One System**

The technical and operational architecture for connecting three AI modalities into a single production pipeline with defined handoffs and quality gates.

---

## **The Source-to-Distribution Pipeline: One Recording, 20 Outputs**

A detailed walkthrough of the highest-ROI multimodal pipeline — converting a single executive recording into a full content distribution package.

---

## **Tool Selection by Modality: The 2026 Stack Recommendations**

Current best-of-class tool recommendations for each modality, with selection criteria and the specific use cases where each tool leads.

---

## **Quality Standards Across Modalities: Maintaining Brand Consistency**

How to set and enforce consistent quality and brand standards when output spans three different AI modalities with different failure modes.

---

## **Building with Make/n8n vs. Claude Agent SDK: When Each Approach Fits**

A practical comparison of no-code workflow automation versus agent SDK for building the routing layer of a multimodal marketing system.

---

## **Measuring Multimodal ROI: Cross-Channel Attribution for AI-Produced Content**

How to measure the business impact of multimodal AI production across channels where AI assets commingle with human-produced content.

---

EXECUTIVE SUMMARY

# The Multimodal AI Marketing Engine: Text, Image, and Voice Working Together

Most marketing teams have adopted AI tools, but they've adopted them in silos. The content team uses ChatGPT for copy. The design team uses Midjourney for visuals. Someone on the demand team experiments with a voice clone for a podcast. Each tool delivers some individual productivity gain, but the gains don't compound — because the tools aren't connected. The multimodal opportunity is not about having more AI tools. It's about connecting them into a production system where the output of one modality becomes the input for another, and where a single piece of source content fans out into an entire channel's worth of distribution assets. This guide is a practical architecture document for building that system — not a survey of what's theoretically possible, but a blueprint for what B2B marketing teams are doing operationally in 2026.

## IN THIS GUIDE

- ✓ Why teams using single-modal AI — text only, image only, voice only — are capturing less than 30% of the available efficiency and output gains
- ✓ The modality map: a decision framework for which AI type to deploy for which marketing task
- ✓ The source-to-distribution pipeline architecture that turns one executive recording into 20 distinct content outputs
- ✓ A curated 2026 tool stack with selection criteria by modality and use case
- ✓ A 90-day plan for scaling from a multimodal pilot to a full production operation

**Who this is for:** B2B marketing directors and CMOs who have adopted individual AI tools but have not yet integrated them into a coordinated production system.

## SECTION 1

# The Multimodal Opportunity: Why Single-Modal AI Teams Are Leaving ROI on the Table

A team that uses only text AI produces more copy faster — a real but bounded gain. A team that uses text, image, and voice AI in a coordinated system produces an entirely different category of output. The difference is not additive; it's multiplicative. When a transcript from a voice recording becomes the source document for an AI writing tool, which produces a long-form article that is then broken into ten social variants, each paired with AI-generated visuals and reformatted for platform-specific dimensions — you've converted a single 30-minute executive conversation into a month's worth of channel content. That pipeline doesn't exist when the tools are isolated. Each modality's output sits in a different folder, produced by a different team member, with no systematic handoff to the next stage. Single-modal AI adoption is better than nothing, but it still requires the human connective tissue that most content teams don't have capacity for at scale.

The ROI gap between single-modal and multimodal operation compounds over time. Single-modal teams see efficiency gains plateau after 60-90 days — they've optimized individual tasks but haven't restructured the production system. Multimodal teams see gains accelerate because every new source asset flows through the same pipeline and produces the same range of derivatives. A team of four running a mature multimodal system typically outproduces a team of twelve running disconnected tools on a content-volume basis. The productivity leverage doesn't come from using more AI — it comes from eliminating the manual handoffs between modalities that consume most of the production time.

- Audit your current AI tool stack: which tools are used by which team members, and which tools' outputs feed other tools
- Identify the three highest-volume manual handoffs in your current content production process
- Calculate how many distinct distribution assets you produce from a single source piece today
- Benchmark your current output rate: how many finished, publish-ready pieces per week per content producer
- Set a multimodal output target: what would 3x that rate look like in terms of pipeline contribution

*Teams with mature multimodal production systems report producing 8-12 distribution-ready assets per source piece — vs. 1-2 for teams using single-modal or disconnected tools.*

**71%**

of B2B marketing teams using AI tools are single-modal or unconnected, according to 2025-2026 practitioner surveys

## SECTION 2

# The Modality Map: What Each AI Type Does Best in Marketing

Not every marketing task benefits equally from every AI modality. Deploying the wrong modality for a task produces mediocre output and wastes configuration effort. The modality map provides a decision framework for matching AI type to marketing application. Text AI (large language models) handles everything that requires reasoning, synthesis, and structured communication: long-form content, email copy, ad variants, landing page copy, SEO content, research synthesis, brief writing, and conversation (chatbots, sales enablement). Text AI is the connective tissue of a multimodal system — it produces the structured language that other modalities can reference or that downstream channels can consume directly. Image AI (diffusion models and image generation) handles visual asset production, background generation, product visualization, brand imagery, ad creative, and social graphics. It does not replace photography for product or people shots, but it eliminates the creative asset bottleneck for abstract, conceptual, or illustrative imagery at high volume.

Voice AI covers three distinct applications: transcription (converting spoken content to text), synthesis (converting text to audio), and conversation (real-time interactive voice agents). Each has different quality thresholds and appropriate use cases. Transcription is production-ready across all major tools and should be part of every content team's workflow for repurposing executive and subject matter expert conversations. Synthesis is production-ready for internal content, podcasts, and specific external applications where the use case doesn't require the authenticity of a human voice. Conversational voice AI for customer-facing applications requires more rigorous design and testing — covered in a separate guide. The modality map decision rule: start with the source asset type (spoken, written, visual), identify the target distribution format, and select the AI modality chain that spans the gap with the fewest manual steps.

- Text AI: Long-form articles, email sequences, ad copy, landing pages, research synthesis, chatbot scripts
- Image AI: Social graphics, ad creative, concept illustrations, background imagery, visual variants for testing
- Voice transcription AI: Converting executive conversations, interviews, and recordings to structured text source material
- Voice synthesis AI: Podcast production, internal training narration, accessible content formats
- Multimodal chain: Source audio → transcription → structured text → long-form content → social variants → paired visuals

*Voice transcription is the most underused modality in B2B marketing — it converts every executive conversation, sales call, and webinar into structured source material for text AI at near-zero marginal cost.*

**3** primary AI modalities (text, image, voice) that cover 95% of B2B marketing content production use cases

### SECTION 3

## Orchestration Architecture: Connecting Text, Image, and Voice AI into One System

Orchestration is the difference between a tool stack and a production system. Orchestration means the output of one tool automatically becomes the input for the next, with defined quality gates, routing logic, and exception handling at each handoff. The orchestration architecture for a multimodal marketing system has four layers. Layer one is the source layer: the inputs that seed the pipeline — executive recordings, research documents, campaign briefs, product information. This layer determines pipeline quality more than any other; garbage in, garbage out applies at multimodal scale. Layer two is the transformation layer: the AI tools that convert source material into production-ready content. Voice transcription converts audio to text. Text AI transforms raw transcripts into structured documents. Image AI generates visual assets. Each transformation step has a quality gate — a check that output meets the input requirements for the next stage. Layer three is the distribution layer: the channel-specific formatting, sizing, and packaging of production assets for their destination platforms. Layer four is the routing layer: the workflow automation

(Make, n8n, Zapier) that moves assets between layers and notifies humans when exceptions require intervention.

The most common architecture failure is building the transformation layer without the routing layer. Teams configure great AI tools but still manually move outputs between them. The routing layer can be as simple as a structured folder system with naming conventions that trigger the next step, or as sophisticated as a full workflow automation that handles branching logic based on content type. Start simple: even a documented manual routing protocol with clear handoff checklists is better than ad-hoc movement between tools. The goal of the routing layer is to eliminate the question 'what do I do with this now?' from every production step.

- Document your current source-to-distribution production flow before designing the orchestrated version
- Identify every manual handoff step in the current flow — each is an automation opportunity
- Define the input and output specification for each AI transformation step: what format goes in, what format comes out
- Build quality gates: specific pass/fail criteria that must be met before output moves to the next stage
- Start the routing layer as a documented manual protocol before automating — you'll identify edge cases
- Map exception routing: which exceptions require human review vs. which can be auto-retried
- Version-control your orchestration documentation — it will change as tools evolve

*The routing layer — the workflow automation connecting tools — is typically 10% of the build effort and 60% of the production efficiency gain. Most teams skip it and wonder why their multimodal stack doesn't deliver.*

**4** architecture layers in a functional multimodal marketing system: source, transformation, distribution, routing

#### SECTION 4

## The Source-to-Distribution Pipeline: One Recording, 20 Outputs

The flagship multimodal pipeline in B2B marketing starts with a 20-40 minute recording of an executive, subject matter expert, or customer interview — a resource most organizations already have but systematically underuse. Here is how a single recording produces 20 distribution-ready outputs. Step one: transcription. The recording goes through a high-accuracy transcription tool (Whisper API, Otter.ai, or Deepgram) producing a structured transcript with speaker labels and timestamps. Step two: transcript structuring. A text AI prompt converts the raw transcript into a structured source document: key claims, supporting evidence, narrative arc, notable quotes, and data points. This structured document is the source for all downstream content. Step three: long-form content. The structured document feeds three long-form outputs: a long-form article (800-1200 words), a LinkedIn article variant, and an email newsletter feature. Step four: short-form breakdown. The long-form article is broken into 8-10 social posts (platform-specific length and format), 3-5 pull quotes formatted for social graphic templates, and a 60-90 second audio clip script. That's 14-18 text outputs.

Step five adds the image layer. The 3-5 pull quotes feed an image AI prompt that generates matched visual assets — abstract imagery, typography-forward graphics, or concept illustrations for each quote. A single prompt template, adapted for each quote, produces 3-5 image assets. Step six adds the audio layer. If voice synthesis is in the stack, the article or newsletter section can be converted to an audio version for podcast distribution or an audio clip for social. That adds 1-2 audio assets. Total: 18-22 distribution-ready outputs from a single 30-minute conversation. The key operational requirement: the pipeline needs to be documented and templated so any team member can run it without inventing steps. The templating is what makes the pipeline scalable — it's not a one-off production process, it's a repeatable machine.

- Step 1: Transcribe the recording with speaker labels using Whisper API, Otter.ai, or Deepgram
- Step 2: Run the structured source document prompt to extract key claims, quotes, and narrative arc
- Step 3: Generate the long-form article, LinkedIn variant, and email newsletter feature from the source doc
- Step 4: Break the article into 8-10 social posts and 3-5 pull quotes with platform-specific formatting
- Step 5: Feed pull quotes to image AI to generate matched visual assets for each
- Step 6: Convert article or newsletter section to audio if voice synthesis is in the stack
- Templating the entire pipeline so any team member can run it without inventing steps

***The structured source document in Step 2 is the leverage point of the entire pipeline. Investing 20 extra minutes in prompt quality at this step pays dividends across all 18-22 downstream outputs.***

**20+**

distribution-ready outputs producible from a single 30-minute executive recording in a mature multimodal pipeline

**SECTION 5**

## **Tool Selection by Modality: The 2026 Stack Recommendations**

Tool selection in multimodal AI is a moving target — capabilities shift quarterly. The 2026 recommendations below reflect capability state as of early 2026 and should be reassessed every two quarters. Text AI tier one: Claude 3.5/3.7 Sonnet for long-form content, complex research synthesis, and any task requiring strong reasoning. GPT-4o for content types requiring broader creative range or when OpenAI's ecosystem integrations are valuable. Gemini 1.5 Pro for tasks requiring very long context windows — processing large documents or extended research corpora. Text AI tier two for specialized applications: Writer for brand-governed content production at team scale, Jasper for template-based high-volume content with workflow management, Perplexity for research synthesis with live source verification. Image AI tier one: Midjourney v7 for brand and conceptual imagery requiring high aesthetic quality. DALL-E 3 (via API) for programmatic image generation integrated into automation pipelines. Stable Diffusion (via RunDiffusion or Replicate) for teams requiring fine-tuned models or custom style training.

Voice transcription: Whisper API (OpenAI) for highest accuracy on unstructured recordings at low cost. Deepgram for real-time transcription applications and speaker-labeled outputs. Otter.ai for teams wanting a managed transcription service with collaboration features. Voice synthesis: ElevenLabs for highest quality voice cloning and text-to-speech production. PlayHT for production-volume synthesis with API integration. Suno for AI-generated music and audio branding assets. Orchestration layer: Make (formerly Integromat) for teams wanting no-code visual workflow automation with strong AI tool integrations. n8n for teams wanting self-hosted open-source automation with maximum flexibility. Claude Agent SDK for teams building custom agents requiring complex reasoning within the automation. Stack selection principle: choose tools that have documented API access and webhook support — these are the tools that integrate cleanly into the routing layer. Tools without programmatic access create manual handoff dependencies.

- Text: Claude Sonnet for quality, GPT-4o for ecosystem, Gemini for long context — pick primary based on most common use case

- Image: Midjourney for quality brand work, DALL-E 3 API for automated pipeline integration
- Transcription: Whisper API for cost-sensitive volume, Deepgram for real-time and speaker labels, Otter.ai for managed service
- Synthesis: ElevenLabs for production quality, PlayHT for API integration at scale
- Orchestration: Make for no-code teams, n8n for self-hosted flexibility, Claude Agent SDK for complex reasoning flows
- Selection rule: API access and webhook support are non-negotiable for any tool in the automated pipeline

*The orchestration layer tool (Make, n8n, or SDK) is more important to production efficiency than any individual AI tool choice — connectivity beats capability when you're building a system.*

**Q2**

recommended reassessment cadence for multimodal tool stack — capabilities shift enough every 6 months to warrant a structured review

## SECTION 6

# Quality Standards Across Modalities: Maintaining Brand Consistency

Multimodal production introduces a quality challenge that single-modality operations don't face: each modality has different failure modes, different quality signals, and different brand risk categories. A text quality review framework doesn't catch image brand violations, and neither catches voice synthesis issues. The quality standard system needs to be modality-specific while enforcing a consistent brand identity across all three outputs. Text quality standards: adherence to brand voice (four-layer framework covered in Guide 13), factual accuracy for any claims referencing data or research, SEO optimization for content intended for search, and channel-specific format compliance. Image quality standards: brand palette adherence (specific color values, not approximations), visual style consistency across a campaign, absence of distorted or uncanny elements that damage credibility, and legal clearance for any image elements that might raise copyright questions. Voice quality standards: pronunciation accuracy for brand and product names, natural prosody that doesn't undermine credibility, and disclosure compliance for AI-generated audio used in customer-facing contexts.

The cross-modality quality requirement is visual and narrative coherence: when text and images ship together (social posts, ads, emails), they need to tell the same story at the same register. This is the most commonly broken quality standard in multimodal production because text and image are typically reviewed in separate workflows by separate reviewers. Build a joint review step for paired assets — reviewing text and image together, not in sequence. The tool for this review is simple: a review deck that presents each paired asset side by side with the brief, making coherence gaps immediately visible. Multimodal quality metrics to track: paired-asset coherence score (subjective reviewer rating), image brand compliance rate, voice accuracy rate, and the cross-modality revision rate.

- Build modality-specific quality rubrics: text (voice, accuracy, format), image (brand palette, style, distortion), voice (pronunciation, prosody, disclosure)
- Add a joint review step for paired assets — never review text and its paired image in separate workflows
- Create a review deck format that presents brief, text, and image together for coherence review
- Set specific color values for image AI prompts — not 'navy blue' but '#010F3B navy'
- Build a pronunciation guide for AI voice synthesis: product names, brand names, industry terms
- Track cross-modality revision rates by asset type to identify the highest-drift production steps

*Paired asset coherence — whether the image and text tell the same story — is the most commonly broken quality standard in multimodal production and requires an explicit joint review step to catch.*

**38%**

of multimodal marketing teams report frequent coherence mismatches between AI-generated text and image pairs in their first 90 days of production

## SECTION 7

# Building with Make/n8n vs. Claude Agent SDK: When Each Approach Fits

The routing layer — the automation that moves assets through the pipeline — can be built with visual no-code tools (Make, n8n) or with code-based agent frameworks (Claude Agent SDK,

LangGraph). The choice matters because it determines who can maintain the system, how much it costs to extend, and what level of decision-making complexity the routing logic can handle. Make and n8n are the right choice for the majority of multimodal marketing pipelines. They cover 80-90% of routing use cases through visual workflow builders, have robust AI tool integrations, handle branching logic, error routing, and notification triggers without code, and can be maintained by non-engineering team members with a day of training. Make has superior out-of-the-box integration coverage (5000+ apps) and a polished interface. n8n has the advantage of self-hosting for teams with data residency requirements and a more flexible data transformation layer. The practical ceiling of no-code tools: they struggle with complex conditional logic, multi-step reasoning about content quality, and tasks that require understanding context across the pipeline rather than just routing data between defined fields.

The Claude Agent SDK (and similar frameworks like LangGraph) is the right choice when the routing decisions themselves require intelligence. If the system needs to assess whether a transcript contains high-value insights before triggering the full pipeline, or route content to different output formats based on the topic classification of the source material, or make multi-step quality judgments that can't be reduced to field-matching rules — that's where agent architecture adds value that no-code tools can't replicate. The practical decision rule: start with Make or n8n. Document the cases where the no-code routing breaks down or requires human workarounds. Those specific cases are the candidates for agent logic. Build the agent layer to handle exceptions and complex decisions; keep the no-code layer handling the predictable volume. Hybrid architectures — where a Make workflow calls a Claude agent for specific decision points — are often the most practical production pattern.

- Use Make for: high-volume predictable routing, integrations with established SaaS tools, teams without engineering resources
- Use n8n for: self-hosted requirements, complex data transformation, teams comfortable with JSON and light technical work
- Use Claude Agent SDK for: routing decisions requiring content intelligence, quality assessment logic, complex conditional branching
- Hybrid pattern: Make/n8n for volume routing, Claude agent called at specific decision points requiring reasoning
- Start with no-code, document where it breaks down, then build agent logic only for those specific failure points
- Factor maintenance into the decision: agent-based routing requires an engineer; no-code can be maintained by a capable marketing ops person

***The hybrid architecture — no-code tools for volume routing with agent calls at specific intelligent decision points — is the most production-practical pattern for most B2B marketing teams.***

# 85%

of multimodal marketing routing use cases can be handled with no-code tools like Make or n8n without requiring custom agent development

## SECTION 8

# Measuring Multimodal ROI: Cross-Channel Attribution for AI-Produced Content

Multimodal ROI measurement is complicated by two factors: AI-produced and human-produced content are typically distributed through the same channels, making attribution messy, and the primary value of multimodal production is output volume and velocity — which is easy to measure — while the downstream business impact requires attribution infrastructure that most teams don't have at full fidelity. Start with the metrics you can measure cleanly. Production metrics: assets produced per week per content producer, time from source to published, ratio of AI-to-human revision time per asset. These establish the efficiency baseline. Engagement metrics by modality: track engagement rates separately for AI-generated visuals vs. designed assets, AI-written copy vs. human-written copy, and AI-produced audio vs. recorded audio. The engagement comparison doesn't require perfect attribution — it requires consistent tagging of asset origin in your CMS and social scheduling tools. Pipeline metrics: for B2B organizations with content-influenced pipeline attribution, the highest-value measurement is whether AI-produced content at scale is producing more pipeline influence per content producer than before the multimodal system was in operation.

The attribution framework for cross-channel AI content: tag every AI-produced asset at creation with modality type (text, image, voice) and generation tool. Carry that tag through to the publishing platform. Report on engagement, click, and conversion metrics segmented by AI-produced vs. human-produced, and by modality combination (text only, image only, text+image). Over time, the modality combination data tells you which AI asset types are performing and which need quality investment. The ROI summary metric: compare the cost per pipeline-influenced touch before and after multimodal system implementation. Account for the tooling cost (typically \$500-1500/month for a full stack) and the internal configuration and training time. Most teams with a mature multimodal system see a 40-60% reduction in cost per pipeline touch within 6 months — primarily from volume expansion with flat headcount.

- Track production efficiency: assets per producer per week, time from source to published, revision ratio

- Tag all AI-produced assets at creation with modality type and generation tool — carry tags to publishing platform
- Report engagement metrics segmented by AI vs. human-produced and by modality combination
- For B2B: measure pipeline-influenced touches per content producer before and after multimodal implementation
- Calculate true tool stack cost: licensing + configuration time + training time, amortized over 12 months
- Set 6-month ROI target at go-live: most teams hit 40-60% reduction in cost per pipeline touch

*Volume is a vanity metric without performance tracking. Tag every AI-produced asset at origin and track engagement by modality combination — this data tells you where to invest your next quality improvement cycle.*

**40–60%**

typical reduction in cost per pipeline-influenced touch for B2B teams with a mature multimodal production system vs. pre-AI baseline

## SECTION 9

# Scaling from Pilot to Full Operation: The 90-Day Plan

Most multimodal AI initiatives stall at the pilot stage — not because the technology doesn't work, but because the pilot succeeds with one source type and one output format and then nobody builds the broader system. The 90-day plan is designed to prevent that stall by sequencing the build in a way that produces early wins, validates the architecture, and builds team competency before expanding scope. Days 1-30: the pilot. Select one source type (executive recording is recommended), one transformation chain (transcription → structured document → article → 5 social variants), and one team member to operate it. Do not automate at this stage. The goal is to validate the source-to-distribution flow, identify where the AI tool outputs need refinement, and measure the time savings against your baseline. Document every manual step. Days 31-60: expand and automate. Add the image modality — pair AI-generated visuals with the social variant outputs. Build the routing layer for the validated pipeline: automate the handoffs you've been doing

manually, starting with the highest-volume step. Add a second content producer to the system. Train on the quality standards established in the pilot.

Days 61-90: expand to full operation. Add your second source type — if the pilot used executive recordings, add customer interview recordings or research documents as inputs. Configure the quality gates as systematic checkpoints rather than ad-hoc reviews. Build the measurement dashboard tracking your production efficiency and engagement metrics by modality. Conduct the first 30-day retrospective: what's working, what's failing, what needs configuration adjustment. At day 90, the system should be handling at least 70% of your weekly content production volume with consistent quality. The remaining 30% — high-stakes, highly customized, or channel-leading flagship content — stays primarily human-produced with AI assistance. That ratio shifts over time as team competency grows, but maintaining a human-primary category for flagship content is a deliberate quality and brand-accountability decision, not a failure of AI adoption.

- Days 1-30: Pilot one source type, one transformation chain, one operator — no automation, full documentation
- Days 1-30: Establish the time-from-source-to-published baseline and the revision rate baseline
- Days 31-60: Add image modality pairing to pilot outputs and validate visual quality standards
- Days 31-60: Build the routing layer for the pilot pipeline — automate the highest-volume manual handoff first
- Days 61-90: Add a second source type and expand to full content team operation
- Days 61-90: Build the measurement dashboard and conduct the 30-day retrospective
- Day 90: Target 70% of weekly content production volume running through the multimodal system

*The 90-day plan deliberately delays automation until day 31 — running the pilot manually first reveals the edge cases and quality gaps that would break an automated pipeline and force costly rebuilds.*

**Day 31**

is when automation of the routing layer should begin — after the pilot has validated the production flow manually and documented all edge cases

# Multimodal AI Marketing System Implementation Checklist

## Phase 1 — Foundation

---

- Audit existing AI tool usage: which tools, which team members, which content types, which handoffs are manual
- Map the three highest-volume manual handoffs in current content production
- Select pilot source type (executive recording recommended) and pilot transformation chain
- Choose tools for each modality: text AI, image AI, voice transcription
- Document the input/output specification for each transformation step
- Set production efficiency baseline: assets per producer per week, time from source to published
- Run 30-day manual pilot and document every step and edge case

## Phase 2 — Launch

---

- Add image AI pairing to the pilot pipeline and validate visual quality standards
- Build the routing layer for the pilot pipeline using Make, n8n, or documented manual protocol
- Deploy content-type-specific quality rubrics for text, image, and paired-asset coherence
- Add a second content producer and deliver pipeline training
- Configure asset tagging system: modality type and generation tool tracked from creation to publishing
- Launch measurement dashboard: production efficiency and engagement by modality combination

## Phase 3 — Optimize

---

- Expand to second source type and validate pipeline handles new input format
- Conduct 30-day retrospective: production metrics, quality rates, engagement performance

- Update system prompts, image prompt templates, and routing logic based on retrospective findings

---

- Set 90-day target: 70% of weekly content production through the multimodal system

---

- Schedule quarterly tool stack review — capabilities shift enough every 6 months to warrant reassessment

---

**NetWebMedia**

# Build the Multimodal System, Not Just the Tool Stack

NetWebMedia architects and implements multimodal AI marketing systems for B2B teams — from source-to-distribution pipeline design through tool selection, routing layer automation, and quality governance. We've built multimodal systems on Make, n8n, and the Claude Agent SDK, and we can match the architecture to your team's technical resources and operational requirements. If you're running disconnected AI tools and not seeing the compounding efficiency gains, the system design is the problem — and that's exactly what we fix.

AI Marketing Automation

AEO & AI-First SEO

Autonomous AI Agents

Paid Media + AI Creative

CRM + AI Workflows

[netwebmedia.com/contact](https://netwebmedia.com/contact)