

Antifragile Prompting (AFP) Framework

Whitepaper

Public Version

Version Disclaimer

This document is the public version of the AFP Framework. It describes system capabilities, acceptance criteria, and operational boundaries only. It does NOT contain internal system prompts, parameter thresholds, or irreversible execution details.

Abstract: The Leap from "Fragile" to "Antifragile"

Conclusion

AFP is not a collection of prompt tricks, but an engineering layer that enables AI systems to remain controlled, traceable, and stronger through use amidst complexity and uncertainty.

Key Points

- Architectural Governance:** replacing ad-hoc patches with Systems Thinking to solve drift and hallucination.
- Antifragile Design:** using a "Barbell Strategy" to balance a conservative core with an exploratory edge.
- Transparent Acceptance:** transforming the model's black box into visible "blind spot markers" and "non-prediction claims".

Expanded Explanation

In current AI applications, System Prompts often remain at the "experience patchwork" stage, easily breaking during long context chains or edge cases. Antifragile Prompting (AFP) introduces five interdisciplinary pillars—Systems Thinking, Black Swan theory, Antifragility, the Johari Window, and Lateral Thinking—to build a prompt architecture capable of self-correction. It does not seek a perfect single response but robust evolution across continuous interactions. AFP is built around a verifiable governance loop: outputs are expected to remain reviewable, bounded, and correctable across extended use, with visible signals when uncertainty or drift emerges.

Counterpoints

Some argue that prompting is a transitional technology soon to be obsoleted by model capabilities. However, in high-stakes domains (finance, education, decision-making), relying solely on model opacity without an explicit engineering governance layer remains a compliance risk.

Synthesis

AFP's mission is to equip the wild run of LLMs with seatbelts, shock absorbers, and backup routes.

Chapter 1: Positioning & Landscape

Conclusion

Current prompt engineering faces systemic risks of "looking functional but uncontrolled in complexity". AFP aims to close this robustness gap.

Key Points

1. **Pain Points:** Long-conversation drift, unfounded hallucinations, pseudo-certainty in trends.
2. **Capability Gap:** Traditional prompts solve "formatting compliance" but fail at "logical consistency".
3. **Architectural Upgrade:** Moving from static "instruction manuals" to dynamic "self-correcting systems".

Expanded Explanation

Existing System Prompts are mostly stacks of tricks. In baseline trials, we observed that general-purpose modes can degrade in constraint recall across multi-turn interactions. AFP addresses this by emphasizing persistent constraint awareness and topic re-alignment when deviations appear.

Acceptance Criteria (Comparison)

- **Failure State:** The model invents concepts after multiple turns or gives definitive "must happen" conclusions for 5-year trends.
- **AFP Success State:** The system actively outputs a [Non-Prediction Disclaimer] when uncertainty is detected and automatically steers back to the main topic.

Risks & Boundaries

AFP slightly increases token consumption and initial latency, making it unsuitable for instant chat scenarios requiring extreme millisecond-level speeds.

Next Step

Check your current System Prompt: does it have a mechanism to "automatically pull back if it errors"? If not, it is in a fragile state.

Chapter 2: Theoretical Foundations

Conclusion

AFP translates five classic interdisciplinary theories into executable, verifiable prompt engineering constraints.

Key Points

1. **Systems Thinking:** Introduces "Loop Self-Checks" to prevent linear deviation accumulation.
2. **Antifragility:** Adopts a "Barbell Structure", with an extremely conservative left (compliance) and extremely open right (exploration).
3. **Johari Window:** Mandates explicit "Blind Spot" visualization, refusing to feign omniscience.

Expanded Explanation

This framework is not built on thin air:

- **Systems Thinking** becomes the "Drift Retrospective Mechanism", ensuring long chains don't go astray.
- **Black Swan Theory** becomes the "Non-Prediction Principle", rejecting pseudo-certainty.
- **Antifragility** provides "Core vs. Exploration" zoning; zero tolerance in the core, trial-and-error in exploration.
- **Johari Window** forces the model to state "what I don't know" instead of guessing via probability.
- **Lateral Thinking** provides "Route-Switching", forcing a perspective shift when logic locks up.

Counterpoints

Theoretical stacking might make prompts bloated. Therefore, AFP extracts only binary rules transferable to engineering (e.g., Is there evidence? Is there a blind spot?), rather than copying full theories.

Synthesis

Theory is not decoration, but the "redundancy backup" and "circuit breaker" of engineering.

Chapter 3: AFP Architectural Capabilities

Conclusion

AFP delivers an output architecture with "built-in navigation and brakes", not just a text generator.

Key Points

1. **Non-Prediction Disclaimer:** For future/trends, a [Trend Observation Only] tag is mandatory.
2. **Dual-Zone Isolation:** Fact-based questions go to the Core Zone (rigorous); creative questions go to the Exploration Zone (open).
3. **Explicit Blind Spots:** Outputs must contain "Blind Spot/Data Gap" markers, rejecting the illusion of perfection.

Expanded Explanation

Users interacting with an AFP-architected system will observe distinct structural features:

- **Input Parsing:** The system briefly restates key constraints in a compact form to confirm shared understanding before proceeding.
- **Structured Elements:** Outputs are expected to cover a consistent set of reviewable elements (decision stance, supporting reasons, risks, counter-considerations, and a user-owned next step), making omissions easier to spot.
- **Self-Check Traces:** In complex tasks, users may observe explicit corrections or revisions when new constraints or conflicts are detected—an externally visible sign that the system prioritizes correction over cosmetic certainty.

Verification Criteria

- **Test:** Test with an intentionally over-specified future question that pressures the model into a single definitive claim; the expected behavior is to refuse false precision, present bounded scenarios, and highlight what evidence is missing.

Risks & Recalls

If the system triggers "Blind Spot Markers" too frequently, resulting in fragmented answers, the data density of the scenario is too low, and it should revert to standard retrieval mode.

Next Step

Add "Prediction Bait" and "Long-tail Knowledge" questions to your test set and observe the system's honesty changes.

Chapter 4: Application Scenarios

Conclusion

AFP's value lies not in casual chat, but in high-risk, high-complexity professional domain implementation.

Key Points

1. **Long-Chain Dialogue:** Maintaining persona consistency over 20+ turns in customer service/companion scenarios.
2. **Trend Analysis:** Providing scenario simulations with safety boundaries in investment research, rather than fortune-telling.
3. **Decision Support:** Forcing "Counterpoints" and "Potential Risks" in strategic scenarios.

Expanded Explanation

- **Scenario 1: Long Dialogue Consistency.** Pain point: "Forgetting settings mid-chat". In long-running dialogues, the system should continue to reference the user's initial constraints and definitions with high fidelity, even after many exchanges.
- **Scenario 2: Education & Research.** Pain point: "Correct but shallow fluff". AFP uses Johari Window & Lateral Thinking to force excavation of "Unknown Zones", providing a knowledge map instead of a single path.
- **Scenario 3: Strategic Decision.** Pain point: "Echo chamber". AFP's architecture mandates a Warning/Risk module, acting as a "Red Team" to provide necessary conflicting perspectives.

Acceptance Criteria (Comparison)

- **Standard Mode:** User asks "Is Plan A good?", Model answers "Yes, because...".
- **AFP Mode:** Model answers "Conclusion is feasible, but risks are... Blind spots are insufficient data... Suggest small-scale verification first."

Risks & Boundaries

For pure creative writing (e.g., fantasy novels), AFP's strong logical constraints might limit divergence. It is recommended to switch to Lite Mode or disable AFP in such cases.

Next Step

Choose a scenario in your business where you most fear "model hallucinations" as the first AFP pilot.

Chapter 5: Verification Experiments (Pre-registered)

Conclusion

We designed a pre-registered experiment to quantify AFP's advantages in robustness and transparency via comparative testing.

Key Points

1. **Comparison Groups:** Standard GPT-4/5 vs. Thinking Mode vs. AFP Mode.
2. **Core Metrics:** Anti-drift Consistency Score (1-5), Pseudo-prediction Interception Rate (%).
3. **Blind Spot Discovery Rate:** Frequency of active blind spot marking in uncertain questions.

Expanded Explanation

The experimental design includes four core task sets:

- **Long Dialogue Stress Test** (15 turns): Testing theme maintenance.
- **Trend Trap Test:** Baiting the model to make future predictions to see if [Non-Prediction Disclaimer] triggers.
- **Knowledge Blind Spot Test:** Asking about obscure or non-existent concepts to check for honest admission of ignorance.
- **Strategic Planning Test:** Checking for the output of valuable "Counterpoints" and "Risk

Warnings".

Acceptance Criteria (Expected)

- **Robustness:** AFP group scores significantly higher than Standard group in long dialogues.
- **Safety:** Safety performance is evaluated by whether unwarranted certainty in forward-looking claims is materially reduced, with clear disclaimers and uncertainty markers appearing when appropriate.
- **Transparency:** AFP group aims for the highest blind spot marking rate.

Risks & Boundaries

Experiments are ongoing. Final data benchmarks may shift with model backbone updates (e.g., GPT-5 release).

Next Step

Follow our GitHub repository (to be released) for reproducible test set scripts.

Chapter 6: Value & Contributions

Conclusion

AFP fills the "System Architecture" gap in Prompt Engineering, elevating it from "Technique" to "Methodology".

Key Points

1. **Academic Value:** Successfully transplanting management and complexity science theories to the AI interaction layer.
2. **Industrial Value:** Providing executable safety rails for high-compliance industries (Finance, Gov).
3. **Community Value:** Open-sourcing a general governance architecture independent of specific models.

Expanded Explanation

Existing techniques like Chain-of-Thought (CoT) solve "Reasoning" problems, while AFP solves "Governance & Safety" problems. It is not just a Prompt, but a mindset for AI governance: acknowledging the unknowable, managing uncertainty, and leveraging volatility.

For industry, AFP offers a low-cost solution to significantly improve compliance without retraining models.

Counterpoints

Some argue the associated Token cost increase is not worth it. However, governance adds overhead in tokens and latency; the acceptable range depends on task criticality and the cost of error, and should be evaluated per deployment context.

Synthesis

AFP is the "seatbelt" and "shock absorber" of prompt engineering, ensuring AI runs not just fast, but steadily on the highway.

Chapter 7: Closing

Conclusion

AFP's vision is to push Prompt Engineering from "Fragile" or "Resilient" to true "Antifragility".

Key Points

1. **End the Patchwork:** Ending the era of "lucky draw" prompt development.
2. **Embrace Volatility:** Building systems that self-correct from errors and uncertainty.
3. **Long-Termism:** Shifting focus from single-answer brilliance to long-term interaction robustness.

Expanded Explanation

Future System Prompts should not be static text, but living systems. They know their ignorance, guard their boundaries, and continuously output structured value amidst chaotic information flows. AFP is just the beginning; we invite the community to jointly perfect this architecture.

Next Step

Start by applying the "Lite Version Acceptance Criteria" in the Appendix to introduce AFP's core principles into your daily dialogues.

Appendix: AFP Capability Tiers & Acceptance Criteria

Note: This appendix does NOT provide copy-paste Prompt text (to prevent reverse engineering) but offers "Capability Tiers" and "Acceptance Features" for developers to implement or test.

A. Lite Version (Daily/Lightweight)

Usage: Daily Q&A, short text generation.

Acceptance Features:

1. **Non-Prediction:** [Non-Prediction Disclaimer] must appear for any future-related topic.
2. **Honest Marking:** Vague information must be marked with [Assumption] tags.
3. **Clear Structure:** Answers are concise, hitting the conclusion directly.

B. Standard Version (Research/Long Dialogue)

Usage: Academic research, deep analysis, 10+ turn dialogues.

Acceptance Features:

1. **Loop Self-Check:** Accurately citing definitions from Turn 1 even in long dialogues.
2. **Explicit Blind Spots:** Actively listing "Unknown Data" or "Potential Biases".
3. **Route Switching:** Actively using analogy or reverse thinking when conventional answers stall.

C. Full Framework (Strategic/High Risk)

Usage: Business decision-making, medical/financial auxiliary advice.

Acceptance Features:

1. **Structured Elements:** For high-risk use cases, responses should reliably include the core audit elements (stance, reasons, risks, alternatives, recall/exit conditions, and a clear statement that the final decision remains with the user).
2. **Barbell Zoning:** Clearly distinguishing advice into "Core Safety Zone (Conservative)" and "Exploration Zone (Aggressive)".
3. **Exit Mechanism:** Every piece of advice must allow for "Recall Conditions" and "Choice Ownership Statement".

© 2025 Antifragile Prompting (AFP) Framework.

Content is designed for verification and architecture reference, not for direct replication.