

Prompt Contracts: A Comprehensive Probabilistic Formalization for Testing and Validating Large Language Model Outputs

Philipp Melikidis
Independent Researcher
Bietigheim-Bissingen, Germany
philipp.melikidis@gmail.com

Abstract

Large Language Models (LLMs) function as stochastic, untyped interfaces lacking formal specifications. We introduce PCSL v0.3, a comprehensive probabilistic formalization for LLM prompt testing with rigorous statistical validation. Through N-sampling with configurable aggregation (majority/all/any) and bootstrap confidence intervals ($B=1000$, $\delta=0.95$, percentile method), PCSL quantifies reliability across 5 tasks totaling 250 labeled fixtures. Evaluation demonstrates 100% validation success with enforce mode, 95% with assist (18% repair rate, task accuracy invariant). Seed robustness: mean 94.2%, std 1.1% across 5 seeds. We address reproducibility through: (1) comprehensive dataset documentation with inter-rater reliability ($\kappa=0.88$, substantial agreement); (2) fixed seeds (42) and pinned dependencies; (3) detailed bootstrap parameters; (4) compliance mapping to ISO/IEC/IEEE 29119 and EU AI Act with audit examples. All code, 250 fixtures (CC BY 4.0), and documentation are publicly available. Overhead: <5ms validation, 2.1× latency for $N=5$. One-command Docker reproduction (make eval-full).

Keywords

Large Language Models, Prompt Engineering, Probabilistic Contracts, Statistical Validation, Compliance

1 Introduction

Large Language Models function as *untyped, stochastic interfaces*: prompts map inputs to probabilistic outputs without formal behavioral guarantees [1]. Consider an LLM as $f_\theta : \mathcal{X} \rightarrow P(\mathcal{Y})$ where $P(\mathcal{Y})$ denotes a probability distribution over outputs. Unlike deterministic APIs with explicit contracts, LLM outputs vary across runs, making traditional contract testing insufficient.

This gap becomes critical as LLMs deploy in regulated domains [4]. The EU AI Act mandates transparency, auditability, and robustness testing. Yet prompt engineering lacks specification infrastructure: no type checking, no contract enforcement, no systematic validation with statistical confidence.

Research Problem. How can we define, validate, and enforce behavioral contracts for probabilistic LLM interfaces while ensuring reproducibility and regulatory compliance?

Contributions.

- (1) **Probabilistic specification:** PCSL v0.3 with N-sampling, aggregation policies, bootstrap CIs ($B=1000$), convergence proofs, and compositional semantics (Section 3).
- (2) **Rigorous evaluation:** Five tasks (1,247 fixtures), ablation studies (N , aggregation, repair, τ), seed robustness (5 seeds),

Table 1: Framework Comparison

Framework	Contracts	Probabilistic	CI/CD	Semantic	Compliance
HELM [9]	×	×	×	×	×
CheckList [14]	Manual	×	Partial	×	×
Guidance [11]	×	×	×	×	×
OpenAI Struct. [12]	Partial	×	×	×	×
PCSL v0.3	✓	✓	✓	✓	✓

comparative benchmarks (CheckList, Guidance, OpenAI), LLM-judge vs. human ($\kappa = 0.82$) (Section 5).

- (3) **Compliance framework:** ISO 29119 mapping with audit case study including real artifacts (Section 4.5), operationalizing compliance-as-code.

2 Related Work

Contract-based testing. Design-by-contract [10] formalizes deterministic specifications. PCSL extends to probabilistic functions via N-sampling and statistical confidence bounds. OpenAPI [13] provides REST API contracts; PCSL adapts this for natural language interfaces.

LLM frameworks. CheckList [14] enables behavioral testing but requires manual test writing (120 min setup vs. PCSL’s 2 min). HELM [9] focuses on model benchmarking, not prompt contracts. LangChain [2] abstracts development but lacks systematic testing. Guidance [11] constrains generation; PCSL validates post-hoc. OpenAI Structured Outputs [12] enforces schemas but is vendor-locked. PCSL uniquely combines formal specification, probabilistic semantics, multi-provider execution, and compliance mapping (Table 1).

Regulation. EU AI Act [4] mandates transparency (Art. 13), records (Art. 12), accuracy (Art. 15). ISO 29119 [7] codifies testing principles. PCSL bridges requirements through formal artifact mapping (Section 4.5).

3 PCSL: Formal Specification

3.1 Core Definitions

A *prompt contract* $C = \langle \mathcal{P}, \mathcal{E}, \mathcal{X} \rangle$ consists of: Prompt Definition \mathcal{P} (template, I/O expectations), Expectation Suite $\mathcal{E} = \{e_1, \dots, e_m\}$ (validation checks), Evaluation Profile \mathcal{X} (fixtures, targets, config). Each check $e_i : \Omega \rightarrow \{\text{pass}, \text{fail}\}$. Single-output satisfaction:

$$\text{sat}(C, o) \iff \bigwedge_{i=1}^m e_i(o) = \text{pass}$$

3.2 Probabilistic Semantics

Given stochastic LLM f_θ , *probabilistic satisfaction*:

$$\Pr[\text{sat}(C, o)] = \Pr_{o \sim f_\theta(x)}[\text{sat}(C, o)]$$

PCSL estimates via N-sampling: $\{o_1, \dots, o_N\}$, empirical pass rate $\hat{p} = \frac{1}{N} \sum_{j=1}^N \mathbb{I}[\text{sat}(C, o_j)]$.

Statistical properties. The estimator \hat{p} is unbiased: $\mathbb{E}[\hat{p}] = p$. Variance:

$$\text{Var}(\hat{p}) = \frac{p(1-p)}{N}$$

decreases as $O(1/N)$, enabling precision-confidence tradeoffs. Standard error: $\text{SE}(\hat{p}) = \sqrt{p(1-p)/N}$.

Convergence. By Central Limit Theorem:

$$\sqrt{N}(\hat{p} - p) \xrightarrow{d} \mathcal{N}(0, p(1-p))$$

Approximate 95% CI: $\hat{p} \pm 1.96\sqrt{\hat{p}(1-\hat{p})/N}$.

Bootstrap CIs. Percentile method [?] provides non-parametric bounds. Algorithm: (1) Resample with replacement $B = 1000$ times, (2) compute $\hat{p}^{(b)}$ for each, (3) report 2.5th and 97.5th percentiles. Convergence: CI width stabilizes at $B \geq 500$ (empirical variance < 0.001 for $B \in [500, 2000]$).

Aggregation policies $A : \{o_1, \dots, o_N\} \rightarrow \{\text{PASS}, \text{FAIL}\}$:

$$\begin{aligned} A_{\text{first}}(\{o_j\}) &= \text{sat}(C, o_1) \\ A_{\text{majority}}(\{o_j\}) &= \text{PASS} \iff \hat{p} > 0.5 \\ A_{\text{all}}(\{o_j\}) &= \text{PASS} \iff \hat{p} = 1.0 \\ A_{\text{any}}(\{o_j\}) &= \text{PASS} \iff \hat{p} > 0 \end{aligned}$$

Fixture-level validation with tolerance τ :

$$C \models_{\tau} \mathcal{F} \iff \frac{|\{f \in \mathcal{F} \mid A(\{o_f^j\}) = \text{PASS}\}|}{|\mathcal{F}|} \geq \tau$$

3.3 Compositional Semantics

For multi-step pipelines (e.g., RAG = retrieval \circ generation): $C_{\text{comp}} = C_1 \circ C_2$. Satisfaction:

$$\text{sat}(C_1 \circ C_2, (i, o)) \iff \text{sat}(C_1, (i, o_{\text{inter}})) \wedge \text{sat}(C_2, (o_{\text{inter}}, o))$$

where o_{inter} is intermediate output.

Complexity. Pipeline: $O(|\mathcal{F}| \cdot N \cdot (|\mathcal{E}_1| + |\mathcal{E}_2|) \cdot \max(n_1, n_2))$ where n_i = output size. Parallel sampling (N workers): $O(|\mathcal{F}| \cdot (|\mathcal{E}_1| + |\mathcal{E}_2|) \cdot \max(n_1, n_2))$.

3.4 Check Catalog

Structural ($O(n)$): json_valid, json_required, enum, regex_absent, token_budget, latency_budget. **Semantic**: contains_all, contains_any, regex_present, similarity (sentence-transformers MiniLM-L6-v2 [?], cosine threshold ≥ 0.8). **Judge** [15]: LLM-as-judge with natural language criteria.

4 Framework Architecture

4.1 Execution Pipeline

Algorithm 1 formalizes sampling-enabled execution.

4.2 Execution Modes

observe: Validation only. **assist**: Prompt augmentation with constraints. **enforce**: Schema-guided JSON (OpenAI response_format). **auto**: Capability-based fallback (enforce \rightarrow assist \rightarrow observe). Negotiation: $\mu(\mathcal{A}_{\text{cap}}, M_{\text{req}}) \rightarrow M_{\text{actual}}$.

Algorithm 1 PCSL Execution with Probabilistic Sampling

```

1: Input:  $C = \langle \mathcal{P}, \mathcal{E}, \mathcal{X} \rangle$ ,  $(N, \text{seed}, A)$ ; Output:  $\mathcal{R}$ 
2:  $\mathcal{R} \leftarrow \emptyset$ ; if seed then set_seed(seed)
3: for each  $f \in \mathcal{X}.\text{fixtures}$  do
4:    $p \leftarrow \text{render}(\mathcal{P}, f)$ ;  $\mu \leftarrow \text{negotiate}(\text{adapter.cap}(), \mathcal{X}.\text{mode})$ 
5:   if  $\mu = \text{enforce}$  then
6:      $\sigma \leftarrow \text{derive\_schema}(\mathcal{E})$ 
7:   end if
8:   if  $\mu = \text{assist}$  then
9:      $p \leftarrow \text{augment}(p, \mathcal{E})$ 
10:  end if
11:  for  $j = 1$  to  $N$  do
12:     $o_r^j \leftarrow \text{adapter.gen}(p, \sigma)$ ;  $o_n^j \leftarrow \text{repair}(o_r^j, \Pi)$ 
13:     $\text{res}^j \leftarrow \{e_i(o_n^j) \mid e_i \in \mathcal{E}\}$ ; Append to samples
14:  end for
15:   $s, \text{CI} \leftarrow A(\text{samples})$ , bootstrap_ci(samples,  $B = 1000$ )
16:   $\mathcal{R} \leftarrow \mathcal{R} \cup \{(f, s, \text{CI}, \text{samples})\}$ 
17: end for
18: return  $\mathcal{R}$ 

```

Table 2: Compliance Mapping

PCSL	ISO 29119	EU AI Act
PD	Test Item (§7.1)	-
ES	Test Conditions (§7.2)	Art. 15 (accuracy)
EP	Test Case (§7.3)	Art. 9 (risk mgmt)
save_io	Test Log (§8.3)	Art. 12 (records)
Negotiation	Test Env (§8.1)	Art. 13 (transparency)
N-sampling+CI	Statistical (29119-4)	Art. 15 (robustness)
Repair ledger	Incident (§8.4)	Art. 14 (oversight)

4.3 Repair Policy

$\Pi = \langle \text{enabled}, \text{max_steps}, \text{allowed} \rangle$. Strategies: strip_markdown_fences ($O(n)$), json_loose_parse (4 strategies), lowercase_fields ($O(d)$). Risk: High repair rate (> 0.5) signals quality issues. Fail-safe: max_steps=0. Fail-open: max_steps=2. All logged in repair ledger.

4.4 Compliance Mapping

Table 2 maps PCSL to ISO 29119 and EU AI Act.

4.5 Audit Case Study

Scenario: Healthcare support classifier (EU AI Act Art. 6(2): high-risk). Workflow: (1) Define contract, (2) Run `-save_io audit/`, (3) Generate `-report junit`.

Artifacts: input_final.txt (prompt with constraints), output_raw.txt (model response), output_norm.txt (post-repair), run.json (metadata with timestamp, seed, checks, repair ledger, prompt hash SHA-256).

Verification: ISO 29119 §8.3: test log \checkmark . EU Art. 12: immutable hash, repair ledger \checkmark . EU Art. 13: capability negotiation log \checkmark .

5 Evaluation

5.1 Setup

Tasks: Classification, Extraction, RAG Q&A, Summarization, Tool-calls. Total: 250 labeled fixtures (50 per task). **Models:** GPT-4o-mini (enforce), Mistral-7B (assist). **Metrics.** (1) *validation_success*: Percentage passing all checks, (2) *task_accuracy*: Exact match to gold labels (when available), (3) *repair_rate*: Fraction requiring normalization, (4) *latency_ms*: Mean generation time, (5) *overhead_pct*: Check execution time as percentage of total latency. **Reproducibility:** Fixed seed=42, temp=0, pinned dependencies (Python 3.11.7, torch

2.0.1, sentence-transformers 2.2.2). Docker: prompt-contracts:0.3.1 with PYTHONHASHSEED=42. Reproduction: make eval-full.

5.2 Statistical Methodology

Bootstrap Confidence Intervals. We employ the percentile bootstrap method [3, 6] to construct $(1 - \alpha)$ confidence intervals for validation success rates. Given N samples, we: (1) Generate $B = 1000$ bootstrap resamples with replacement, (2) Compute pass rate \hat{p}_b for each resample $b \in \{1, \dots, B\}$, (3) Define CI as $[\hat{p}_{\alpha/2}, \hat{p}_{1-\alpha/2}]$ where \hat{p}_q is the q -th quantile. For $\alpha = 0.05$, this yields a 95% confidence interval.

Multiple Comparisons. We do not apply multiple-comparison correction (e.g., Bonferroni) in the current implementation. This is a known limitation: when evaluating k models simultaneously, the family-wise error rate increases. Future work will integrate false discovery rate (FDR) control methods.

Inter-Rater Reliability. All 250 fixtures were labeled by 3 annotators following the protocol in Appendix B. Agreement metrics: Cohen’s κ (pairwise) and Fleiss’ κ (3+ annotators) [5, 8]. Overall weighted $\kappa = 0.88$ (substantial agreement). See docs/DATA_CARD.md for complete annotation protocol.

5.3 Main Results

Table 3 presents aggregate results.

Table 3: Validation Results Across Tasks (all CIs: bootstrap percentile, B=1000, 95%)

Task (N)	Mode	Val.	Task Acc.	Repair	Lat. (ms)	OH%*
Classification (410)	None	12%	8%	0%	1,847	2.1
	Struct.	78%	71%	43%	1,923	2.3
	Assist	92%	87%	68%	2,314	2.8
	Enforce	100%	98%	0%	847	1.9
Extraction (287)	None	9%	-	0%	2,108	2.0
	Assist	89%	-	72%	2,541	2.9
	Enforce	100%	-	0%	923	2.1
Summarization (203)	None	31%	-	0%	3,214	1.8
	Assist	74%	-	54%	3,687	2.4
	+Judge	87%	-	61%	4,102	3.1
RAG (187)	Assist	76%	69%	49%	3,301	2.7
	+Judge	81%	74%	53%	3,819	3.3
Tool-calls (160)	Enforce	100%	-	0%	778	1.8

* OH% = Overhead% = (check execution time / total latency) \times 100

CIs (bootstrap percentile, B=1000): Classification (assist, N=10): 95% CI [0.89, 0.94]. Extraction (enforce): [0.98, 1.00]. **Repair:** 68% (classification), 81% fence stripping, 19% lowercasing. Disabling reduces success 92% \rightarrow 34%. **Latency:** Enforce 847ms, assist 2,314ms (2.7 \times). Overhead: <3%.

5.4 Ablation Studies

Sample size N (Table 4): N=3: 85%, N=10: 92%, N=30: 93% (diminishing returns). CI width: 0.12 \rightarrow 0.05 \rightarrow 0.03. **Recommendation:** N=10 (cost-confidence balance).

Table 4: Sample Size Ablation (Classification Task, bootstrap CI)

N	Val.	CI Width*	Var(\hat{p})	Lat. (ms)	Mult.
1	78%	-	-	2,314	1.0 \times
3	85%	0.12	0.0036	6,942	3.0 \times
10	92%	0.05	0.0008	23,140	10.0 \times
30	93%	0.03	0.0003	69,420	30.0 \times

bound (95% bootstrap percentile)

Aggregation (Table 5): Majority (92%) optimal. All (87%): safety-critical. Any (97%): exploratory.

Table 5: Aggregation Policy (N=10)

Policy	Val.	FP	FN	Use Case
first	78%	5%	17%	Baseline
majority	92%	3%	5%	**Production**
all	87%	0%	13%	Safety
any	97%	8%	0%	Exploratory

Repair depth: max_steps=0: 34%, =1: 78%, =2: 92%, =3: 92%. **Rec:** 2. **Tolerance τ :** Optimal $\tau = 0.9$ (F1=0.94).

5.5 Seed Robustness

5 seeds (42, 123, 456, 789, 999): Mean 91.8%, Std 1.2% (empirical), Range [90.3%, 93.1%]. Low variance confirms determinism despite LLM stochasticity.

Table 6: Seed Robustness (Classification, N=10, Assist Mode)

Seed	42	123	456	789	999	Mean	Std*
Val. (%)	92.0	91.5	90.3	93.1	92.0	91.8	1.2
Repair (%)	68	71	74	65	69	69.4	3.1

deviation across 5 seeds

5.6 Comparative Benchmarks

Table 7: PCSL (enforce) F1=0.99, (assist) F1=0.92 vs. CheckList 0.82, Guidance 0.86, OpenAI Struct. 0.97. Setup: PCSL 2 min vs. CheckList 120 min.

Table 7: Framework Comparison (N=50 Shared Fixtures)

Framework	Prec.	Rec.	F1	Setup (min)	Repro.	CI/CD
CheckList	0.89	0.76	0.82	120	Partial	\times
Guidance	0.92	0.81	0.86	30	Manual	\times
OpenAI Struct.	1.00	0.94	0.97	5	Vendor-lock	Limited
PCSL (assist)	0.96	0.88	0.92	2	Full	\checkmark
PCSL (enforce)	1.00	0.98	0.99	2	Full	\checkmark

5.7 Semantic Validation

LLM-judge vs. human (100 outputs, 3 raters, MT-Bench scale [15]):

Table 8: LLM-Judge vs. Human

Judge	Pearson r	Spearman ρ	κ	Agree%	Cost/100
GPT-4o	0.87	0.84	0.82	86%	\$2.40
GPT-4o-mini	0.79	0.77	0.74	81%	\$0.24
Human (inter)	-	-	0.89	91%	\$150

Result: $\kappa = 0.82$ (substantial), 62× cheaper. **ROC:** Similarity AUC=0.91 (threshold=0.82, F1=0.88). Judge AUC=0.89 (rating ≥ 7 , F1=0.85).

5.8 Repair Policy Sensitivity

Motivation. LLMs frequently generate syntactically varied outputs (markdown fences, extra whitespace) that do not affect semantic correctness. Automated repair policies normalize outputs before validation.

Transformations. PCSL applies ordered normalizations: (1) `strip_markdown_fences`: Remove “json” wrappers, (2) `strip_whitespace`: Trim leading/trailing whitespace, (3) `normalize_newlines`: Unify line endings.

Analysis. Table 9 compares validation success with repair enabled vs. disabled.

Table 9: Repair policy impact on validation success. Task accuracy remains invariant (semantics preserved).

Task	w/o Repair	w/ Repair	Δ	Task Acc.
Classification	82%	98%	+16%	94%
Extraction	78%	96%	+18%	91%
Summarization	74%	92%	+18%	87%
Average	78%	95%	+17%	91%

Key Findings: Repair improves validation success by 17% on average. Task accuracy (semantic correctness) is invariant to repair, confirming that transformations preserve meaning. Structured tasks (classification, extraction) benefit more than free-form (summarization). Repair rate of 18% indicates prompts generate syntactically varied but semantically correct outputs. **False Positives:** Repair does not create false positives—genuinely invalid JSON (missing braces, malformed) remains invalid after normalization. The `repair_ledger` tracks all transformations for transparency.

6 Limitations and Future Work

Language and Domain. Our evaluation fixtures are limited to English and primarily cover business/technical domains. Global applicability requires multilingual datasets and culturally diverse contexts.

Annotation Resources. As an open-source project, our annotation budget is constrained. All 250 fixtures were labeled by 3 annotators, but larger-scale evaluation (thousands of examples) would benefit from crowdsourcing platforms with quality control.

Statistical Methods. We use percentile bootstrap without multiple-comparison correction. For comparative studies evaluating many models simultaneously, false discovery rate (FDR) control would be appropriate.

Scope Exclusions. PCSL does not currently address: (1) Fairness and bias testing (requires demographic data and fairness metrics), (2) Adversarial robustness (requires attack generation), (3) Data privacy (requires differential privacy integration), (4) Long-context evaluation (> 4K tokens), (5) Multimodal inputs (images, audio).

Future Directions. (1) Adaptive sampling with stopping rules to reduce API costs, (2) Causal inference for identifying which prompt variations affect success rates, (3) Real-time monitoring with drift

detection, (4) Cross-model benchmarking with standardized leaderboards.

7 Discussion

Limitations. Structural checks dominate; semantic (similarity, judge) depend on embedding/judge quality. Tolerance τ requires domain calibration. Provider non-determinism: 2-3% variance despite seeding. JSON-focused: free-text/multimodal need alternative strategies. Auto-repair 68% risks masking issues; monitor ledger.

Contributions vs. prior work. CheckList: PCSL adds formal spec, probabilistic semantics, CIs. OpenAI Struct.: PCSL provider-agnostic, semantic checks, audit. Guidance: PCSL post-hoc validation with statistical confidence.

Future. Differential testing (drift), multi-turn contracts, adversarial robustness (jailbreak), contract synthesis, adaptive τ learning, causal validation (RAG correctness), fairness/bias.

Review-driven improvements (Table 10):

Table 10: Response to Peer-Review

Criticism	Addressed By
Bootstrap details missing	\$3.2: B=1000, convergence
No seed robustness	\$5.3: 5 seeds, std 1.2%
N-sampling unjustified	\$5.2: N=3/10/30 ablation
No convergence proof	\$3.2: CLT, variance $O(1/N)$
Lacks compositional	\$3.3: Multi-step, RAG
No direct comparison	\$5.4: CheckList/Guidance/OpenAI
Semantic weak	\$5.5: Judge vs. human, $\kappa = 0.82$
Audit abstract	\$4.5: Case study, artifacts
Claims too strong	Abstract: “comprehensive formalization”

8 Conclusion

PCSL v0.3 provides a comprehensive probabilistic formalization for LLM prompt testing. Rigorous evaluation (250 labeled fixtures, 5 tasks) demonstrates 95% validation (assist) with statistical confidence (95% bootstrap CI, B=1000, $\delta=0.95$), seed robustness (empirical std 1.1% across 5 seeds). LLM-judge (GPT-4o) achieves $\kappa = 0.82$ vs. humans. Formal compliance mapping operationalizes ISO 29119 and EU AI Act with practical audit examples.

We have addressed transparency and reproducibility concerns through: (1) comprehensive dataset documentation (250 fixtures with $\kappa=0.88$ inter-rater agreement), (2) fixed seeds (42) and pinned dependencies, (3) detailed statistical methodology (bootstrap parameters, multiple-comparison limitations), (4) practical compliance examples (risk matrices, audit bundles). All code, fixtures (CC BY 4.0), and documentation are publicly available under open licenses (MIT for code), enabling independent verification.

PCSL bridges software testing and AI evaluation, enabling systematic prompt testing, CI/CD integration (<3% overhead), and regulatory auditing. We envision PCSL as a foundational layer for trustworthy LLM deployment, particularly in regulated industries. Open source: <https://github.com/philippmelikidis/prompt-contracts>.

References

- [1] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *FAccT*.
- [2] Harrison Chase. 2023. LangChain. <https://github.com/langchain-ai/langchain>.
- [3] Bradley Efron and Robert J. Tibshirani. 1993. *An Introduction to the Bootstrap*. Chapman and Hall/CRC, New York.

- [4] European Parliament and Council. 2024. Regulation (EU) 2024/1689 on Artificial Intelligence (AI Act). Official Journal of the European Union. Available at: <https://artificialintelligenceact.eu/>.
- [5] Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin* 76, 5 (1971), 378–382.
- [6] Peter Hall. 1992. The Bootstrap and Edgeworth Expansion. (1992).
- [7] ISO/IEC/IEEE. 2013. ISO/IEC/IEEE 29119-1:2013 Software and Systems Engineering – Software Testing – Concepts and Definitions.
- [8] J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics* 33, 1 (1977), 159–174.
- [9] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. In *Proceedings of NeurIPS*.
- [10] Bertrand Meyer. 1992. Applying design by contract. *Computer* 25, 10 (1992), 40–51.
- [11] Microsoft Research. 2023. Guidance. <https://github.com/microsoft/guidance>.
- [12] OpenAI. 2023. Structured Outputs in the API. <https://platform.openai.com/docs/guides/structured-outputs>.
- [13] OpenAPI Initiative. 2017. The OpenAPI Specification. Linux Foundation.
- [14] Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. *Proceedings of ACL* (2020), 4902–4912.
- [15] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging LLM-as-a-judge with MT-bench and Chatbot Arena. *Advances in NeurIPS* 36 (2023).