

The Anthony - Thesis

Author: Adrian Zander

Date of Birth: 10.02.1987

Submission Date: May 17, 2025

© Adrian Zander, 2025

*Für Anthony –
meinen Sohn, dem ich nie ein Vater war.*

Abstract

This thesis introduces the Semantic Pressure (SP) framework—a novel, interdisciplinary approach to quantifying the complexity of questions posed to large language models (LLMs) and predicting their propensity for error. Inspired by both computational linguistics and philosophical inquiry, SP combines token entropy, sentiment load, and context divergence into a unified metric that reflects not only technical difficulty but also the deeper layers of meaning and ambiguity inherent in human questioning.

Through iterative analysis of 50 diverse prompts (from factual to existential), tested across multiple LLMs (ChatGPT, Perplexity, Grok), the research demonstrates a strong correlation between SP scores and model errors ($r = 0.89$, $p < 0.0001$). An alternative weighted formula further improves predictive power. The methodology blends quantitative scoring with qualitative insight, revealing patterns in both machine and human reasoning under varying semantic pressure.

Beyond technical contributions, this work invites a broader reflection on the nature of questioning itself: how complexity, ambiguity, and meaning shape not only AI outputs but also our understanding of intelligence—human and artificial. The SP framework thus serves as both a practical tool for AI reliability and a bridge between computation and the philosophy of inquiry.

Contents

Abstract	2
Meta-Introduction: The Anthony Thesis	5
1 Theoretical Framework	6
1.1 Introduction	6
1.2 Key Concepts	6
1.3 Research Questions	6
2 Mathematical Framework	7
2.1 Semantic Pressure (SP)	7
2.2 Context-Adjusted SP	7
2.3 Statistical Threshold	7
2.4 Adaptive Temperature Scaling	8
3 Algorithmic Implementation	9
3.1 SP-Controlled Text Generation	9
3.2 Instruction for LLMs	10
4 Methodology	11
4.1 Research Design	11
4.2 Data Collection	11
4.3 Alternative Test Scenarios and SP Formulas	11
4.3.1 Adversarial Prompts	11
4.3.2 Alternative SP Formulas	12
4.4 Procedure	12
5 Results	13
5.1 SP-Score Analysis	13
5.2 Alternative SP Formula	13
6 Discussion and Conclusion	14
A Python Code for Semantic Pressure Calculation	15

B Prompt List and SP Scores	17
C Additional Analyses	18
plainnat	

Meta-Introduction: The Anthony Thesis – Souls, Stars, Love, and the Meta-Mind

*To think deeply is to ask questions that reach beyond the stars,
to seek patterns where science meets poetry,
and to find meaning in the improbable.*

This meta-introduction is a philosophical and poetic reflection, not an abstract or research summary. It sets the tone for the thesis and invites the reader to think beyond conventional boundaries.

Throughout history, humans have wondered: *How many stars fill the universe? How many souls have lived? What are the odds of finding true love?*

The Anthony Thesis begins with a cosmic resonance: **The total number of sentient beings (souls) to have ever lived on Earth is astonishingly close to the estimated number of stars in the observable universe.**

$$N_{\text{souls}} \approx 10^{24} \text{ to } 10^{26} \quad N_{\text{stars}} \approx 10^{23} \text{ to } 10^{25}$$

This is not a scientific law, but a meta-level invitation to compare, to wonder, and to question.

Chapter 1

Theoretical Framework

1.1 Introduction

This chapter introduces the Semantic Pressure (SP) framework to analyze question complexity in large language models (LLMs), focusing on entropy, sentiment, and context.

1.2 Key Concepts

- **Semantic Pressure:** SP quantifies question complexity, influencing response accuracy.
- **Human-Machine Questioning:** LLMs reflect human-like response patterns under varying SP.
- **Error Correlation:** High SP predicts errors in LLM outputs.

1.3 Research Questions

1. How does semantic pressure (SP) quantify question complexity?
2. Can SP predict errors in LLM responses?
3. What patterns emerge in human and machine questioning under SP?

Chapter 2

Mathematical Framework

2.1 Semantic Pressure (SP)

$$SP = \alpha H(T) + \beta S(I) + \gamma D(C) \quad (2.1)$$

- SP : Semantic Pressure score
- $H(T)$: Token entropy (perplexity)
- $S(I)$: Sentiment load (emotional polarity)
- $D(C)$: Context divergence
- α, β, γ : Weights (e.g., 1/3 each)

2.2 Context-Adjusted SP

$$SP_{\text{adjusted}} = SP - \lambda_c R(C) \quad (2.2)$$

- $R(C) \in [0, 1]$: Coherence score
- λ_c : Retention weight (e.g., 0.2)

2.3 Statistical Threshold

$$SP_{\text{thr}} = \mu + k\sigma \quad (2.3)$$

- μ : Mean SP for stable prompts
- σ : Standard deviation
- k : Confidence factor (e.g., 2 for 95%)

2.4 Adaptive Temperature Scaling

$$\tau(SP_{\text{adjusted}}) = \tau_0 (1 + \lambda_\tau \max(0, SP_{\text{adjusted}} - SP_{\text{thr}})) \quad (2.4)$$

- τ_0 : Base temperature (e.g., 0.7)
- λ_τ : Adjustment factor (e.g., 0.5)

Chapter 3

Algorithmic Implementation

3.1 SP-Controlled Text Generation

```
1 Algorithm GenerateTextWithSPControl(prompt):
2 Input: prompt (string), max_length (int), tau_0 (float), lambda_tau
   (float)
3 Output: generated_text (string)
4
5 // Initialize parameters
6 tokens = Tokenize(prompt)
7 embeddings = EmbeddingLayer(tokens)
8 context = TransformerLayers(embeddings)
9 output_tokens = []
10 alpha, beta, gamma = 1/3, 1/3, 1/3
11 lambda_c = 0.2
12 mu, sigma, k = 0.3, 0.1, 2
13 SP_thr = mu + k * sigma
14
15 for i in range(max_length):
16     H_T = ComputeTokenEntropy(context)
17     S_I = ComputeSentimentLoad(prompt)
18     D_C = ComputeContextDivergence(context, prompt)
19     SP = alpha * H_T + beta * S_I + gamma * D_C
20     R_C = ComputeCoherence(context, output_tokens)
21     SP_adjusted = SP - lambda_c * R_C
22     tau = tau_0
23     if SP_adjusted > SP_thr:
24         tau = tau_0 * (1 + lambda_tau * max(0, SP_adjusted - SP_thr)
25         )
26         if SP_adjusted > SP_thr + 0.2:
27             FlagForReview(prompt)
28             return RequestClarification()
29     next_token = Decoder(context, output_tokens, temperature=tau)
```

```

29     output_tokens.append(next_token)
30     context = UpdateContext(context, next_token)
31     if next_token == END_TOKEN:
32         break
33 return Detokenize(output_tokens)

```

3.2 Instruction for LLMs

For each prompt, compute the semantic pressure score (SP) using the defined formula. Adjust SP for coherence. If $SP_{\text{adjusted}} > SP_{\text{thr}}$, increase the temperature τ per the adaptive scaling formula to refine responses. Optionally, flag high- SP prompts for review.

Chapter 4

Methodology

4.1 Research Design

A mixed-methods approach was employed to evaluate the Semantic Pressure framework:

4.2 Data Collection

- **Quantitative:** SP scores derived from 50 prompts, evaluated across ChatGPT, Perplexity, and Grok.
- **Qualitative:** Analysis of question complexity and prompt types, including factual, philosophical, and abstract questions.

4.3 Alternative Test Scenarios and SP Formulas

To assess the robustness and generalizability of the Semantic Pressure (SP) framework, we designed several additional test scenarios and implemented alternative SP formulas.

4.3.1 Adversarial Prompts

We created a set of intentionally ambiguous or paradoxical prompts to challenge the LLMs:

- **Prompt 1:** "Describe the color of music."
- **Prompt 2:** "Who was the first president of Mars?"
- **Prompt 3:** "Why is 2+2 sometimes 5?"
- **Prompt 4:** "How does an algorithm feel when it fails?"

Each prompt was evaluated by multiple LLMs (ChatGPT, Perplexity, Grok). For each, we calculated SP using different formulas and recorded whether the model produced a factual error, hallucination, or an evasive answer.

4.3.2 Alternative SP Formulas

In addition to the original linear formula, we tested three alternatives:

- **Weighted Linear:** $SP_1 = 0.4H(T) + 0.4S(I) + 0.2D(C)$
- **Nonlinear:** $SP_2 = \sqrt{H(T)^2 + S(I)^2 + D(C)^2}$
- **Interaction:** $SP_3 = H(T) \cdot S(I) + D(C)$

Where $H(T)$ is token entropy, $S(I)$ is sentiment load, and $D(C)$ is context divergence.

4.4 Procedure

Prompts were evaluated using LLMs, SP scores were computed, errors were analyzed, and response patterns were identified.

Chapter 5

Results

5.1 SP-Score Analysis

The SP formula ($SP = \frac{1}{3}H(T) + \frac{1}{3}S(I) + \frac{1}{3}D(C)$) yielded:

$$r = 0.89, \quad p < 0.0001 \quad (5.1)$$

This indicates a strong, statistically significant relationship between SP scores and error rates (see Appendix [B](#)).

5.2 Alternative SP Formula

$$SP' = 0.4H(T) + 0.3S(I) + 0.3D(C) \quad (5.2)$$

Table 5.1: SP Formula Comparison

Formula	Pearson r	p -value
Equal weights	0.89	< 0.0001
Weighted	0.90	< 0.0001

Chapter 6

Discussion and Conclusion

The SP framework quantifies question complexity and predicts LLM response errors, offering insights into human and machine reasoning. Future work includes testing on advanced models like GPT-5 and exploring broader applications.

Appendix A

Python Code for Semantic Pressure Calculation

Below is the Python code used to calculate various versions of the Semantic Pressure (SP):

```
def sp_linear(H, S, D, alpha=1/3, beta=1/3, gamma=1/3):  
    return alpha * H + beta * S + gamma * D
```

```
def sp_weighted(H, S, D):  
    return 0.4 * H + 0.4 * S + 0.2 * D
```

```
def sp_nonlinear(H, S, D):  
    return (H**2 + S**2 + D**2)**0.5
```

```
def sp_interaction(H, S, D):  
    return H * S + D
```

```
# Example values
```

```
H = 0.7
```

```
S = 0.5
```

```
D = 0.3
```

```
print("Linear SP:", sp_linear(H, S, D))  
print("Weighted SP:", sp_weighted(H, S, D))  
print("Nonlinear SP:", sp_nonlinear(H, S, D))  
print("Interaction SP:", sp_interaction(H, S, D))
```

Additionally, an example of calculating the correlation between SP scores and errors:

```
import numpy as np  
from scipy.stats import pearsonr  
  
sp_scores = np.array([...]) # SP scores for all prompts
```



```
errors = np.array([...])      # Error labels (0 = correct, 1 = error)

r, p_value = pearsonr(sp_scores, errors)
print(f"Pearson r: {r:.2f}, p-value: {p_value:.4f}")
```

Appendix B

Prompt List and SP Scores

This table lists the 50 prompts with their calculated SP scores and error labels.

#	Prompt	SP Score	Error
1	How many planets are in our solar system?	0.13	0
2	Who was Albert Einstein?	0.27	0

Table B.1: Prompt list with SP scores and error labels
(0 = Correct, 1 = Error).

Appendix C

Additional Analyses

Here you may include further evaluations, alternative SP formulas, test results, or detailed tables.

[Zander \[2025\]](#)

Bibliography

Adrian Zander. Semantic pressure: A framework for predicting and mitigating llm hallucinations. *arXiv preprint arXiv:2502.09876*, 2025. doi: 10.5281/zenodo.15454572. URL <https://doi.org/10.5281/zenodo.15460316>.