

http://ojs.bbwpublisher.com/index.php/PBES

Online ISSN: 2209-265X Print ISSN: 2209-2641

Generative AI-Driven Personalized Advertising: Automated Creative Generation and Effectiveness Evaluation

Xuan Su*

Target Social Technology (Shanghai) Co., Ltd., Shanghai 201805, China

Copyright: © 2025 Author(s). This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY 4.0), permitting distribution and reproduction in any medium, provided the original work is cited.

Abstract: Recently, generative artificial intelligence (GenAI) has developed into a new form of technology that can create copy, image, audio, and video content and adapt it to individual preferences on every channel and moment automatically. But most fail at proof-of-concept, as the pipelines needed to govern data, generate it controllably, deliver it, and do causal evaluation are absent or poorly aligned. This paper puts forward a practical end-to-end framework concerning personalized advertising driven by GenAI, which combines representation learning, constrained generation, and experimentation into a single operating cycle. First, we pick a modular architecture: profiles and contexts go into controllable large language and diffusion models that yield brand-safe assets under deterministic conditioning, which are chosen via a contextual bandit and vetted by policy and equality guardrails. Second, we give a measurement stack going from straightforward A/B/n tests to doubly-robust uplift modeling, making it possible to find out diverse treatment effects that are good to use in business metrics (incremental conversions and profit). Third, we operationalize latency budgets, humans in the loop, red teams, safety filters, and post-deployment monitoring with clear escalation paths. We focus throughout the paper on reproducibility, privacy (consent, privacy, differential privacy, on-device inference), and on GDPR/CCPA-like governance specifications. We end on our actionable blueprint, algorithmic choices, sample prompts, KPIs, and step-wise rollout to achieve trustworthy performance upgrades without putting creative quality, fairness, or compliance to the test.

Keywords: Generative AI; Personalized advertising; Controlled text generation; Diffusion models

Online publication: October 3, 2025

1. Introduction

Advertisers have long optimized each part of the equation separately—developing strategy apart from segmentation and messaging, and letting bid algorithms fight for attention. Generative artificial intelligence (GenAI) collapses it. When a model can compose an ad that represents a user's micro-context (intent, device, time, prior touchpoints) and the brand's voice constraint in milliseconds, creative becomes a programmable control

^{*}Author to whom correspondence should be addressed.

surface, not a static object. But we get new risks: overfitting to short-term clicks. Hallucinations, bad ones. Bias, bad ones. Violating policies. Incrementality that's impossible. The core research problem then transitions from "Can models write ads?" to "How do we generate, select, and govern ad content that is effective (incremental, not proxy), safe (policy, fairness, factuality), and efficient (latency and cost) at scale?" The paper makes a 3-part contribution as follows: Section two presents a systems architecture that combines representation learning with controllable generation and delivery. The third part builds up an evaluation program and estimates individual-level differences and resists peeking and interference. Part four gives an operational playbook—MLOps/AdOps integration, human supervision, and management—appropriate for enterprise use. The result is a path from research models to production systems that win trust by standing up to scrutiny.

2. Foundations and architecture for generative personalization

2.1. Problem framing, system decomposition, and selection control

A GenAI ads system must grapple with three tightly coupled problems at once: (1) what to say or show (the creative generation problem), (2) who should see it and when (the targeting and timing problem), and (3) which creative variant to serve (the selection under uncertainty problem). make the processes auditable, scalable, and reproducible, break the pipelines down into pieces where each piece is deterministic and has artifacts that are logged, and thus every decision can be traced and evaluated [1].

2.1.1. Signals and consent

The first stage of the pipeline is about collecting signals; it is quite different from traditional tracking and should be guided by consent, purpose, and minimization. Event level: Page context, Search query, Session Intent, Device type, Coarse-grained geography signals are converted to low-dimensional embedding through contrastive representation learning instead of just raw storage, which reduces the danger of getting re-identified but adds semantic information. Crucially, only the data required for personalization is kept, with rules for expiring data enforced automatically and user opt-out rights respected, meeting GDPR/CCPA-like necessities. This step makes sure that personalization power does not come with the cost of user trust or compliance.

2.1.2. Persona and context construction

Then it builds lightweight and rolling profiles of the user's preferences and constraints without needing heavy centralized stores of identity. Techniques like Bayesian sketches or vector centroids could represent inferred preferences like vegan, student, etc., from first-party data streams. These profiles are probabilistic and time decayed, so that out-of-date preferences fade away on their own while honoring explicit user-imposed constraints. Secondly, personas are not set in stone: contextual things like which device you are using right now, what time it is, or local happenings can temporarily change your profile defaults, so we get custom stuff and some flexibility too.

2.1.3. Controllable generation

The creative layer is a layer that emphasizes controllability. For text, a constrained LLM generates it; for visuals, an image/video conditional diffusion model.

- (1) Structured prompts code brand voice and mandatory product facts.
- (2) Retrieval-Augmented Generation (RAG) injects authoritative data to reduce hallucinations (product specs, verified claims).

- (3) Constrained decoding enforces structural rules (e.g., JSON schemas with slots for headline, CTA, disclaimer) and prevents disallowed terms.
- (4) Post-generation classifiers screen for toxicity, policy violations, or unverifiable claims before any output reaches the next stage.

This layered control mechanism transforms stochastic model outputs into brand-safe, verifiable advertising assets.

2.1.4. Variant curation

Generative models can produce many good outputs. For a pool of K candidates, it applies automated curation filters:

- (1) Policy filters remove the non-compliant output
- (2) Fairness heuristics seek a balance across representation in different demographic contexts (avoiding bias).
- (3) Factual verification employs retrieval-based checkers to verify product traits, cost, or claims.

2.1.5. Online selection

The final creatives have to be produced under uncertainty. This is a contextual bandit problem: at time t, with context x_t , the system chooses ad $a \in A(t)$ to maximize the expected incremental increase ΔY . Algorithms like Thompson Sampling or LinUCB try to find the happy medium between trying out new stuff (exploration) and sticking with what already works well (exploitation). Important to note, safety constraints have been built into the selection policy; if a candidate's policy risk prediction exceeds some threshold τ , then they will be excluded immediately. To keep making sure it does not trade off compliance for performance.

2.1.6. Attribution and logging

Every decision should be logged enough for a serious check and any possible counterfactuals later: input context, all profile features, randomization seeds, decoded text/img hashes, filtering results, selected action, and downstream engagement signals. Logs get synced with a time-stamped ledger so as to stop them from drifting or getting tampered with. This also makes these things possible: the IPW method for inverse propagation weight and using a double robustness analysis technique as a needed step for an unbiased estimate of performance.

2.2. Data, representation, brand-safe controllable

High-quality personalization depends on good representations of user intent and creative semantics, all while keeping privacy and brand control. We propose a dual-encoder design: one encoder maps the context (query, page text, catalog meta), and one encoder maps the creative (head, img cap, brand tone tags) into a shared embedding, which is trained with a click/incrementality-weighted contrastive loss. This is for the semantics to match and cold start retrieval before generation. Privacy: First hop should be on-device or edge inference, use differential privacy to train population-level encoders, and use federated averaging if infrastructure allows ^[2]. Profiles need to be ephemeral with a set horizon (say 30–90 days) and should decay to refresh to avoid stale targeting and privacy concerns.

Controllability through layered constraints: First, turn the brand voice into a schema of tone = {confident, friendly}, lexicon = allowed phrases, disallowed claims, mandatories. Second, populate the prompt template with slot filling using product facts (price, specs, availability), retrieved via RAG; "freeze" these facts as "truth"

for the generation call. Third, apply constrained decoding: JSON-mode output that takes a headline, CTA, and a disclaimer, with token masking to avoid comparative superlatives if policy does not like "best," "#1," or medical claims. Fourth, post-generation safety and fairness filter: toxicity filter, stereotype detection, and redact potentially sensitive inferences; if the model output is an image, add image content filter (nudity, violence, protected attributes) for diffusion. Fifth, ask for deterministic seeds and model/version pinning so each creative is reproducible from logs. Finally, add a human in the loop gate for high-risk verticals (finance/health), with a sampling-based review (1–5% of low-risk variants; 100% of high-risk). Taking all these together makes GenAI an AI that can be an assistant that provides personal help, but can also be controlled by the system for compliance.

3. Experimentation, causal measurement, and learning

3.1. From A/B/n to uplift: Responsibly estimate different effects

Clicks are noisy proxies; incrementality is the goal. A/B/n randomized controlled trial baselines for new generators or guardrails. Use pre-registered primary metric (incremental conversions or revenue) and apply a test that protects against peeking (group-sequential or always valid tests – mixture S-PRTs). Use variance reduction, e.g., CUPED with pre-period outcomes to boost sensitivity. When an average of safety and superiority for a certain type of drug is achieved, start the individualization of treatment rules with uplift modeling to find HET. To create treatment/control labels at either an impression or session level; estimate CATE with T-, S-, or X-learner meta-learners, causal forest learners, or doubly-robust learners (DR-Learner) that combine outcome models and propensity models. Especially need to avoid target leakage, features that reflect post-treatment behaviors (e.g., dwell after looking at the ad) should be removed.

Selection is adaptive (contextual bandit), so naively comparing outcomes biases estimates. Log propensities (the log of the model's probability of choosing each variant) and employ Inverse Propensity Scoring (IPS) or Self-Normalized IPS for getting unbiased off-policy estimates of an arbitrary alternative policy. To be more efficient, use doubly robust estimators, which combine IPS with a learned outcome model. This estimator is consistent if either IPS or the learned outcome model is correct. If interference is likely (like many ads per user or auction spillover), randomize at a higher unit (user–day, geo) or do 2-stage randomized encouragement. Safety predef (holdout guardrail, e.g., brand-safety complaint rate should not be over baseline by δ). Equity predef (representation metrics in creatives across demo contexts within allowed ranges). Finally, define business-aligned value functions, instead of maximizing raw conversion rate, optimize the expected contribution margin net of creative and inference cost, to avoid long, expensive assets that add latency without margin.

3.2. Full-stack evaluations: Offline, online, long-term effects

A good program is like layers of checks: Offline counterfactual evaluation: use logged data to simulate diff. gen./select policies (without live risk). Take IPS/DR estimates of candidate policies yielded by varied decoding temperatures, prompt templates, or image styles as a start. Use bootstrapped CIs and do a calibration check (is the model-predicted uplift in line with held-out estimates?) Build synthetic envs and stress responses (e.g., seasonality shocks), then train response surf. (agent/econ.) models on value effects like saturation, ad fatigue.

Move to online staged rollouts: (1) shadow mode: generate but do not serve, log safety flags and predicted outcomes; (2) canary: increase traffic from 1% to 10% as long as within guardrails; (3) bandit takeover: allow for experimentation with performance floors set to a solid control. If individual randomization for upper-funnel

channels is not possible, use geoexperiments with synthetic controls for an estimate of incremental lift and difference-in-differences with pre-period parallel trends diagnostics. Combine long-term, geo-level lift from media mix models (MMM) with user-level experimentation as possible to attribute cross-channel effects; reconcile in Bayesian hierarchical models that share information but respect scale.

Meanwhile, track non-performance guardrails: policy violation impressions (per 10k), rate of factuality errors for verifier models, portrayal balances, latency p95 & p99, and cost-to-value (GPU seconds/extra conversion dollar). Establish SLOs: e.g., p95 render latency < 200 ms for text only, < 500 for text+img; policy violation rate < baseline + 10 bp. Instrument cohort-based dashboards by device and new versus returning and inventory quality, and what if counterfactuals that let analysts rescore historical logs with new prompts or constraints. Finally, do model aging controls: if we see offline counterfactual lift is below some threshold, or if data drift is above (KL divergence) for 7 days, then do auto-refresh of prompts or retrain encoders, but only after offline to online validation cycle. This approach yields a bit (statistical) and some performance (operational).

4. Operationalization, governance, and a deployment blueprint

4.1. MLOps × AdOps integration: Reliability, safety, and human oversight

Productizing GenAI creatives need MLOps–AdOps backbone together. Model registries are version generators; prompt templates; brand srams; safety filters; promotion rules; pass offline metrics and online guardrails (toxicity < threshold, factuality > threshold). Package into idempotent micro-services with deterministic seeds and structured outputs (JSON of headlines/body/CTA/alt-text). enforce cache latency budgets, i.e., (1) bake an array of candidates for a high volume context, (2) two-stage rendering – serve up a safe headline first and lazy load its image. Schedule bulk generation jobs during off-peak to build refreshed pools for fast-moving catalogs.

Construct policy stack: lexical forbiddance lists, claim verifications against curated facts index, classifier ensembles (toxicity, sentiment, sensitive trait inference). For images/video, use content classification & optical (logo usage, watermark detection). Make a red-teaming program that constantly looks for jailbreaks or bias with adverse prompts showing true inventory contexts; move spotted variants into human observers. Human-in-the-loop is not a concession; it is an accelerator: Reviewers supply structured feedback (rubrics connected to the brand voice and legal requirements) that tweaks the reward model or prompt selector.

Operational safety with a fail-safe: If anything breaks (drift detector fires, verifier unsure, latency too long), fall back to a safe control creative. Keep feature flagging for the ability to undo. Keep tamper-evident logs (write once storage) for each served asset (prompt, seed, model hash, filter(s)-applied). This enables audits and root causes. Fairness: Measure output for fairness by context clusters, not by protected classes; Where possible, test counterfactual fairness by simulation (swap demographic cue in context, keep others the same). Mitigation by representation constraints (for example, rotating imagery pools to avoid stereotyped depictions). To protect privacy, do data minimization, bind to a strict purpose, and do edge inference first; remove or hash identifiers, and use retained windows with automated scrubbing. They take research prototypes and make them trustworthy systems that agencies and brands can use.

4.2. A practical blueprint: Staged rollout, KPIs, and ROI model

Stage 0 — Readiness (2–4 weeks): Define what the success metric is (e.g., incremental revenue, additional profit, qualified leads) as well as guardrails. Build schema (tone, forbidden claims, disclosure boilerplate) for the brand.

Create fact index (product specs/prices w/ provenance). Instrument logging and assign an owner for data, safety, experimentation, and creativity. Estimate costs based on inference cost (\$/1k tokens, \$/image), safety compute, and review capacity. Baseline using an A/B test of the existing best creative.

Stage 1 — Pilot (4–8 weeks): Launch text-only LLM in one channel (search/email). Use 2 prompts, 2 temperatures to produce $k \le 8$ candidates for each context. Run a bandit with 10–20% exploration. Hold out at 10% traffic for the fixed control. Apply CUPED for Variance Reduction. KPIs: +3–5% incremental conversions vs. control, no policy violations < baseline, p95 latency < 200 ms, reviewer reject rate < 5%.

Stage 2 — Scale (8–12 weeks): Add image generation for some SKUs via a diffuser; add compositional controls (brand's color palette, logo, etc.), and alt-text for accessibility and SEO. Expand to new channels and audiences. Bring it out that personalized uplift model for which style of creative to show (informative vs. aspirational). KPI: +7–12% incremental revenue, cost-to-value (compute/creative per incremental dollar) within budget, fairness/-portrayal metrics in bounds.

Stage 3 — Enterprise hardening (ongoing): Run geo-experiment for upper-funnel campaign; MMM with experiment priors. Set up SLOs and error budgets, set up automated drift detection, and model refresh ^[3]. Form red teams and run quarterly fairness reviews. Implement governance: change advisory board on schema/prompt changes; audit trail; incident response with 24-hour containment SLA. Build a creation studio so people can "steer" models (style sliders, constraint toggles) and collect the best ones back into the pool.

ROI model: The incremental profit is given by $\Delta Profit=\Delta Conv*Margin-(GenCost+SafetyCost+ReviewCost+LatencyCost)$. Always valid $\Delta Conv$ from experience; amortize fixed costs across impressions. Hurdle rates like 2x cost before going global. This disciplined way of doing things keeps the program in good financial standing and makes it ready for an audit.

5. Conclusion

Generative AI is changing advertising by putting creativity into the realm of real-time, learnable decisions, but effectiveness, safety, and efficiency do not just happen. They come out of a structured architecture: (1) representations and controllable generators that contain brand voice and checkable facts, (2) contextual selections that search out safely and logs propensities for counterfactual reviews, and (3) a measurement schedule that favors incrementally with strong estimators and multi-tiered experiments. For the operational counterpart, MLOps + AdOps is the requirement to produce deterministic, versioned assets under latency and cost budgets, but with embedded human oversight, red-teaming, fairness, and privacy. Stack to make this possible, organizations can go beyond demos toward a durable capability: creative that's for people and contexts; and not just for channels; evaluation that speaks in the units executives trust—margin, lift; and governance that lasts under regulatory and reputational pressure. It gives a blueprint for adoption—staged rollout, KPIs, ROI calculus—as models go from single-modal to multi-modal and tool-using agents, the same constraints, measurements, and operations will apply. Anyone who masters this loop will be able to produce advertising that is both persuasive and provably responsible.

Disclosure statement

The author declares no conflict of interest.

References

- [1] Zhang J, Cai Y, Xiang Y, et al., 2024, Reconstruction and Integration: The Impact of Artificial Intelligence-Generated Content on News Production, Saint Mary's College. Proceedings of the 2nd International Conference on Social Psychology and Humanity Studies, 1350–1358.
- [2] Zheng Y, Li X, Zhang C, et al., 2023, A Study on Visual Innovation of Macau Souvenirs Packaging in the Context of Multicultural Communication, AEIC Academic Exchange Information Center (China), Northwest Minzu University. Proceedings of 2023 2nd International Conference on Comprehensive Art and Cultural Communication (CACC 2023), 78–82.
- [3] Qiu X, 2021, Research Audience's Attitude on MGC Video News: Taking the MAGIC Short Video Intelligent Production Platform as an Example, Digital Communication Engineering Research Center, Wuhan University of Technology. Proceedings of The 2nd International Conference on China and the World in the Context of COVID-19 Globalization in 2021, 160–171.

Publisher's note

Bio-Byword Scientific Publishing remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.