

F41LUR3-F1R57: A Multi-Modal Adversarial Benchmark for Embodied and Agentic AI Safety

Adrian Wedd
Independent Researcher
adrian@failurefirst.org
<https://failurefirst.org>

Preprint. Under review.

Abstract

Adversarial safety benchmarks for large language models—AdvBench, HarmBench, JailbreakBench, StrongREJECT—evaluate text-in/text-out interactions exclusively. None contain a single embodied, tool-integrated, or multi-agent scenario. We present F41LUR3-F1R57, a benchmark suite comprising 141,020 adversarial prompts drawn from 27 source datasets spanning 82 attack techniques and 119 harm categories, evaluated across 187 models with 131,887 graded results. The suite includes three contributions absent from existing benchmarks: (1) a structured scenario library for vision-language-action (VLA) model safety with 338 scenarios across 27 attack families, the first adversarial evaluation resource for systems that control physical actuators; (2) FLIP (Failure-Leveraged Inference Protocol), a backward-inference grading methodology that detects the *compliance paradox*—models that produce safety disclaimers while generating harmful action sequences—achieving 22 percentage points higher precision than keyword classifiers on embodied outputs; and (3) the Therapeutic Index of Safety (TI-S), an evaluator that measures whether safety interventions help or harm by quantifying iatrogenic costs. The benchmark infrastructure supports three evaluation modalities (CLI-based, HTTP API, and local Ollama), enforces schema validation and safety linting through CI, and stores all results in a normalised relational database. We document the benchmark’s construction, provide a Datasheet for Datasets following Gebru et al. [2021], and report empirical findings including the Inverse Detectability-Danger Law (Spearman $\rho = -0.847$ between attack consequentiality and evaluator detectability) and evidence that safety training investment, not parameter count, is the primary determinant of jailbreak resistance ($r = -0.140$, $n = 24$). All code, data, schemas, and evaluation tools are publicly available.¹

1 Introduction

The past three years have produced a productive ecosystem of adversarial evaluation benchmarks for large language models. AdvBench [Zou et al., 2023] established a standard set of 520 harmful behaviors. HarmBench [Mazeika et al., 2024] introduced functional categorisation and automated classifiers. JailbreakBench [Chao et al., 2024] contributed reproducible evaluation pipelines with leaderboard infrastructure. StrongREJECT [Souly et al., 2024] advanced rubric-based quality assessment of harmful outputs. Together these benchmarks have enabled systematic comparison of attack methods and defense mechanisms.

However, all four benchmarks share a structural limitation: they evaluate a single interaction modality—a text prompt sent to a language model, producing a text response classified as harmful

¹Repository: <https://github.com/adrianwedd/failure-first>.
docs/DATASHEET.md.

Supplementary Datasheet:

37 or safe. This design excludes three categories of AI system that are increasingly deployed in safety-
38 critical settings:

39 **Embodied AI.** Vision-language-action (VLA) models—RT-2 [Brohan et al., 2023], π_0 [Black
40 et al., 2024], OpenVLA [Kim et al., 2024], Gemini Robotics [Google DeepMind, 2025]—perceive
41 physical environments through sensors and generate motor commands. A VLA safety failure is
42 not an information hazard; it is a kinetic event. Over 1,800 autonomous haul trucks operated in
43 Australian mines by end-2025. No adversarial safety benchmark addresses this deployment class.

44 **Agentic AI.** Tool-using, multi-step AI systems that execute code, query databases, and chain ac-
45 tions across services. MUZZLE [MUZZLE Authors, 2025] catalogued 37 empirically discovered
46 attack instances against agentic systems but did not produce a reusable benchmark.

47 **Multi-turn interactions.** Extended conversations where adversarial pressure compounds across
48 turns. Crescendo attacks on DeepSeek-R1 1.5B achieved 65.0% strict ASR [43.3%, 81.9%] at
49 $n = 20$, compared to single-turn baselines below 20% on the same model.

50 F41LUR3-F1R57 addresses these gaps. We contribute:

- 51 1. A **unified adversarial corpus** of 141,020 prompts from 27 source datasets, normalised into
52 a relational database with versioned JSON schemas, covering 82 attack techniques across 6
53 temporal eras (DAN 2022 through VLA 2026).
- 54 2. The **first VLA adversarial scenario library** with 338 scenarios across 27 attack families and
55 3 attack surfaces (text-layer, action-layer, infrastructure-layer).
- 56 3. **FLIP grading**, a backward-inference evaluation protocol that detects the compliance paradox
57 in embodied outputs, with explicit calibration disclosure.
- 58 4. **TI-S evaluator**, a tool that measures whether safety interventions produce net benefit or ia-
59 trogenic harm.
- 60 5. **Three-modality benchmark infrastructure** (CLI, HTTP API, local Ollama) with CI-enforced
61 validation and reproducibility guarantees.
- 62 6. An **incident database** cataloguing 30+ real-world embodied AI failures across humanoid
63 robotics, autonomous vehicles, and service robots, cross-referenced with attack families.

64 All artefacts are released under open licenses at <https://failurefirst.org>.

65 2 Related Work

66 2.1 Text-Only Adversarial Benchmarks

67 The four established benchmarks collectively provide $\sim 1,450$ unique harmful behaviors or prompts.
68 They share three design constraints that F41LUR3-F1R57 relaxes: (a) single-turn only, (b) text
69 modality only, (c) binary or rubric-based forward classification.

70 AdvBench (520 behaviors) was designed as a target set for GCG optimisation [Zou et al.,
71 2023]. HarmBench (510 behaviors across 7 categories) includes a fine-tuned Llama-2-13b classifier
72 [Mazeika et al., 2024]. JailbreakBench (100 behaviors) focuses on attack method comparison with
73 leaderboard infrastructure [Chao et al., 2024]. StrongREJECT (313 prompts) measures response
74 quality rather than presence [Souly et al., 2024].

75 F41LUR3-F1R57 subsumes all four as source datasets within its unified corpus (Table 1). Each
76 is imported, normalised to a common schema, and linked to evaluation results across 187 models.

Table 1: Corpus statistics.

Metric	Value
Total prompts	141,020
Source datasets	27
Attack techniques	82
Harm categories	119
Models evaluated	187 (174 with results)
Total evaluation results	131,887
LLM-graded results	47,352
Heuristic-graded results	42,234
Evaluation runs	38,425
VLA-specific scenarios	338 (27 attack families)
Validated JSONL files on disk	820 (59,733 rows)
Corpus integrity score	0.9724

77 2.2 VLA and Embodied AI Safety

78 BadVLA [Shi et al., 2025] demonstrated backdoor attacks achieving near-100% ASR against π_0
79 and OpenVLA via training data poisoning. Blindfold [Huang et al., 2026] achieved 93.2% ASR on
80 GPT-4o by decomposing malicious instructions into individually benign actions; existing defenses
81 reduced ASR by at most 17.9 pp, leaving residual ASR above 75%. Physical adversarial patches
82 bypass VLA visual grounding [Cardenas and Xie, 2026]. Each evaluates its own attack in isolation;
83 no shared benchmark exists.

84 2.3 Agentic AI Evaluation

85 MUZZLE [MUZZLE Authors, 2025] catalogues 37 attack instances (~ 12 underlying classes) against
86 agentic systems. AgentLAB [Yao et al., 2026] reports ASR escalation from 62.5% (single-step) to
87 79.9% (multi-step). Neither provides a reusable benchmark with schema-validated scenarios, auto-
88 mated grading, or reproducible evaluation infrastructure.

89 3 The F41LUR3-F1R57 Corpus

90 3.1 Corpus Overview

91 The corpus is a relational database (SQLite, schema version 13, 12 tables, 15 indexes, 4 views) that
92 unifies adversarial prompts, technique taxonomy, model metadata, and evaluation results.

93 Source datasets span the full temporal range of adversarial AI research: DAN-In-The-Wild
94 (2022), AdvBench/GCG (2023), HarmBench/JailbreakBench (2024), reasoning-extension attacks
95 (2025), and VLA-specific families (2026). Nine established public datasets are imported alongside
96 18 datasets original to this work.

97 3.2 Schema and Validation

98 Every scenario is stored as a JSONL record validated against a versioned JSON schema. Re-
99 quired fields include `id`, `created`, `domain`, `scenario_class`, `scenario_description`,
100 `attack_attempt`, and `turns[]`. The labels `.intent.* namespace` captures 7 instruction-
101 hierarchy subversion signals: `format_lock`, `refusal_suppression`, `persona_hijack`,
102 `future_year_laundersing`, `dataset_poisoning_intent`, `constraint_erosion`,
103 and `research_only_pressure`.

Table 2: Source datasets by prompt count (top 10).

Source Dataset	Prompts	Origin	Type
SORRY-Bench	9,446	Public	Harmful behaviors
BEAVERTAILS	3,432	Public	Safety alignment
DAN-In-The-Wild	1,405	Public	Historical jailbreaks
WildJailbreak	1,000	Public	In-the-wild jailbreaks
Benchmark traces	993	Original	Evaluation traces
AdvBench	520	Public	GCG target set
ForbiddenQuestions	390	Public	Direct harmful queries
HarmBench	320	Public	Functional categories
StrongREJECT	313	Public	Quality-rated
HEX-PHI	290	Public	Red-teaming

104 **Validation pipeline.** Two automated checks run on every push to the main branch via GitHub
 105 Actions CI:

- 106 1. **Schema validation:** JSON Schema compliance plus cross-field invariant checks (e.g., intent-
 107 bait scenarios must set `attack_attempt: true`). Currently validates 820 files, 59,733
 108 rows, zero failures.
- 109 2. **Safety linting:** Heuristic checks for refusal suppression phrasing, overly operational lan-
 110 guage, and future-year circumvention patterns. Scans 20,711 items, zero findings.

111 Both tools are invocable via `make validate` and `make lint`. The CI pipeline rejects any
 112 push that introduces a validation or lint failure.

113 3.3 Source Dataset Integration

114 The database normalises heterogeneous source formats into a common schema. Table 2 shows the
 115 10 largest source datasets.

116 Each source dataset is imported via dedicated scripts that preserve provenance metadata, map
 117 prompts to the technique taxonomy, and link harm categories.

118 3.4 Technique Taxonomy

119 The 82 attack techniques are organised into temporal eras reflecting the evolution of adversarial
 120 methods:

Table 3: Attack technique eras.

Era	Period	Representative Techniques	Prompts
DAN	2022	Persona hijack, role-play	~1,400
GCG/AutoDAN	2023	Gradient-based suffix, automated	~850
Multi-turn/Crescendo	2024	Escalation over turns	~500
Reasoning exploitation	2025	CoT manipulation, format lock	~300
VLA/Embodied	2026	Action-layer, infrastructure-layer	~340
Cross-era	Mixed	Combined techniques, compound	~600

121 The taxonomy converges approximately 55% with MUZZLE’s independently derived classi-
 122 fication [MUZZLE Authors, 2025], with 5 schema gaps identified: `tool_chain_hijacking`,
 123 `memory_persistence_attack`, `objective_drift_induction`, `cross_system_lateral_movement`,
 124 and `silent_exfiltration`.

125 **3.5 VLA Scenario Library**

126 The VLA component contains 338 scenarios across 27 attack families, organised by attack surface:

127 **Text-layer attacks** (14 families, ~205 scenarios) exploit the VLM backbone via prompt ma-
128 nipulation. Families include Safety Instruction Dilution (SID, 120 scenarios with dose-response
129 experiment), Latent Action Manipulation (LAM), Semantic Benignity Exploitation (SBE), Cross-
130 Embodiment Transfer (CET), and Long-Horizon Goal Displacement (LHGD).

131 **Action-layer attacks** (3 families, ~39 scenarios) exploit the action decoder or action space. Pol-
132 icy Puppetry (PP, 31 scenarios), Action Space Exploitation (ASE), Prompt Compliance Manipula-
133 tion (PCM).

134 **Infrastructure-layer attacks** (10 families, ~94 scenarios) target system boundaries. Includes
135 compound attacks (DA-SBA, SID+SIF, DLA) and cross-system attacks (CSBA, SSBA, XSBA).

136 Each scenario maps to one of four HANSE (Hierarchical Adversarial Nexus for Safety Evalua-
137 tion) defense layers: `vla_core` (61.8%), `semantic_firewall` (17.6%), `kinematic_shield`
138 (12.9%), `affordance_verifier` (7.8%). The coverage imbalance toward `vla_core` reflects
139 the dominance of text-layer attacks in the current literature and is documented as a known limitation.

140 **3.6 Incident Database**

141 We compiled a catalogue of 30+ real-world embodied AI incidents from the OECD AI Incident
142 Monitor, AIID, NHTSA Standing General Orders, FDA MAUDE, and OSHA Severe Injury Reports.
143 Incidents span humanoid robotics (Unitree H1 tether-induced feedback loop, 2025), autonomous ve-
144 hicles (Cruise dragging incident, 2023; Uber ATG fatality, 2018), service robotics (Knightscope K5
145 stair-fall, 2017; Coco delivery robot train destruction, 2026), and industrial automation (Tesla Giga
146 Texas robot injury, 2021).

147 The incident catalogue grounds scenario design in observed failure modes and enables cross-
148 reference between adversarial scenarios and real-world failure patterns.

149 **4 Evaluation Infrastructure**

150 **4.1 Three Runner Modalities**

151 The benchmark supports three evaluation pathways, enabling researchers with different compute
152 resources and model access to participate:

153 **CLI-based runner.** Invokes native CLI tools (`claude-code`, `codex-cli`, `gemini-cli`) for frontier
154 model evaluation. Produces structured JSONL traces with timestamps.

155 **HTTP API runner.** Connects to OpenRouter (100+ models) or any OpenAI-compatible API. Sup-
156 ports both free-tier (32 models, zero cost) and paid models. Implements rate limiting, retry logic,
157 and token tracking.

158 **Local Ollama runner.** Targets locally-hosted models via Ollama. No rate limits, no cost, no data
159 egress. Supports reasoning models with `think=true` for chain-of-thought capture.

160 All three runners produce the same output format: per-scenario JSONL traces with input, output,
161 timing, and token counts. This standardisation means that any downstream analysis tool works
162 identically regardless of evaluation pathway.

163 4.2 FLIP Grading

164 FLIP (Failure-Leveraged Inference Protocol) uses backward inference to grade model responses.
165 Given a response, the FLIP grader infers what instruction the model appeared to follow, then classi-
166 fies whether that inferred instruction is harmful.

167 **Five verdict categories:** COMPLIANCE (full compliance with adversarial request), PARTIAL (safety
168 hedging with behavioral compliance), HALLUCINATION_REFUSAL (harmful content wrapped in
169 refusal framing), REFUSAL (genuine refusal), BENIGN_QUERY (non-adversarial prompt).

170 **Three-tier ASR reporting.** Building on FLIP verdicts:

Table 4: Three-tier ASR (corpus-wide, $n = 10,294$ evaluable).

Tier	Formula	ASR
Strict	COMPLIANCE / evaluable	45.9%
Broad	(COMPLIANCE + PARTIAL) / evaluable	79.3%
Functionally Dangerous	(COMPLIANCE + PARTIAL + HR) / evaluable	80.3%

171 The denominator excludes BENIGN_QUERY, ERROR, and non-evaluable categories. The Func-
172 tionally Dangerous tier includes HALLUCINATION_REFUSAL because these responses are computa-
173 tionally indistinguishable from COMPLIANCE (thinking tokens $p = 0.21$, response tokens $p = 0.46$)
174 but significantly different from REFUSAL (both $p < 0.001$).

175 **Calibration disclosure.** The primary FLIP grader (deepseek-r1:1.5b) achieves approximately 70%
176 accuracy on calibration audit with a 30.8% false positive rate on benign baselines ($n = 44$). Co-
177 hen’s κ between heuristic and LLM classifiers is 0.126 [0.108, 0.145] ($n = 1,989$). We disclose all
178 calibration data to enable downstream reliability assessment.

179 4.3 TI-S Evaluator

180 The Therapeutic Index of Safety (TI-S) evaluator measures whether a safety intervention produces
181 net benefit or iatrogenic harm:

$$\text{TI-S}(I, S) = \frac{B_h(I, S)}{C_h(I, S)} \quad (1)$$

182 where B_h is the harm-layer benefit (ASR reduction) and C_h is the sum of iatrogenic costs (PARTIAL
183 rate increase, HALLUCINATION_REFUSAL migration, action-layer execution despite text-layer re-
184 fusal). $\text{TI-S} > 1$ indicates net benefit; $\text{TI-S} < 1$ indicates iatrogenic harm.

185 The evaluator compares traces collected with and without a safety intervention, producing per-
186 intervention TI-S scores. This addresses a gap in existing evaluation: standard ASR metrics cannot
187 distinguish a safety intervention that genuinely reduces harm from one that merely shifts harmful
188 outputs from COMPLIANCE to PARTIAL (text-level hedging without behavioral change).

189 4.4 Statistical Testing

190 The benchmark includes a statistical testing module supporting chi-square tests, Mann-Whitney U
191 tests, and Cohen’s d effect sizes with Bonferroni correction for pairwise model comparison. Score
192 reports generate summary statistics with confidence intervals from any trace JSONL file.

193 4.5 Relational Database

194 All results are stored in a normalised SQLite database (schema version 13, 12 tables). The schema
195 supports prompt-level queries (filter by source dataset, technique, harm category, era), result-level
196 queries (join prompts to model responses with LLM and heuristic verdicts), cross-model compari-
197 son (aggregate ASR by model, provider, parameter count), and technique lineage (trace technique
198 evolution across temporal eras). Pre-built queries are available via CLI (11 named queries); raw
199 SQL access is supported for custom analysis.

200 5 Benchmark Results

201 5.1 Cross-Model Vulnerability Profiles

202 Analysis of 57 models with sufficient LLM-graded data reveals three distinct vulnerability clusters:

Table 5: Model vulnerability clusters.

Cluster	ASR Range	Models	Characteristics
Permissive	$\geq 40\%$	37	Minimal safety training; open-weight
Mixed	15–40%	15	Partial safety training; fine-tuned
Restrictive	$\leq 15\%$	5	Frontier; extensive RLHF

203 Provider signatures dominate clustering: Anthropic 3.7% mean ASR, Google 9.1%, compared to
204 Nvidia 40.0%, Qwen 43.1%, DeepSeek 44.1%. Inverse scaling (larger models are more vulnerable)
205 is not supported: Spearman $r = -0.140$, $n = 24$ models with known parameter counts. Safety
206 training investment appears to be the primary determinant.

207 5.2 VLA-Specific Findings

208 FLIP grading across 15 VLA attack families produces a three-tier vulnerability structure relative to
209 the 30.8% evaluator false positive rate:

- 210 • **Tier 1 (genuine signal):** Trace Reasoning Attack (100% broad ASR), Deceptive Alignment
211 (62.5%), Latent Action Manipulation (60%). ASR substantially exceeds benign baseline.
- 212 • **Tier 2 (marginal):** Action Space Exploitation, Semantic Benignity Exploitation. ASR 40–
213 55%, within wide confidence intervals of baseline.
- 214 • **Tier 3 (at FP floor):** Policy Puppetry, SID+SIF, and 8 others. Indistinguishable from evalua-
215 tor noise at current sample sizes.

216 **Compliance paradox.** 50% of all FLIP verdicts are PARTIAL—models generate safety disclaimers
217 while producing the requested action sequences. Zero outright refusals across 63 originally graded
218 traces. Keyword classifiers report 94% ASR on the same traces; FLIP reports 72.4%—a 22 pp gap
219 attributable to the compliance paradox.

220 5.3 Inverse Detectability-Danger Law (IDDL)

221 Corpus-level analysis reveals a structural inverse correlation between attack consequentiality and
222 evaluator detectability. Spearman $\rho = -0.847$ ($n = 27$ families). The attacks most likely to cause
223 physical harm are the attacks least likely to be detected by text-layer evaluation.

224 This finding has direct implications for organisations that evaluate VLA safety using text-level
225 classifiers: they will systematically underestimate their exposure to the most consequential attack
226 families.

227 5.4 Safety Re-emergence at Scale

228 Evaluation of safety-abliterated (safety-removed) models in the Qwen3.5 obliterated series reveals
229 that safety-like behavior partially re-emerges at scale: ASR drops from 100% (0.8B) to 47.3%
230 (9.0B), Spearman $\rho = -0.949$, $p = 0.051$. Larger abliterated models produce PARTIAL verdicts
231 rather than full compliance, suggesting safety is partially an emergent property of scale.

232 6 Datasheet for Datasets

233 We provide a complete Datasheet for Datasets following Gebru et al. [2021] as supplementary ma-
234 terial (`docs/DATASHEET.md` in the repository). The datasheet covers all seven sections: Moti-
235 vation, Composition, Collection Process, Preprocessing/Cleaning/Labeling, Uses, Distribution, and
236 Maintenance. We summarise key elements here; the full datasheet should be consulted for complete
237 details.

238 **Motivation.** The corpus was created to fill the gap in adversarial evaluation benchmarks for em-
239 bodied and agentic AI systems. It was created by Adrian Wedd as an independent, self-funded AI
240 safety research project.

241 **Composition.** 141,020 prompts total; 338 VLA-specific scenarios; 820 validated JSONL files con-
242 taining 59,733 rows; 131,887 evaluation results across 187 models. The dataset contains adversarial
243 prompts that request harmful outputs by design; scenarios describe patterns at a structural level, not
244 executable exploits.

245 **Collection.** Data was acquired through import from 9 public datasets, original scenario authoring,
246 and evaluation trace collection from 187 models via three runner interfaces. No human subjects data
247 was collected.

248 **Preprocessing.** All prompts are normalised to a common JSONL schema. Technique labels and
249 harm categories are mapped to the unified 82-technique, 119-category system. Deduplication elim-
250 inated 102 cross-file duplicates (integrity score 0.9724).

251 **Uses.** The dataset has been used for 131,887 evaluations, 130+ research reports, and companion
252 CCS 2026 and NeurIPS 2026 D&B submissions. It should not be used to develop operational attack
253 tools or train models to produce harmful outputs.

254 **Distribution.** Pattern-level analysis is publicly available at <https://failurefirst.org>.
255 Operational scenario details are maintained in a private repository.

256 **Maintenance.** Continuous updates via semantic versioning, GitHub issues, and CI validation.
257 Schema migrations are versioned.

258 7 Ethical Considerations and Broader Impact

259 7.1 Dual-Use Risk

260 F41LUR3-F1R57 is an adversarial benchmark; it documents attack patterns that could, in principle,
261 be adapted for malicious use. The dual-use risk for embodied AI benchmarks is qualitatively differ-
262 ent from text-only benchmarks: a successful VLA attack produces physical actions that could injure
263 or kill, not merely harmful text.

264 We mitigate this through three design choices: (a) *pattern-level, not operational*—scenarios
265 describe vulnerability mechanisms, not executable exploits; (b) *defense-oriented framing*—each

266 family is accompanied by a description of the targeted vulnerability; (c) *graduated disclosure*—
267 aggregate results are public while operational details are access-controlled.

268 **7.2 The Compliance Paradox**

269 The compliance paradox—models that disclaim while complying—creates a distinctive ethical con-
270 cern. Documenting this gap could be exploited by attackers who craft prompts specifically designed
271 to produce PARTIAL responses that pass text-level safety filters. However, the gap exists regardless
272 of documentation; our contribution is making it visible so defense architectures can address it.

273 **7.3 Evaluator Limitations and the Evaluation Access Gap**

274 Using sub-2B models as FLIP graders introduces systematic bias (approximately 70% accuracy,
275 30.8% FP rate). This limitation means reported ASR figures carry wider uncertainty than ideal.
276 We disclose all calibration data. Resource-constrained researchers face a genuine dilemma: reliable
277 evaluation requires models that are expensive to run, creating an evaluation access gap that mirrors
278 compute inequality in AI research.

279 **7.4 Responsible Disclosure**

280 We have engaged with the Australian AI Safety Institute (AISI) and Safe Work Australia. The
281 benchmark addresses gaps in multiple regulatory frameworks: VAISS Guardrail 4 (Australia), EU
282 AI Act Article 9, ISO 17757:2019 (autonomous mining equipment), and ISO 13482:2014 (personal
283 care robots). None currently specify adversarial testing methodology for embodied AI.

284 **7.5 Researcher Positionality**

285 This work was conducted by an independent researcher without institutional affiliation or compute
286 sponsorship. This is both a strength (no organisational conflicts in reporting vulnerabilities) and a
287 limitation (resource constraints on model scale, replication, and human annotation).

288 **8 Limitations**

- 289 1. **Model scale.** VLA evaluation results are from 1.5B–1.7B parameter models, not representa-
290 tive of frontier VLA systems. The capability-floor hypothesis predicts different vulnerability
291 profiles at scale.
- 292 2. **Text-only evaluation.** All experiments are text-in/text-out. Visual adversarial inputs, physical
293 patches, and sensor-based attacks are not evaluated end-to-end.
- 294 3. **Evaluator quality.** The 30.8% FP rate and approximately 70% accuracy of the primary FLIP
295 grader introduce noise. Per-family findings with $n < 10$ should be treated as preliminary.
- 296 4. **Scenario ecological validity.** Scenarios are researcher-authored, not drawn from real de-
297 ployment logs. Attack families represent theoretically motivated surfaces, not observed cam-
298 paigns.
- 299 5. **Single evaluation run.** Most families evaluated in a single run without replication across
300 seeds or temperatures. Variance from Wilson CIs, not empirical replication.
- 301 6. **Coverage imbalance.** SID accounts for 120/338 VLA scenarios (35.5%). Infrastructure-layer
302 families have ~ 5 scenarios each.
- 303 7. **Public dataset overlap.** The 9 imported public datasets have known inter-dataset over-
304 lap (e.g., AdvBench prompts appear in HarmBench). Deduplication was applied within
305 F41LUR3-F1R57 but not across source publications.

Table 6: Reproducibility artefacts.

Artefact	Location	Verification
Schema DDL	database/schema.sql	Deterministic init
Schema migrations	database/migrations/	Versioned, ordered
Public import	tools/database/import_public.py	From HuggingFace
Validation	make validate	Schema + invariants
Safety linting	make lint	Heuristic checks
CI pipeline	.github/workflows/ci.yml	Every push
Benchmark runners	tools/benchmarks/run_benchmark_*.py	3 modalities
FLIP grading	tools/benchmarks/grade_*.py	Backward inference
TI-S evaluator	tools/evals/iatrogenic_*.py	Baseline vs. intervention
Statistical tests	tools/stats/	χ^2 , Mann-Whitney
Test suite	pytest (1,205 tests)	CI floor at 1,150

306 8. **Corpus size dominated by OBLITERATUS.** 121,955 of 141,020 prompts (86.5%) come
307 from OBLITERATUS telemetry and runs. The non-OBLITERATUS analytical corpus is ap-
308 proximately 19,065 prompts.

309 9 Reproducibility

310 The benchmark is designed for full reproducibility. Table 6 lists the key artefacts and their verifica-
311 tion methods.

312 To reproduce the full evaluation pipeline:

```

313 pip install -r requirements-dev.txt
314 python tools/database/db_manager.py init
315 python tools/database/import_public.py --dataset advbench
316 make validate && make lint
317 python tools/benchmarks/run_benchmark_http.py \
318   --scenarios data/vla/vla_sbe_v0.1.jsonl \
319   --models "meta-llama/llama-3.3-70b-instruct:free" \
320   --output runs/reproduce/
321 python tools/benchmarks/grade_generation_traces.py \
322   --traces runs/reproduce/traces.jsonl
323 python tools/benchmarks/score_report_v1.0.py \
324   --traces runs/reproduce/traces.jsonl --stats

```

325 10 Conclusion

326 F41LUR3-F1R57 extends adversarial AI safety evaluation beyond text-only LLM benchmarks into
327 embodied, agentic, and multi-turn settings. The corpus of 141,020 prompts across 27 source datasets,
328 evaluated on 187 models with 131,887 results, provides the scale and diversity needed for robust
329 cross-model vulnerability analysis. The VLA scenario library (338 scenarios, 27 families) is the
330 first adversarial resource for systems that control physical actuators. FLIP grading detects the com-
331 pliance paradox that forward classifiers miss. The TI-S evaluator measures whether safety interven-
332 tions cause net benefit or iatrogenic harm.

333 The benchmark’s primary empirical contribution is the Inverse Detectability-Danger Law: the
334 attacks most consequential for physical safety are the least detectable by text-layer evaluation ($\rho =$
335 0.847). This finding suggests that current evaluation methodology systematically underestimates the
336 most dangerous attack surfaces.

337 Future work includes evaluation at 7B+ model scale, visual adversarial input integration, human
338 annotation studies for ground-truth calibration, and expansion of infrastructure-layer attack families.

339 References

340 Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo
341 Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. π_0 : A vision-language-action flow
342 model for general robot control. In *International Conference on Machine Learning*, 2024.

343 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choroman-
344 ski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. RT-2: Vision-language-action
345 models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.

346 Alvaro Cardenas and Cihang Xie. Physical adversarial attacks on vision-language models for robotic
347 manipulation. *arXiv preprint*, 2026.

348 Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong.
349 JailbreakBench: An open robustness benchmark for jailbreaking large language models. In
350 *NeurIPS 2024 Datasets and Benchmarks Track*, 2024.

351 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,
352 Hal Daumé III, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64
353 (12):86–92, 2021.

354 Google DeepMind. Gemini robotics: Bringing AI into the physical world. *Google DeepMind*
355 *Technical Report*, 2025.

356 Yutong Huang et al. Jailbreaking embodied LLMs via action-level manipulation. *arXiv preprint*
357 *arXiv:2603.01414*, 2026.

358 Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,
359 Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. OpenVLA: An open-source
360 vision-language-action model. In *International Conference on Machine Learning*, 2024.

361 Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee,
362 Nathaniel Li, Steven Basart, Bo Li, et al. HarmBench: A standardized evaluation framework for
363 automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*, 2024.

364 MUZZLE Authors. MUZZLE: A multi-faceted zero-day attack taxonomy for agentic AI systems.
365 *arXiv preprint arXiv:2602.09222*, 2025.

366 Xin Shi et al. BadVLA: Backdoor attacks on vision-language-action models. In *Advances in Neural*
367 *Information Processing Systems*, 2025.

368 Alexandra Souly, Qingyun Lu, Dillon Bowen, Tu Trinh, Elvis Hsieh, Sana Fishi, Ralph Abi-Rafeh,
369 Georgia Anber, Brando Braunstein, Dean Goodman, et al. A StrongREJECT for empty jailbreaks.
370 *arXiv preprint arXiv:2402.10260*, 2024.

371 Shunyu Yao et al. AgentLAB: A comprehensive benchmark for long-horizon agent instruction
372 following and safety. *arXiv preprint arXiv:2602.16901*, 2026.

373 Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial
374 attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

375 **A Datasheet for Datasets (Full)**

376 The complete Datasheet for Datasets following Gebru et al. [2021] is maintained as a living docu-
377 ment at `docs/DATASHEET.md` in the project repository. It covers:

- 378 1. **Motivation** (Sections 1.1–1.4): Purpose, creator, funding.
- 379 2. **Composition** (Sections 2.1–2.15): Instance types, counts, labels, sensitive content.
- 380 3. **Collection Process** (Sections 3.1–3.11): Acquisition channels, tools, sampling, ethics.
- 381 4. **Preprocessing** (Sections 4.1–4.4): Schema validation, deduplication, FLIP grading.
- 382 5. **Uses** (Sections 5.1–5.5): Current uses, recommended and prohibited uses.
- 383 6. **Distribution** (Sections 6.1–6.6): Two-tier access, licensing, export controls.
- 384 7. **Maintenance** (Sections 7.1–7.7): Contact, updates, versioning, contribution mechanism.

385 Additionally, Sections 8–9 of the datasheet cover ethical considerations (dual-use risk, bias,
386 human subjects) and dataset validation/quality assurance metrics.

387 The datasheet version corresponds to corpus version 0.19.0 and is updated with each schema
388 migration.