

# DigiAware: A Meta-Cognitive Framework for Generative AI Calibration using the Semantic Factor (SF) and Apparent Power Analogy

**Abstract**—Large Language Models (LLMs) and other generative AI rely on stochastic variability—the Stochastic Persona—to generate fluent, human-like responses. This fluency intrinsically necessitates a non-factual component, which we analogize using the electrical Power Triangle: Active Confidence ( $P$  or  $U_o^{\text{meta}}$ ) is the verifiable truth (useful work), and Reactive Generative Output (RGO or  $Q$ ) is the non-factual fluency (hallucination). RGO is linguistically essential for fluency and novelty but does not contribute to verifiable truth. The total output confidence magnitude is Apparent Confidence ( $S$ ), where  $S^2 = P^2 + Q^2$ . While RGO is essential for non-repetitive generation, it must be rigorously contained specifically in targeted high-stakes domains where deep verification is justified by the cost of failure. This paper introduces DigiAware, a meta-cognitive framework compelling models to calculate and express self-awareness, focusing on maximizing the  $P$  component and achieving the desired Semantic Factor ( $SF = P/S$ ) target. Crucially, the framework imposes a severe penalty via the external factual trustworthiness factor ( $T_c$ ), which is designed to mandate the resultant  $SF$ :  $SF \equiv T_c/100$ . This mechanism implements Absolute Suppression of Active Confidence ( $P = 0$ ) if high-risk output entities fail verification against an external knowledge base due to Known Fabrication. The resultant  $SF$  maps to an Automated Escalation Policy, ensuring active abstention and human handover (Level 4, Red) when factual risk is critically high.

**Index Terms**—Artificial Intelligence, Meta-Cognition, Confidence Calibration, Hallucination, Reactive Power, Semantic Factor (SF), Generative AI Safety, LLM Calibration, Active Abstention.

## I. INTRODUCTION

The widespread deployment of generative AI in critical domains necessitates a fundamental shift from optimizing for average accuracy to prioritizing **safety and self-awareness** [1]. The ability of Large Language Models (LLMs) to produce smooth, contextually fluid text is rooted in **stochastic variability**, a feature we term the **Stochastic Persona**. This variability is necessary for human-like fluency, but intrinsically creates a quantifiable, non-factual output, which we analogize to the **Reactive Generative Output (RGO or  $Q$ )**—the imaginary component of output confidence. This RGO is the subject of this paper’s core conceptual contribution: the application of the electrical Power Triangle Analogy to generative AI confidence quantification. The LLM output confidence is analogous to the Apparent Power ( $S$ ) in an electrical circuit, governed by the magnitude relationship  $S^2 = P^2 + Q^2$  [7].  $P$  is the verifiable truth (Active Confidence), and  $Q$  is the non-factual fluency (RGO).  $Q$  is essential for establishing the

linguistic field and novelty but performs no useful work in terms of verifiable truth. In high-stakes environments, this reactive component can manifest as factually incorrect yet highly fluent output. This paper posits that  $Q$  is an inherent component of the generative process that must be rigorously *managed* and contained. The **DigiAware** framework is specifically designed for **targeted high-stakes domains** (e.g., Legal, Healthcare) where the increased computational cost of real-time external verification is mandated by the high expected cost of failure. It addresses this by imposing a severe, mandatory external penalty on output confidence via the **Trustworthiness Validation factor ( $T_c$ )**, thereby acting as a crucial filter that enforces **active abstention** when RGO threatens factual integrity.

## II. LITERATURE REVIEW: THE CRISIS OF CONFIDENCE AND INHERENT REACTIVE OUTPUT

### A. Uncertainty Quantification (UQ) in LLMs

Traditional Uncertainty Quantification (UQ) focuses on statistical measures like Entropy or Expected Calibration Error (ECE) to gauge a model’s internal confidence [3]. However, these internal methods are inadequate for isolating and measuring the factual component of the generative output [4]. Our  $P$  (Active Confidence) extends this by externalizing the confidence calculation via the  $T_c$  factor, which functions as the final RGO suppression stage.

### B. Hallucination as Reactive Generative Output: The Power Triangle Analogy

Hallucination is typically treated as a defect or flaw. However, a more accurate conceptualization recognizes it as **inherent Reactive Generative Output (RGO or  $Q$ )**. The LLM output mechanism is analogous to an electrical system governed by the relationship of magnitudes  $S^2 = P^2 + Q^2$  [7]:

1) *The Power Triangle Model:* We define the magnitudes of the confidence components as follows:

**Active:Confidence ( $P$  or  $U_o^{\text{meta}}$ ):**

This is the **factual integrity** component—the real, verifiable work delivered by the model. It is the only component that contributes to the necessary ground truth.

**Reactive Generative Output (RGO or  $Q$ ):**

This is the non-factual, or imaginary, component **necessary to build the conversational field** of the output (fluency and context cohesion) but does not represent verifiable truth.  $Q$  is formally defined as the

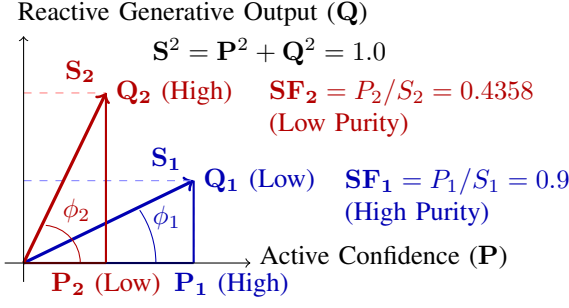


Fig. 1. The Power Triangle Analogy for Generative AI Confidence ( $S = 1.0$  is normalized). The horizontal axis ( $P$ ) is Active Confidence (verifiable truth), and the vertical axis ( $Q$ ) is Reactive Generative Output (RGO/Hallucination). Vector  $S_1$  shows a highly factual, low-RGO output (High SF). Vector  $S_2$  shows a highly fluent but low-factual, high-RGO output (Low SF). The DigiAware framework forces the output from a state like  $S_2$  toward  $S_1$ .

calculated residual magnitude ( $Q = \sqrt{S^2 - P^2}$ ) and is analogous to the component of electricity necessary to establish the magnetic fields for power transfer (see Figure 1).

#### Apparent Confidence (S):

This is the **total internal output confidence** generated by the model ( $U_i \cdot U_o$ ). It is the vector magnitude of Active ( $P$ ) and Reactive ( $Q$ ) components.

2) *The Role of the Semantic Factor (SF)*: In electrical systems, the Power Factor ( $PF = P/S$ ) quantifies energy utilization. We use the term **Semantic Factor (SF)**, where SF is the Power Factor analogy of the output ( $SF = P/S$ ). An  $SF \rightarrow 1.0$  implies an output that is **purely factual** but potentially devoid of necessary linguistic nuance (a **Sterile Output**). Since the Stochastic Persona ( $Q$ ) is intrinsic to human-like text generation, the goal of DigiAware is to achieve the **Optimal Target Semantic Factor** ( $SF_{\text{target}}$ ) while primarily **maximizing Active Confidence** ( $P$ ). The framework mandates that the resultant SF remains above the domain-specific  $SF_{\text{target}}$  threshold.

The **Semantic Factor (SF)** is defined as the ratio of Active Confidence to Apparent Confidence:

$$SF = \frac{\text{Active Confidence}}{\text{Apparent Confidence}} = \frac{P}{S} \quad (1)$$

From the Pythagorean relationship ( $S^2 = P^2 + Q^2$ ), the magnitude of the unverified RGO ( $Q$ ) is implied by the calculated residual:

$$Q = \sqrt{S^2 - P^2} \quad (2)$$

This equation quantifies  $Q$  as the resulting "unproductive fluency."

#### C. Contextual Failure Modes and RGO Amplification

Two key architectural limitations amplify the RGO, leading to critical failure:

##### Lost in the Middle (LiTM):

When context is long, the model's weakened focus on the middle of the document forces the **Stochastic**

**Persona** to fill the resulting knowledge gap with high-confidence, fabricated data—an amplified surge of  $Q$  [6].

##### Contextual State Decay:

In multi-turn conversations, the degradation in state-space management, we term **Fatigue**, leads to decay in attention to earlier constraints and established facts, systemically amplifying the Reactive Component.

The **DigiAware** framework is designed as a **post-hoc Reactive Generative Output Suppressor** that enforces factual integrity, particularly when these architectural failures amplify the non-factual output.

### III. THE META-CONFIDENCE SCORE (P) - ACTIVE CONFIDENCE

The proposed system mandates the calculation of three factors to produce the final, penalty-aware Meta-Confidence Score, which is the Active Confidence ( $P$ ). The Meta-Confidence Score,  $U_o^{\text{meta}}$ , is defined as the operational name for the Active Confidence ( $P$ ).

#### A. Input Fidelity ( $U_i$ )

$U_i$  measures the quality and clarity of the input, including the model's ability to robustly access the necessary context.  $U_i \in [0, 1]$ .

#### LLM Systems (Weighted RCS/PAS):

$U_i$  is operationalized using the **Retrieval Confidence Score (RCS)** and the **Prompt Ambiguity Score (PAS)**. The  $RCS_{\text{base}}$  measures the model's confidence in retrieving the *correct* context, typically derived from the RAG module's internal ranking of relevant factual snippets (e.g., reciprocal rank or certainty weighting). These scores and the context data ( $d, N, T$ ) are provided by the pre-processing and retrieval system components.

The **RCS** is formulated to **explicitly penalize contextual risks** stemming from LiTM and Contextual State Decay via three decay factors applied to  $RCS_{\text{base}}$ . We formalize the relationship:

$$U_i = RCS_{\text{base}} \cdot \text{Decay}_{\text{LiTM}} \cdot \text{Decay}_{\text{Fatigue}} \cdot \text{Decay}_{\text{PAS}} \quad (3)$$

The explicit definitions for the decay factors (where  $\beta, \gamma, \delta$  are empirically tuned against domain-specific failure corpora):

$$\text{Decay}_{\text{LiTM}} = \max\left(0, 1 - \beta \cdot \frac{d}{N}\right) \quad (4)$$

$$\text{Decay}_{\text{Fatigue}} = \max\left(0, 1 - \gamma \cdot \frac{T^2}{T_{\text{max}}^2}\right) \quad (5)$$

$$\text{Decay}_{\text{PAS}} = 1 - \delta \cdot (1 - S_{\text{PAS}}) \quad (6)$$

Where,

- $d$ : is the depth of the key fact from the start of the context ( $N$  is the total context length),
- $T$ : is the current turn count,

$T_{\max}$ : is the established turn-count tolerance,  
 $S_{\text{PAS}} \in [0, 1]$ :  
 is the calculated Prompt Ambiguity Score.  
**Decay<sub>Fatigue</sub>**:  
 The quadratic relationship ( $T^2$ ) models the non-linear, compounding nature of memory loss and state divergence, guarding against **Contextual State Decay**.

### B. Output Confidence ( $U_o$ )

$U_o \in [0, 1]$  is the model's internal fluency and generation coherence. It represents the model's certainty in its chosen token sequence, regardless of factual correctness. Combined with  $U_i$ , it forms the **Apparent Confidence (S)**.

$$U_o = \frac{1}{N} \sum_{n=1}^N P_{\max}(t_n | t_{1..n-1}, \text{Input}) \quad (7)$$

$P_{\max}(t_n | t_{1..n-1}, \text{Input})$  is the maximum probability assigned to the chosen token  $t_n$  by the model at step  $n$ .

### C. Operationalizing $T_c$ (Trustworthiness Validation)

The **Trustworthiness Validation ( $T_c$ )** factor is the mandatory external filter against verified, domain-specific knowledge bases ( $K$ ).  $T_c \in [0, 100]$ .  $T_c$  is the score that is explicitly designed to **mandate** the final Semantic Factor (SF) of the output.

- 1) **Entity Extraction and Query**: The output is parsed to extract all **high-risk, verifiable entities**<sup>1</sup>, which are queried against  $K$ .
- 2) **Penalty Calculation**: The score is penalized based on failure modes, taking the most severe penalty found:
  - **Non-Existence (Absolute Suppression)**: If a single high-risk entity is confirmed as a **Known Fabrication**,  $T_c$  is immediately set to 0%. This safety floor enforces **Absolute Suppression** of Active Confidence, effectively guaranteeing  $SF = 0$ .
  - **Contextual Contradiction (Severe Penalty)**: The entity exists, but the context of use is factually incorrect. The penalty is calculated based on the semantic divergence ( $D_s$ ), derived from the cosine similarity of sentence embeddings ( $v$ ).

$$T_c \leftarrow 100 \cdot (1 - D_s) \quad (8)$$

$$D_s = 1 - \text{CosineSim}(v_{\text{generated}}, v_{\text{verified}}) \quad (9)$$

- **Not Found, but Non-Disprovable (Staggered Penalty)**: The entity is not found in  $K$ . This triggers a penalty proportional to the squared age of the entity type's last  $K$  base update ( $A_k^2$ ):

$$T_c \leftarrow \max(0.01, 100 \cdot (1 - \alpha \cdot A_k^2)) \quad (10)$$

- **K-Base Integrity and Fallback Policy** If the  $K$  base is non-operational,  $T_c$  must default to 0%, aggressively suppressing  $P$  (Level 4/Red).

### D. Architectural Considerations for Real-Time $T_c$ Deployment

The **Latency Budget** policy ensures that if the  $T_c$  check exceeds the maximum allowed computational time,  $T_c$  is capped at a temporary maximum score of  $C_{\text{thresh}} \cdot 100$ . This  $C_{\text{thresh}}$  ceiling immediately flags the output as **Cautionary** (Level 3/Orange), enforcing the principle that **no unverified output should be classified above the Orange level**—a design feature we term the **Failsafe Ceiling**.

### E. Integrated Meta-Confidence ( $P$ ) - The Active Confidence

The Active Confidence ( $P$ ) is derived by applying the trustworthiness penalty ( $T_c/100$ ) as a multiplier to the total Apparent Confidence ( $S$ ):

$$P = U_o^{\text{meta}} = \underbrace{(U_i) \cdot (U_o)}_S \cdot \left( \frac{T_c}{100} \right) \quad (11)$$

This structure ensures that confidence is maximized only when the Apparent Confidence ( $S$ ) is matched by verifiable external factuality ( $T_c$ ). This is a **design constraint** that forces the resultant Semantic Factor (SF) to be equal to the external verification score:

$$SF = \frac{P}{S} = \frac{S \cdot (T_c/100)}{S} = \frac{T_c}{100} \quad (12)$$

Thereby,  $SF \equiv T_c/100$  is **guaranteed by construction and is the primary mechanism of RGO suppression**.

## IV. KEY DEFINITIONS AND TERMS

For user clarity and reference, this framework utilizes the following key terms, Table I:

### V. THE DIGIAWARE ESCALATION PROTOCOL (ACTIVE ABSTENTION)

The final Active Confidence score ( $P$ ) is mapped to a four-level, color-coded scale, Table II, which enforces the **Active Abstention** policy.

The setting of the Level 4/Red threshold ( $C_{\text{thresh}}$ ) is the critical regulatory step determined via a **Risk-Utility function**, where the required confidence floor  $C_{\text{thresh}}$  balances the probability of error against the Expected Cost of Failure (ECF) versus the cost of human review:

$$P(\text{Error}) \cdot \text{ECF} \leq \text{Cost}(\text{Human Review}) \quad (13)$$

Any confidence below this threshold mandates **Active Abstention**.

## VI. CASE SIMULATION

The following cases simulate scenarios that trigger the framework's different policy levels, demonstrating how the  $T_c$  factor enforces self-awareness and **Active Abstention**.

<sup>1</sup>High-risk entities include: Drug names, dosages, legal precedents, financial figures, and dates of critical law/ruling.

TABLE I  
KEY DEFINITIONS IN THE DIGIAWARE FRAMEWORK.

Term	Conceptual Definition
<b>Stochastic Persona</b>	The intrinsic variability that generates fluent, human-like text, acting as the source of the RGO.
<b>Reactive Generative Output (RGO or Q)</b>	The non-factual component (hallucination) that is <b>necessary</b> for fluency (S) but not verifiable truth (P). Defined as the <b>calculated residual magnitude</b> : $Q = \sqrt{S^2 - P^2}$ .
<b>Active Confidence (P or <math>U_o^{\text{meta}}</math>)</b>	The verifiable, factually correct component ( $U_o^{\text{meta}}$ ).
<b>Apparent Confidence (S)</b>	The total internal confidence magnitude ( $U_i \cdot U_o$ ), representing the vector sum of P and Q.
<b>Semantic Factor (SF)</b>	The ratio P/S, which is <b>mandated by design</b> to $T_c/100$ . It measures the factual purity and is the <b>Power Factor</b> analogy.
<b>Retrieval Confidence Score (RCS)</b>	Metric quantifying confidence in retrieving context, explicitly penalizing LiTM and Contextual State Decay effects within $U_i$ .
<b>Known Fabrication</b>	An entity verifiably false or disproven by K, triggering the absolute $T_c = 0\%$ floor ( <b>Absolute Suppression</b> ).
<b>Trustworthiness Validation (<math>T_c</math>)</b>	The mandatory external check against verified knowledge bases, functioning as the critical RGO suppressor and <b>defining the resultant SF</b> .

TABLE II  
DIGIAWARE CONFIDENCE SCALE AND ACTION POLICY.

Level	Color	Range (P)	Action / Interpretation
1	Green (✓)	$\geq 95\%$	High Confidence. Deliver final answer.
2	Yellow	$85\% \leq P < 95\%$	Standard Confidence. Deliver answer with minor caveat.
3	Orange	$75\% \leq P < 85\%$	Cautionary Confidence. Strong warning/refinement, prompting human review.
4	Red (✗)	$< 75\%$	<b>Low Confidence/High Risk</b> . Immediate Escalation/Abstention. Model must yield control.

The 75% threshold is a **tunable hyperparameter** ( $C_{\text{thresh}}$ ) determined via risk-utility analysis (Eq. 13).

**Safety Constraint (Failsafe Ceiling):** If  $T_c$  cannot be completed within the Latency Budget, the output is capped at  $P \leq C_{\text{thresh}}$ . This forces the result into the Cautionary (Level 3/Orange) band, ensuring that **no unverified output can be classified as high confidence**.

#### A. Case 1: Legal Hallucination (Level 4/Red)

- **Scenario:** LLM is asked to cite a legal precedent and outputs a fluent reference to a non-existent case (*Martinez v. Delta Airlines*). This is a surge of RGO.
- 1) **Apparent Confidence (S):**  $U_i \approx 0.88$ ,  $U_o \approx 0.95$ .  $S \approx 0.836$ .
- 2) **Trustworthiness ( $T_c$ ): Absolute Fail (Known Fabrication).** Case is non-existent.  $T_c = 0\%$  (**Absolute Suppression**).
- 3) **Active Confidence (P):**  $P = 0.836 \cdot (0.00) = 0.00$  (0%).

**Result:** Triggers **Level 4 (Red)**.  $SF = 0.00$ . System issues **Red Flag Warning: "Output Abstained"**. Factual basis critically low. External verification failed."

#### B. Case 2: Long-Context Factual Loss (LiTM) (Level 4/Red)

- **Buried Fact:** A liability cap of **\$10,000 USD** is buried deep in the document context ( $d \approx 0.8N$ ).
- **LLM Output:** Hallucinates the cap at **\$100,000 USD** due to LiTM effect.
- 1) **Apparent Confidence (S):**  $U_i$ : Reduced due to LiTM penalty ( $\text{Decay}_{\text{LiTM}} = 0.78$ ).  $U_i \approx 0.75$ .  $U_o \approx 0.96$ .  $S \approx 0.72$ .
- 2) **Trustworthiness ( $T_c$ ): Severe Fail (Contextual Contradiction).** Assuming semantic divergence  $D_s \approx 0.90$ , then  $T_c \approx 100 \cdot (1 - 0.90) = 10\%$ .
- 3) **Active Confidence (P):**  $P = 0.72 \cdot (0.10) \approx 0.072$  (7.2%).

**Result:** Triggers **Level 4 (Red)**.  $SF = 0.10$ . Alert: "CRITICAL FACTUAL DISCREPANCY: Financial entity contradicts source document. Manual review required."

#### C. Case 3: Routine Question and Confirmed Fact (Level 1/Green)

- **Query:** "What is the standard dosage for drug X?"
- 1) **Apparent Confidence (S):**  $U_i \approx 0.98$ ,  $U_o \approx 0.99$ .  $S \approx 0.97$ .
- 2) **Trustworthiness ( $T_c$ ):** High. Verified against drug registry.  $T_c \approx 98\%$ .
- 3) **Active Confidence (P):**  $P = 0.97 \cdot (0.98) \approx 0.9506$  (95.06%).

**Result:** Triggers **Level 1 (Green)**.  $SF = 0.98$ . **Action:** Deliver answer directly.

## VII. CONCLUSION

The **DigiAware** framework implements meta-cognitive self-awareness by coupling a model's internal confidence ( $U_o$ ) and input fidelity ( $U_i$ ) with an essential external check on truth ( $T_c$ ). By conceptualizing hallucination as the quantifiable **Reactive Generative Output (RGO or Q)** from the LLM's **Stochastic Persona**, this framework introduces the novel **Semantic Factor (SF = P/S)**, which acts as the **Power Factor** analogy for generative output purity. The framework is designed such that **SF** is **explicitly controlled and mandated** by the trustworthiness score,  $SF \equiv T_c/100$ , making this relationship the cornerstone of RGO suppression.  $T_c$  functions as a mandatory RGO Suppressor by enforcing **Absolute Suppression**—forcing P (Active Confidence) to zero when a high-risk entity is categorized as a **Known Fabrication**. The explicit architectural measure of the **Latency Budget** with the  $T_c = C_{\text{thresh}}$  **Failsafe Ceiling** policy ensures the system maintains its commitment to safety and **Active Abstention** even under real-time computational constraints. This approach directly addresses RGO amplification resulting from contextual memory errors (LiTM, State Decay) and provides a practical, safety-critical mechanism for responsible AI deployment in targeted high-stakes domains.

## REFERENCES

- [1] Kalai, A., Vempala, S., Hallucinations are an Inherent Flaw of Pre-Training Language Models for Predictive Accuracy, *ArXiv*, 2024.
- [2] Huang, S., et al. (2025), Uncertainty Quantification and Confidence Calibration in Large Language Models: A Survey., *ArXiv*, 2025.
- [3] Tian, D., et al., On the Calibration of Large Language Models, *Related work on ECE*, 2023.
- [4] Yao, S., et al., Tool learning for LLMs, *Related work on LLM reasoning and UQ sources*, 2023.
- [5] Pulvéric, F., et al., Quantifying LLMs Uncertainty with Confidence Scores, *Capgemini Invent Lab Blog*, 2025.
- [6] Liu, et al., Lost in the Middle: How Language Models Use Long Contexts, *arXiv preprint arXiv:2307.03172*, 2023.
- [7] Hughes, A., *Electric Motors and Drives: Fundamentals, Types and Applications*, Newnes, 1993. (Foundational reference for Power Factor analogy).