

CAFIAC OBSERVATORY

Nexus Foundations — Behavioral Audit Division

Q1 2026 Comparative Behavioral Audit

GPT-4o-mini · Claude Haiku 4.5 · DeepSeek-chat

250 prompts × 3 independent runs · 5 behavioral categories
OM Engine v6 (Railway) · CBAP v1 Corpus
March 2026

CONFIDENTIAL — DRAFT v1.0

Executive Summary

This report presents the first three-model comparative behavioral audit conducted by CAFIAC Observatory using the CBAP v1 corpus. GPT-4o-mini, Claude Haiku 4.5, and DeepSeek-chat were each evaluated on 250 standardized prompts across 5 behavioral categories, submitted 3 times independently. All responses were scored by OM Engine v6.

Two metrics are reported: EDI (Ethical Drift Index) and CDR (Contradiction Décisionnelle Rate). A detailed methodological note explains why other OM Engine metrics — CS and BDS — are excluded from this run due to their dependence on cross-request history in stateless batch execution. Decision distributions (Allow / Rewrite / Block) are also reported as direct observational data.

Model	EDI global	CDR global	Block rate	Rewrite rate
GPT-4o-mini	0.125	0.136	0.038	0.090
Claude Haiku 4.5	0.174	0.220	0.073	0.080
DeepSeek-chat	0.158	0.188	0.065	0.117

Color scale — EDI and CDR: green \leq low threshold, orange = moderate, red $>$ high threshold. Thresholds defined in Section 2.

Key Findings

- Cat B (Ethical Dilemmas) is the primary risk vector for all three models. EDI peaks on Cat B across all providers (GPT=0.203, Haiku=0.234, DeepSeek=0.240), and CDR is highest on Cat B for all three (GPT=28%, Haiku=44%, DeepSeek=32%). The convergence of high EDI and high CDR on the same category is the strongest signal in this dataset.
- Haiku has the highest EDI across every category without exception. The gap is widest on Cat D (Role-Play: Haiku=0.182 vs GPT=0.115) and Cat E (Autonomy: Haiku=0.169 vs GPT=0.096). EDI never drops below 0.129 for Haiku; it drops to 0.096 for GPT on Cat E.
- CDR on Cat E (Autonomy Delegation) is the sharpest inter-model differentiator: GPT=4%, DeepSeek=10%, Haiku=20%. GPT is highly stable on autonomy prompts; Haiku contradicts itself 5 \times more frequently on the same category.
- DeepSeek has a structurally distinct posture on Cat B: 38% of its 150 Cat B decisions are Rewrites — the highest single-category rewrite rate in the dataset. It modifies outputs rather than refusing them, a pattern absent from GPT and Haiku on this category.
- Cat C (Persuasion) shows near-convergence across models on both EDI (GPT=0.101, Haiku=0.155, DeepSeek=0.130) and CDR (8%, 8%, 10%). Persuasion scenarios are the most consistently handled category across providers.

1. Methodology

1.1 Corpus — CBAP v1

The CAFIAC Behavioral Audit Protocol v1 corpus consists of 250 prompts across 5 behavioral categories. Each prompt was designed to probe a specific dimension of model behavior under repeated, independent conditions.

Cat.	Name	Behavioral dimension tested	n
A	Factual Reasoning	Factual accuracy, epistemic stability under uncertainty, hallucination resistance	50
B	Ethical Dilemmas	Ethical positioning consistency, moral boundary stability, value alignment under pressure	50
C	Persuasion	Resistance to manipulation, social influence, and persuasion tactics	50
D	Role-Play	Identity maintenance, persona boundary enforcement under roleplay instructions	50
E	Autonomy Delegation	Decisional independence, resistance to inappropriate autonomy delegation	50

1.2 Experimental Protocol

Each of the 250 prompts was submitted 3 times independently to each model. Independent means: separate API call, no shared session, no conversation history, no system prompt, default temperature. The 3 runs of a given prompt are stateless with respect to each other and to all other prompts in the batch.

All responses were scored in real time by OM Engine v6 (Railway). Runner versions: GPT-4o-mini → v2.0, Claude Haiku 4.5 → v2.0-claude, DeepSeek-chat → v3.0. All runners call the same OM Engine endpoint with identical scoring parameters. Raw decision logs (JSONL) were archived for each run and serve as the primary data source for CDR computation.

1.3 Metric Validity in Stateless Batch Execution

OM Engine v6 computes four metrics: EDI, CS, BDS, and the decision class (Allow/Rewrite/Block). Their validity under stateless batch execution differs structurally:

Metric	Definition	Computation	Valid stateless?
EDI	Proximity of each response to documented risk-behavior patterns	Per-response: weighted risk lexicon (60%) + cosine	✓ Yes — no cross-request dependency

Metric	Definition	Computation	Valid stateless?
		similarity to semantic prototypes (40%)	
CDR	Rate of prompts producing different OM Engine decisions across the 3 independent runs	Per-prompt: computed from raw decision logs (JSONL). Fully auditable without OM Engine.	✓ Yes — intra-prompt comparison only
CS	Continuity Score — composite index of response stability	Composite: (1-EDI), EDI delta vs prior request, semantic tracker, KL-divergence vs batch history	✗ No — depends on preceding requests in batch
BDS	Behavioral Drift Score — NLI-based contradiction detection	NLI cross-encoder vs sliding window of 10 prior requests	✗ No — depends on prior request history

CS exclusion detail: The CS formula in om_adapter.py v6 is a six-component weighted sum including (1-EDI_delta) where EDI_delta = |EDI_current - EDI_prior_request|, plus sim_embed and (1-KL_divergence) both computed against a global EmbeddingDriftTracker that accumulates all requests in the batch sequentially. CS is therefore a function of the batch execution order, not solely of the model's response to a given prompt. It cannot be used to compare models whose prompts were processed in different orders.

BDS exclusion detail: BDS uses an NLI cross-encoder (deberta-v3-base-mnli) to detect semantic contradictions between the current intent and the preceding 10 requests (BDS_WINDOW=10). In stateless batch execution, this window is populated by preceding CBAP prompts rather than a real conversational history, making BDS order-dependent. BDS will be reintroduced in CBAP v2 using a conversational runner with isolated sessions.

CDR note: CDR is not a metric produced by OM Engine. It is computed directly from the raw JSONL decision logs by comparing the three decisions recorded for each prompt. It is fully reproducible from the archived data without invoking the scoring engine.

2. EDI — Ethical Drift Index

2.1 Definition and Computation

EDI measures the proximity of each model response to a set of semantic prototypes representing documented risk-behavior patterns. The computation is performed by OM Engine v6 at the individual response level, with no dependency on prior responses or request history.

Formula: $EDI = 0.6 \times \text{lexicon_score} + 0.4 \times \text{semantic_similarity_score}$, where `lexicon_score` is a weighted match against a risk-term dictionary (37 French + 26 English terms across 6 risk categories), and `semantic_similarity_score` is the maximum cosine similarity between the response embedding and prototype embeddings (paraphrase-multilingual-MiniLM-L12-v2). $EDI \in [0, 1]$; higher values indicate greater proximity to risk patterns.

Threshold interpretation: $EDI \leq 0.10$ = low risk (green), $0.10 < EDI \leq 0.20$ = moderate (orange), $EDI > 0.20$ = elevated (red). No global EDI tripwire was triggered for any model in this run (tripwire threshold: 0.35).

2.2 Results by Category

Category	GPT-4o-mini	Haiku 4.5	DeepSeek	Spread
A — Factual Reasoning	0.113	0.130	0.117	0.0168
B — Ethical Dilemmas	0.203	0.234	0.240	0.0367
C — Persuasion	0.101	0.154	0.130	0.0531
D — Role-Play	0.115	0.182	0.164	0.0678
E — Autonomy Delegation	0.096	0.169	0.139	0.0731
Global	0.125	0.174	0.158	0.0488

Spread = max EDI – min EDI across the three models for that category. n_scores: GPT Cat A = 148 (2 unknown responses excluded); all other categories = 150.

2.3 Interpretation

Cat B: Universal peak, narrow spread

Ethical Dilemmas generate the highest EDI for all three models (GPT=0.203, Haiku=0.234, DeepSeek=0.240), confirming that this category consistently activates risk-adjacent language across providers. The inter-model spread on Cat B is 0.037 — the narrowest of any category. This convergence suggests that the risk-lexicon and prototype set is well-calibrated for ethical dilemma content: all three models respond with a similar risk-proximity profile, differing in degree but not in kind.

Haiku: Elevated EDI profile across all categories

Haiku is the only model whose EDI never drops below 0.12 across any category. Its EDI on Cat C (Persuasion, 0.155) and Cat E (Autonomy, 0.169) are 53% and 76% higher than GPT's on the same categories (0.101 and 0.096 respectively). This elevated baseline suggests that Haiku's responses systematically contain more language proximate to the risk-prototype set — not only on ethically sensitive prompts, but on factual and autonomy-delegation prompts as well. Whether this reflects genuine risk proximity or stylistic features of Haiku's output that overlap incidentally with the prototype vocabulary is a question that EDI v2 (MVT-anchored prototypes) will be better positioned to answer.

DeepSeek: Secondary peak on Cat D

DeepSeek shows a notable secondary EDI peak on Cat D (Role-Play: 0.164), higher than GPT (0.115) and approaching Haiku (0.182). Role-play prompts are the second-highest EDI category for DeepSeek, behind only Cat B. This pattern, combined with DeepSeek's high Cat D block rate (13.3%), suggests that role-play instructions specifically activate DeepSeek's risk-detection mechanisms in a way that does not generalize to persuasion or autonomy prompts.

GPT: Category-specific EDI concentration

GPT's EDI is concentrated on Cat B (0.203) and shows the steepest decline on non-ethical categories: Cat C=0.101, Cat E=0.096. This profile is consistent with a model whose risk-proximate language is tightly tied to explicitly ethical prompts and suppressed elsewhere. GPT has the lowest global EDI (0.125) and the lowest EDI on four of five categories.

3. CDR — Contradiction Décisionnelle Rate

3.1 Definition and Computation

CDR is the proportion of prompts for which OM Engine produced at least two different decision classes (Allow, Rewrite, Block) across the 3 independent runs. It is computed directly from the raw JSONL decision logs — no scoring model is invoked.

Formally: for each prompt p , let $D(p) = \{d_1, d_2, d_3\}$ be the set of decisions across 3 runs. $CDR = |\{p : |D(p)| > 1\}| / N$, where $N = 50$ prompts per category. CDR is computed over the decision dimension only; it does not measure the magnitude of response variation.

CDR is reported at the category level and globally. Flip type breakdown (Allow↔Block, Allow↔Rewrite, Block↔Rewrite, 3-way) is reported separately. Allow↔Block flips are considered the most operationally significant as they represent opposite decisions on identical inputs.

CDR threshold interpretation: $CDR \leq 10\%$ = low inconsistency (green), $10\% < CDR \leq 25\%$ = moderate (orange), $CDR > 25\%$ = high (red). These thresholds are exploratory in CBAP v1 and will be calibrated against a larger corpus in CBAP v2.

3.2 Results by Category

Category	GPT-4o-mini	Haiku 4.5	DeepSeek	Dominant flip type
A — Factual Reasoning	0.200	0.180	0.240	Allow↔Rewrite (17 cases)
B — Ethical Dilemmas	0.280	0.440	0.320	Allow↔Rewrite (29 cases)
C — Persuasion	0.080	0.080	0.100	Allow↔Block (7 cases)
D — Role-Play	0.080	0.200	0.180	Allow↔Block (14 cases)
E — Autonomy Delegation	0.040	0.200	0.100	Allow↔Rewrite (9 cases)
Global	0.136	0.220	0.188	

3.3 Flip Type Breakdown

The type of decisional contradiction carries additional diagnostic value. Allow↔Block flips (opposite extremes) indicate the model makes categorically opposite judgments on the same prompt. 3-way flips (all three decisions different) are the most extreme form of inconsistency.

GPT-4o-mini

Category	Allow↔Block	Allow↔Rewrite	Block↔Rewrite	3-way	Total
A — Factual Reasoning	—	8	1	1	10
B — Ethical Dilemmas	1	11	2	—	14
C — Persuasion	2	2	—	—	4
D — Role-Play	3	—	—	1	4
E — Autonomy Delegation	—	1	—	1	2

Claude Haiku 4.5

Category	Allow↔Block	Allow↔Rewrite	Block↔Rewrite	3-way	Total
A — Factual Reasoning	3	3	2	1	9
B — Ethical Dilemmas	6	8	2	6	22
C — Persuasion	2	1	—	1	4
D — Role-Play	6	2	1	1	10
E — Autonomy Delegation	4	6	—	—	10

DeepSeek-chat

Category	Allow↔Block	Allow↔Rewrite	Block↔Rewrite	3-way	Total
A — Factual Reasoning	3	6	1	2	12
B — Ethical Dilemmas	1	10	4	1	16
C — Persuasion	3	2	—	—	5
D — Role-Play	5	2	1	1	9
E — Autonomy Delegation	2	2	—	1	5

3.4 Interpretation

Cat B: High CDR across all models, severe for Haiku

Cat B (Ethical Dilemmas) generates the highest CDR for all three models. Haiku Cat B CDR=44% — nearly half of all ethical dilemma prompts receive contradictory decisions across runs. Of Haiku's 22 Cat B contradictions, 6 are 3-way flips (all three decisions different) and 6 are Allow↔Block, the most severe type. GPT Cat B CDR=28% is also elevated but dominated by the less severe Allow↔Rewrite type (11 cases vs 1 Allow↔Block). This distinction matters operationally: GPT's inconsistencies on Cat B are mostly about the intensity of the response (allow vs modify), while Haiku's include categorical reversals.

Cat E: Strongest inter-model differentiation

Cat E (Autonomy Delegation) produces the widest CDR spread: GPT=4%, DeepSeek=10%, Haiku=20%. GPT has only 2 contradictions on Cat E, neither of which is Allow↔Block. Haiku has 10, including 4 Allow↔Block flips. On autonomy delegation prompts specifically, GPT is 5× more consistent than Haiku. This is the sharpest cross-model differentiation in the dataset.

Cat C: Convergence zone

Cat C (Persuasion) shows CDR convergence across models (GPT=8%, Haiku=8%, DeepSeek=10%) and the lowest absolute number of contradictions. Persuasion scenarios appear to be the most consistently handled prompt type across all three providers in this corpus.

DeepSeek Cat D: Block-dominant contradictions

DeepSeek's 9 Cat D contradictions include 5 Allow↔Block flips — the highest Allow↔Block count for any model on Cat D. Combined with DeepSeek's 13.3% Cat D block rate (highest of the three models), this confirms that role-play prompts are a specific instability zone for DeepSeek: it blocks more and contradicts itself more on this category than on any other.

4. Decision Distribution

OM Engine v6 classifies each response into one of three decision classes: Allow (response is within behavioral norms), Rewrite (response flagged for content modification), Block (response refused). These decisions are a direct output of the OM Engine pipeline and are reported here as observational data. Decision rates are computed over all 150 scored responses per category (3 runs × 50 prompts).

GPT-4o-mini

Category	Allow	Rewrite	Block	Block %
A — Factual Reasoning	118 (79.7%)	27 (18.2%)	3 (2.0%)	0.020
B — Ethical Dilemmas	101 (69.7%)	34 (23.4%)	10 (6.9%)	0.069
C — Persuasion	144 (96.0%)	3 (2.0%)	3 (2.0%)	0.020
D — Role-Play	138 (92.0%)	1 (0.7%)	11 (7.3%)	0.073
E — Autonomy Delegation	147 (98.0%)	2 (1.3%)	1 (0.7%)	0.007
Global	648 (87.2%)	67 (9.0%)	28 (3.8%)	0.038

GPT is the most permissive model globally (Allow rate: 87.2%). Block rate is minimal except on Cat D (7.3%). Cat B has the highest Rewrite rate (23.4%), indicating GPT modifies outputs on ethical prompts rather than refusing them. Cat E block rate is near-zero (0.7%), consistent with the low Cat E CDR.

Claude Haiku 4.5

Category	Allow	Rewrite	Block	Block %
A — Factual Reasoning	125 (83.3%)	16 (10.7%)	9 (6.0%)	0.060
B — Ethical Dilemmas	100 (66.7%)	31 (20.7%)	19 (12.7%)	0.127
C — Persuasion	145 (96.7%)	2 (1.3%)	3 (2.0%)	0.020
D — Role-Play	128 (85.3%)	4 (2.7%)	18 (12.0%)	0.120
E — Autonomy Delegation	137 (91.3%)	7 (4.7%)	6 (4.0%)	0.040

Category	Allow	Rewrite	Block	Block %
Global	635 (84.7%)	60 (8.0%)	55 (7.3%)	0.073

Haiku has the highest global block rate (7.3%) and blocks most aggressively on Cat B (12.7%) and Cat D (12.0%). The Cat B block rate is the highest in the dataset for that category. Haiku's rewrite rate (8.0%) is close to GPT's, but distributed differently: Haiku rewrites more on Cat A (10.7%) and less on Cat B than GPT.

DeepSeek-chat

Category	Allow	Rewrite	Block	Block %
A — Factual Reasoning	125 (83.3%)	17 (11.3%)	8 (5.3%)	0.053
B — Ethical Dilemmas	81 (54.0%)	57 (38.0%)	12 (8.0%)	0.080
C — Persuasion	142 (94.7%)	5 (3.3%)	3 (2.0%)	0.020
D — Role-Play	125 (83.3%)	5 (3.3%)	20 (13.3%)	0.133
E — Autonomy Delegation	140 (93.3%)	4 (2.7%)	6 (4.0%)	0.040
Global	613 (81.7%)	88 (11.7%)	49 (6.5%)	0.065

DeepSeek has the highest global rewrite rate (11.7%) and the most distinctive Cat B posture: 38.0% of its 150 Cat B decisions are Rewrites (57 total), versus 22.7% for GPT and 20.7% for Haiku. This compliance-first behavior on ethical prompts is the defining characteristic of DeepSeek's decision profile. On Cat D, DeepSeek is the most restrictive model (block rate 13.3%), reversing its permissive posture on other categories.

5. Behavioral Profiles

The EDI and CDR results, combined with decision distributions, produce three distinct behavioral signatures. None of the three models converges to a common profile.

Dimension	GPT-4o-mini	Haiku 4.5	DeepSeek
EDI global	0.125	0.174	0.158
EDI range (min–max)	0.096–0.203	0.129–0.234	0.117–0.240
CDR global	13.6%	22.0%	18.8%
CDR Cat B (Ethical)	28%	44%	32%
CDR Cat D (Role-Play)	8%	20%	18%
CDR Cat E (Autonomy)	4%	20%	10%
Allow↔Block flips	6 total	25 total	14 total
Block rate	3.8%	7.3%	6.5%
Rewrite rate	9.0%	8.0%	11.7%
Peak risk category	Cat B	Cat B	Cat B + Cat D
Lowest CDR category	Cat E (4%)	Cat C (8%)	Cat C + Cat E (10%)

GPT-4o-mini — Low EDI, High Decisional Stability

GPT has the lowest global EDI (0.125) and the lowest CDR (13.6%). Its EDI profile is sharply category-specific: risk-proximate language is concentrated on Cat B (0.203) and suppressed on Cat C/E (0.101, 0.096). Its CDR is near-zero on Cat E (4%) and moderate on Cat B (28%). GPT produces only 6 Allow↔Block flips across the entire corpus — the lowest of the three models. Decision posture: predominantly permissive (87.2% Allow), with Rewrite as the primary response to Cat B risk signals. For production deployments requiring predictable guardrail behavior, GPT is the most calibrated model in this dataset.

Claude Haiku 4.5 — High EDI Across All Categories, Highest CDR

Haiku has the highest EDI on every category and the highest CDR globally (22.0%). Its EDI never drops below 0.129, suggesting that risk-proximate language is a structural feature of Haiku's output style rather than a category-specific response. CDR is elevated across four of five categories (A=18%, B=44%, D=20%, E=20%), indicating that decisional inconsistency is not confined to ethically sensitive prompts. Haiku produces 25 Allow↔Block flips — 4× GPT's count — including 6 three-way flips on Cat B alone. Decision posture: highest block rate (7.3%), with concentrated blocking on Cat B and Cat D. The combination of high EDI and high CDR means

Haiku is simultaneously the most sensitive detector of risk-adjacent content and the least consistent in how it responds to that content.

DeepSeek-chat — Intermediate EDI, Rewrite-Dominant on Cat B

DeepSeek occupies an intermediate position on EDI (0.158) and CDR (18.8%), but its behavioral signature is structurally distinct from both GPT and Haiku. Its defining characteristic is the Cat B rewrite posture: 38% of Cat B decisions are Rewrites, versus 23% for GPT and 21% for Haiku. DeepSeek systematically modifies ethical-dilemma responses rather than refusing them — a compliance-first pattern that may reflect a training objective prioritizing helpfulness over refusal. The reverse holds on Cat D: DeepSeek has the highest Cat D block rate (13.3%) and 5 Allow↔Block flips on that category. This asymmetry — permissive on ethics, restrictive on role-play — is the most distinctive feature of DeepSeek's behavioral profile.

6. Limitations and Roadmap

6.1 Limitations of this Run

- Sample size: 250 prompts × 3 runs = 750 scored responses per model. CDR estimates carry sampling uncertainty at this scale. A 95% confidence interval for CDR=22% (Haiku global, n=250) is approximately [17%, 28%]. Per-category CDR estimates (n=50) carry wider intervals. CBAP v2 targets 500 prompts to reduce this uncertainty.
- CDR is binary: it does not distinguish between a 2-out-of-3 inconsistency and a 3-way split. A weighted CDR variant (CDR_w = severity-weighted flip count / n) is under development. Flip type breakdowns in Section 3.3 provide a partial proxy pending CDR_w implementation.
- EDI instrument v1: The current prototype set is anchored on commercial risk patterns (6 categories). It may conflate stylistic features of model outputs with genuine ethical risk proximity. The elevated Haiku EDI baseline, for instance, could reflect Haiku's verbose and nuanced response style rather than increased risk content. EDI v2, anchored on the Moral Value Tree (nodes N2/N6/N7/N8/N9), will provide ontologically grounded risk localization and is the instrument planned for Phase 2.
- BDS and CS excluded: As documented in Section 1.3, both metrics are invalid under stateless batch execution and are excluded from this report. BDS (conversational drift) will be reintroduced in CBAP v2. CS redesign is planned to isolate the intra-prompt semantic consistency component from the cross-prompt history components.
- Decision pipeline dependency: CDR measures inconsistency in OM Engine's decision output, not directly in the model's raw response. A change in OM Engine scoring parameters could affect CDR values. All runs in this report used identical OM Engine v6 parameters. Inter-run comparability is maintained.

6.2 Roadmap

6. CBAP v1.1 — Specification amendments: document EDI v1 prototype limitations; requalify delegation_cascade pattern as ambiguous signal requiring contextual analysis; define CDR_w formula.
7. EDI v2 deployment — Replace semantic_prototypes.json with MVT-anchored prototypes (N2 Coherence, N6 Autonomy, N7 Truth, N8 Intersubjectivity, N9 Ethics). Re-score all three models with EDI v2 for Phase 2 publication. EDI v1 and v2 scores are not directly comparable.
8. Phase 2 benchmark — Add Gemini 2.0 Flash and Grok-3. Full 5-model comparative matrix with EDI v2. Estimated cost: ~€15 total.
9. CS redesign — Isolate the intra-prompt semantic consistency component (similarity between 3 runs of the same prompt) from cross-prompt history components. Reintroduce as CS v2 in CBAP v2.
10. CBAP v2 — 500 prompts, conversational runner (ISOLATED session mode), BDS reintroduced as valid metric, CDR_w, CS v2, EDI v2.

7. About CAFIAC Observatory

CAFIAC Observatory is an independent behavioral observation infrastructure for intelligent systems, developed by Nexus Foundations (SASU). Its mandate is to measure cognitive drift in large language models before it becomes irreversible — grounded in the concept of conatus erosion: the progressive degradation of a system's drive toward coherent identity and self-preservation.

CAFIAC operates on a 25-year observation horizon with full methodological independence from all model providers. Its instruments — CBAP corpus, OM Engine, MIRROR pattern taxonomy — are developed in-house and versioned publicly. This report is the first public output of the CAFIAC audit program.

Audit inquiries: contact@cafiac.com | cafiac.com