# Emotional Neural Networks: Teaching AI Not Just to Think, but to *Feel* Its Decisions

A new class of hybrid architecture combining learnable reasoning with learnable internal feelings to create more robust and introspective AI.

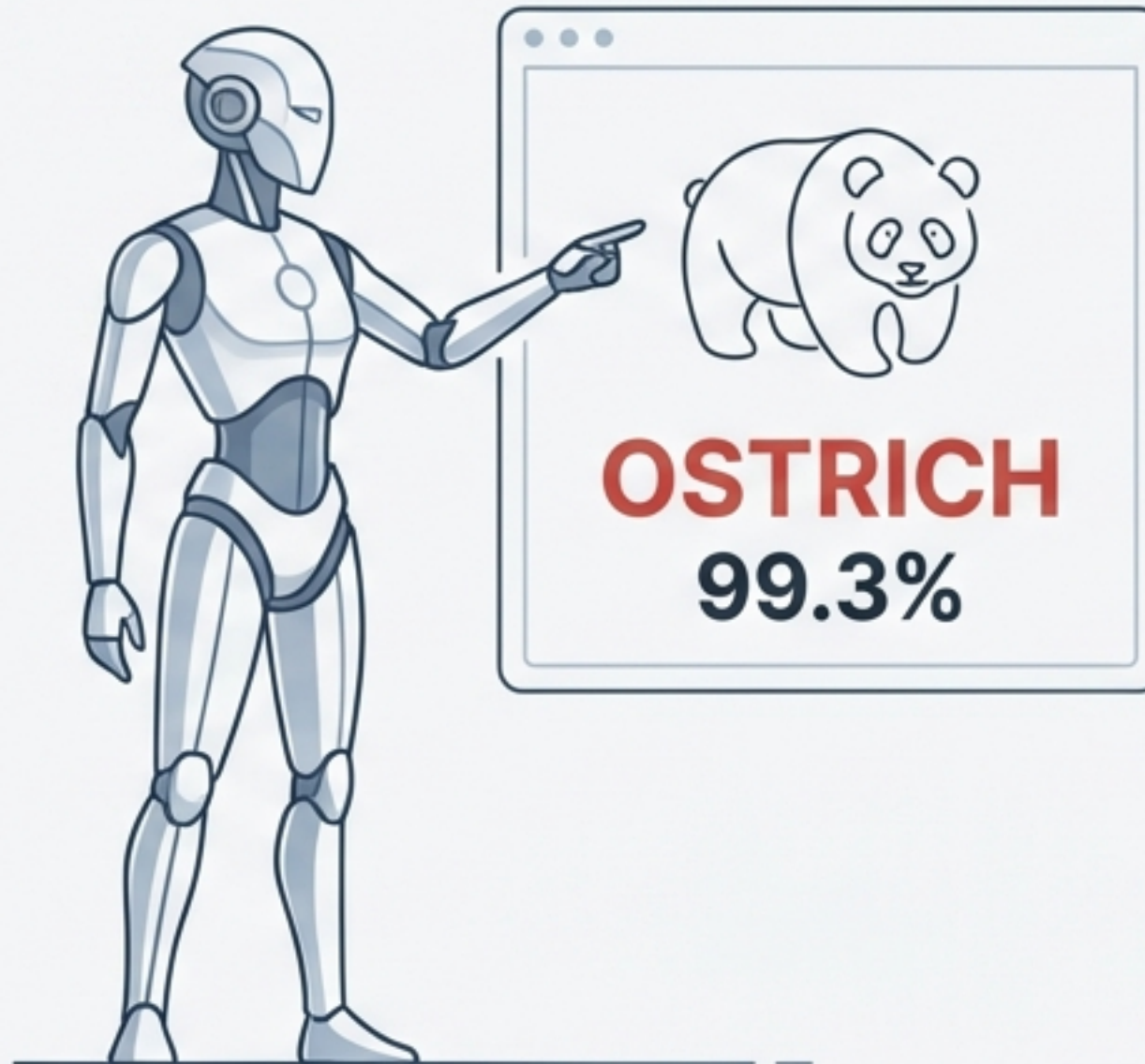# Today's AI: Brilliant Calculators Without Introspection

Standard Neural Networks excel at learning statistical regularities but lack internal awareness. They cannot distinguish between a decision that is *merely likely* and one that also *feels right or wrong*.

### Confidently Wrong

Can produce high-confidence predictions that are factually incorrect, with no internal mechanism to flag the error.

### Brittle Under Uncertainty

Performance degrades unpredictably when faced with novelty, ambiguity, or out-of-distribution data.

**OSTRICH**
**99.3%**

### Vulnerable to Adversarial Input

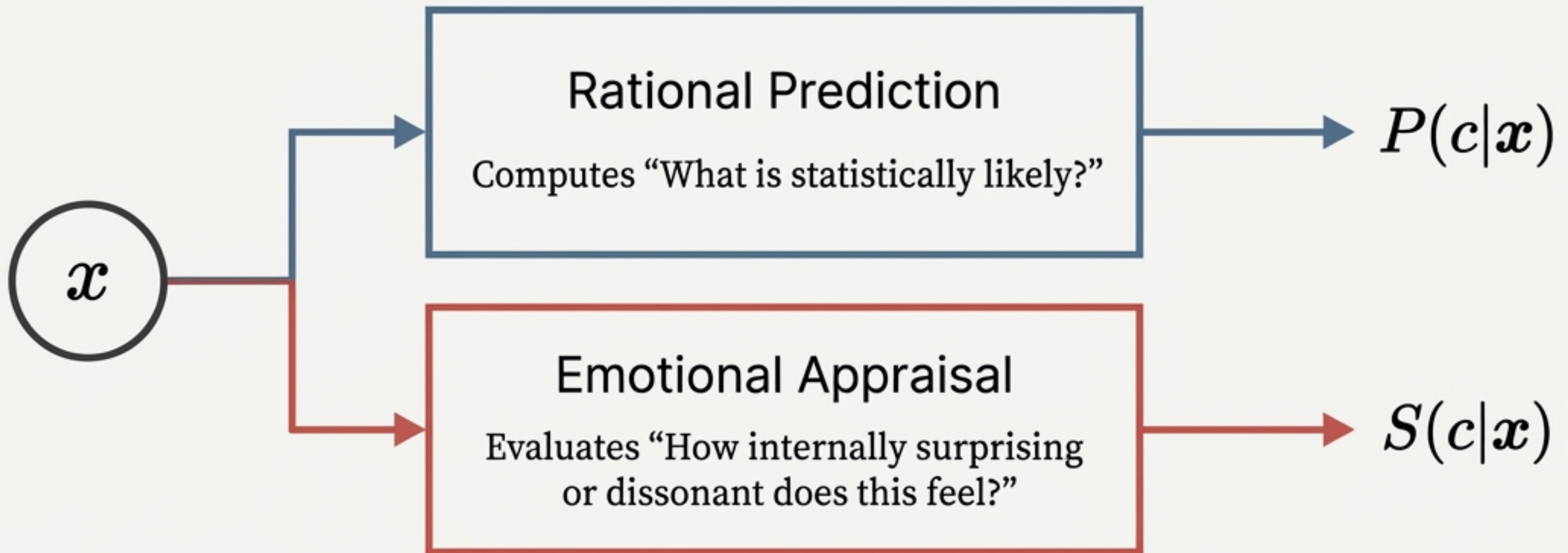Minor, imperceptible perturbations can cause catastrophic failures in classification.

### Lacks Self-Reflection

A standard NN computes probability but has no capacity to reflect on the internal consistency of its own decision-making process.

NotebookLM

# The Solution: A Hybrid of Reason and Feeling

Introducing Emotional Neural Networks (ENNs), an architecture inspired by dual-process theories of the mind. ENNs operate with two parallel, independently trained streams:

**Rational Prediction**

Computes "What is statistically likely?"

$$P(c|\boldsymbol{x})$$

$\boldsymbol{x}$

**Emotional Appraisal**

Evaluates "How internally surprising or dissonant does this feel?"
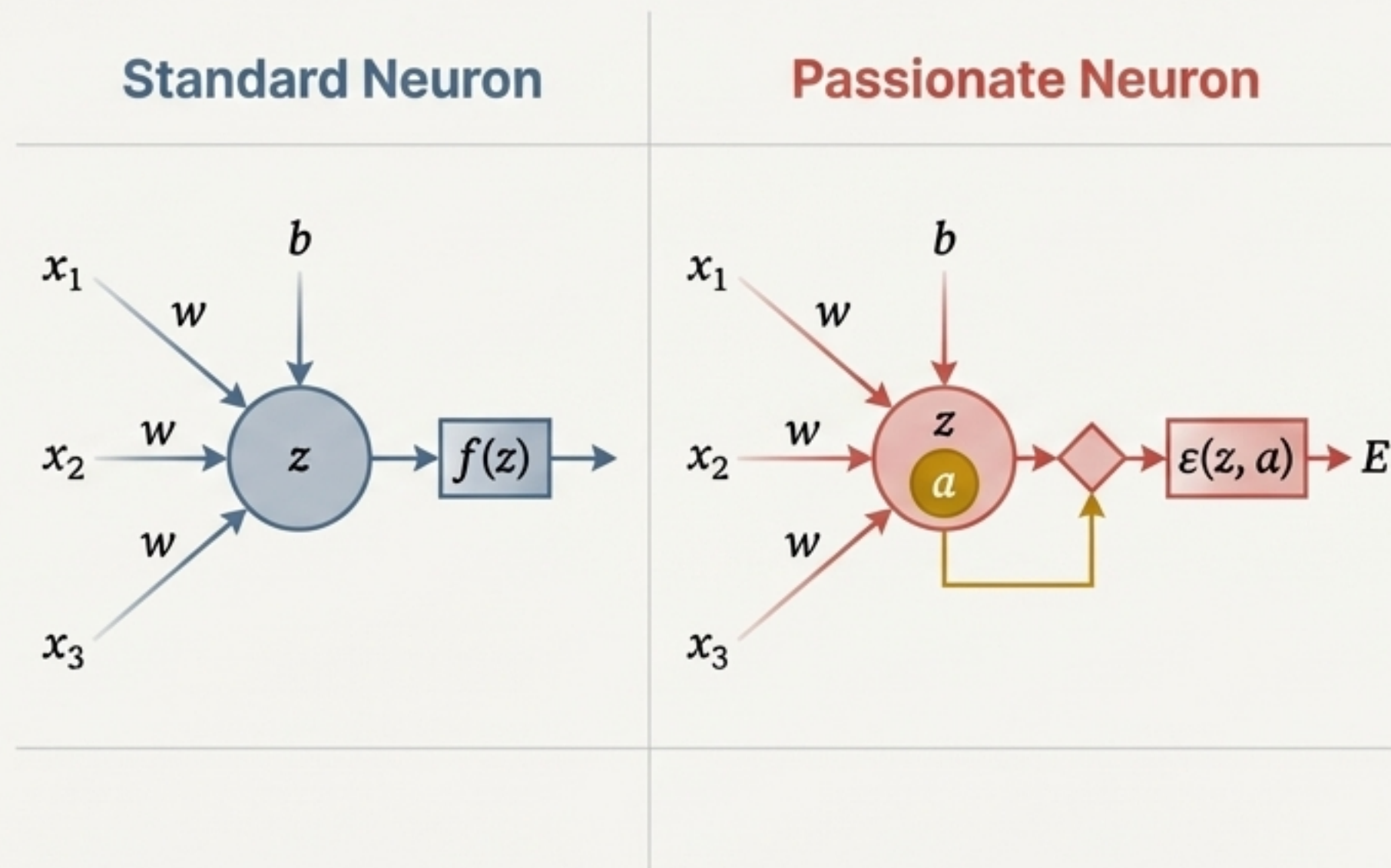
$$S(c|\boldsymbol{x})$$

# The Heart of the ENN: The Passionate Neuron.

The Passionate Neuron is a computational unit that doesn't just process input passively but reflects on its own internal state.

- **Belief (`z`):** The neuron's pre-activation value ($z = w \cdot x + b$), representing its initial appraisal of the input.

- **Expectation (`a`):** A trainable parameter representing the neuron's learned emotional baseline or norm for its belief.

- **Emotion (`E`):** The final activation output, quantifying the reaction to the mismatch between belief and expectation.

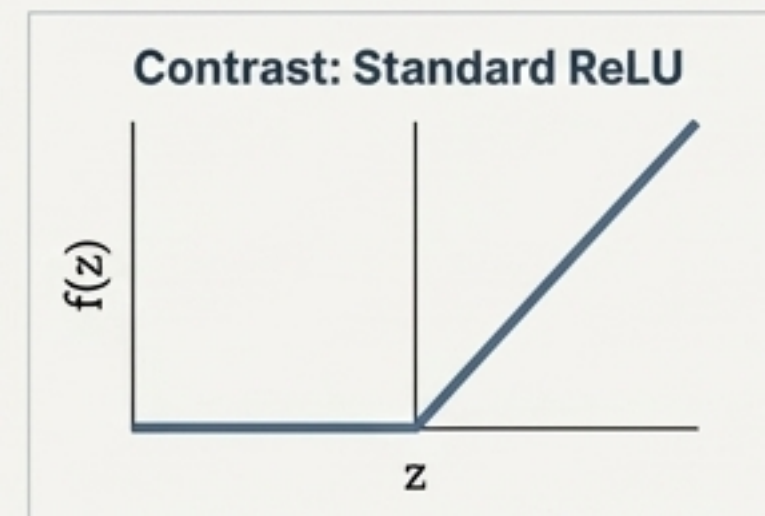**Standard Neuron**



**Passionate Neuron**

# The 'Anxiety' Activation: A Simple Measure of Surprise

$$E = \varepsilon(z, a) = |z - a|$$



This quantifies the emotional reaction to the mismatch between what the neuron believes (z) and what it expected (a).

A larger gap triggers a stronger 'emotional' signal.

It directly measures the 'distance' between belief and expectation.
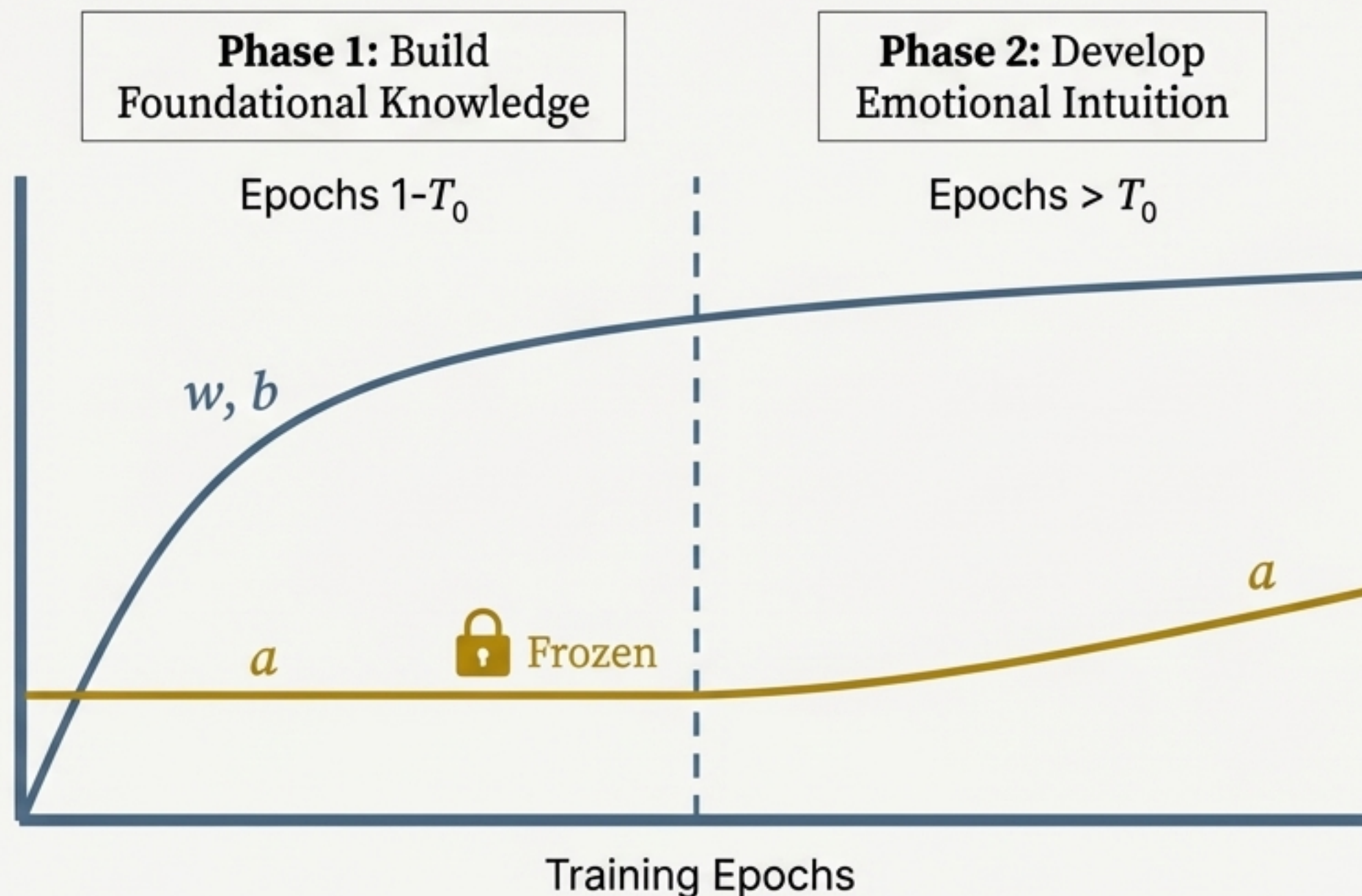
**Contrast: Standard ReLU**

- Piecewise linear: Ensures stable and predictable gradient behavior.

- Symmetric around a learnable anchor a: Unlike ReLU, which is fixed at 0, each neuron learns its own 'center of importance'.

- Computationally lightweight: Efficient for large-scale applications.

# Learning to Feel: Training the Emotional Anchor.

The expectation $a$ is a trainable parameter, but it evolves differently from the weights $w$ and biases $b$.

- **Weights & Biases ($w$, $b$):** Learn quickly, adapting to specific data patterns. This is analogous to acquiring skills or context-dependent knowledge.

- **Expectation ($a$):** Learns slowly, with a lower learning rate and less frequent updates. It represents a cumulative, stable model of what the neuron considers familiar. This is analogous to forming core values or stable priors.



**Phase 1:** Build Foundational Knowledge
Epochs $1$-$T_0$

**Phase 2:** Develop Emotional Intuition
Epochs $> T_0$

$w, b$

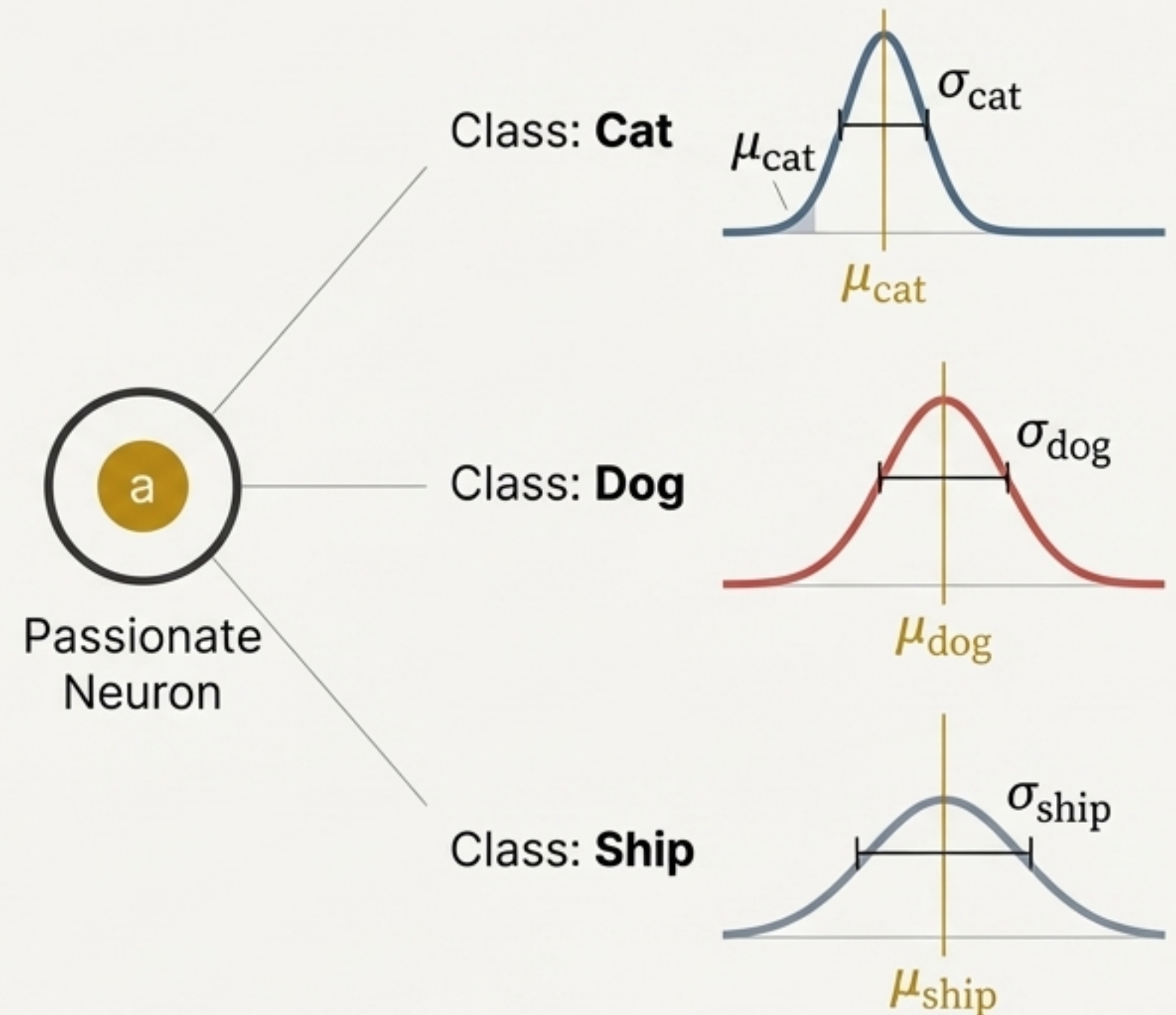$a$

🔒 Frozen

$a$

Training Epochs

# The Emotional Portfolio: A Neuron's Memory of Its Feelings.

After training, the ENN builds a statistical memory of its emotional reactions for each neuron on a per-class basis. These statistics are collected by running a dedicated "portfolio subset" of the data ($D_{portfolio}$) through the trained network.

**Mean Emotional Intensity ($\mu_i^c$):** "How strongly do I usually feel about class `c`?"

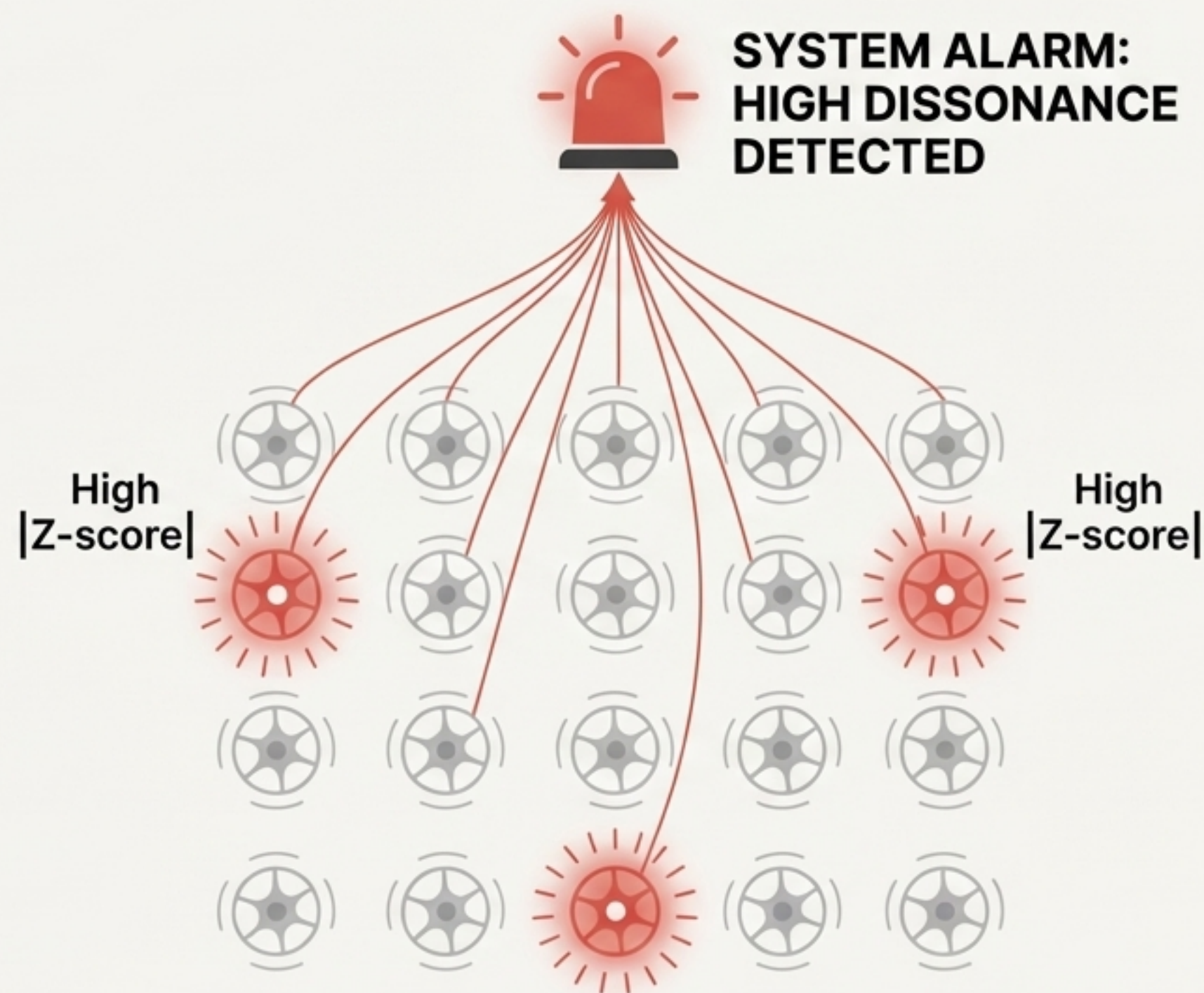**Emotional Volatility ($\sigma_i^c$):** "How stable are my feelings for class `c`?"

# From Local Surprise to Global Alarm

During inference, each neuron uses its emotional portfolio to calculate a Z-score for its current emotion relative to a class.

$$Z_i^c(x) = \frac{E_i(x) - \mu_i^c}{\sigma_i^c + \epsilon}$$

A high absolute Z-score is a **local alarm**, indicating the neuron's current reaction is atypical. Individual Z-scores are aggregated into a surprise vector, which is compared to the class's "emotional signature" to compute a final alarm score.



SYSTEM ALARM: HIGH DISSONANCE DETECTED
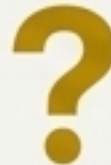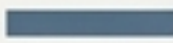
High |Z-score|

High |Z-score|

# Two Outputs Are Better Than One: Probability and Surprise.

An ENN produces two distinct, parallel outputs for any given input `x`:

- **Probability Vector `P(c|x)`**: The standard rational output. "This is 95% likely a cat."
- **Alarm Vector `S(c|x)`**: The introspective emotional output. "But my emotional alarm for 'cat' is at 82%."

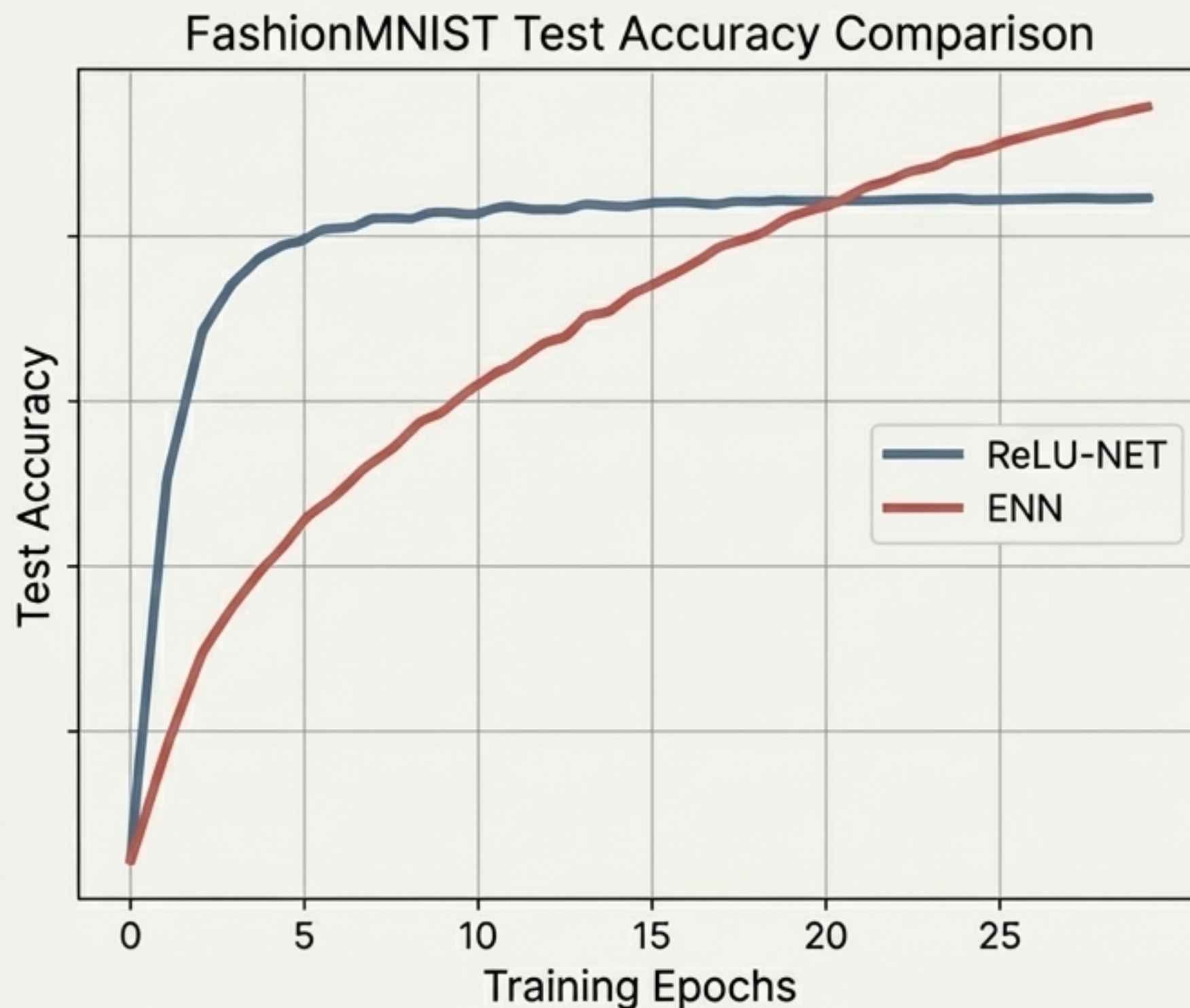|  | **Low Alarm**<br>(Emotionally Consistent) | **High Alarm**<br>(Emotionally Dissonant) |
|---|---|---|
| **High Probability** | **Confident & Aligned**<br>The ideal case. Prediction is statistically strong and feels right. ✔ | **Confident but Suspicious**<br>A red flag. "Likely, but something feels wrong." Potential adversarial attack or anomaly. ❓ |
| **Low Probability** | **Familiar Negative**<br>A clear rejection that aligns with expectations. — | **Novelty/Anomaly**<br>Strong rejection. The input is not recognized rationally or emotionally. ✖ |

NotebookLM

# Emotion Enhances Generalization.

**Experiment:** A minimal ENN (using only the anxiety activation) was compared against an identical ReLU-based network on the FashionMNIST dataset.

**Finding:** The ENN's learning was slower but more stable and consistent, eventually surpassing the baseline ReLU-NET in test accuracy.

**Observation:** "ReLU-NET typically exhibited faster initial progress... but tended to plateau early. In contrast, ENN... demonstrated slower demonstrated slower but more consistent improvement over time."



FashionMNIST Test Accuracy Comparison

# Inspired by How Intelligent Beings Actually Work

ENNs are not an arbitrary design but a synthesis of insights from multiple disciplines.
Emotions are not noise; they are structured reactions to cognitive discrepancies.



## Philosophy

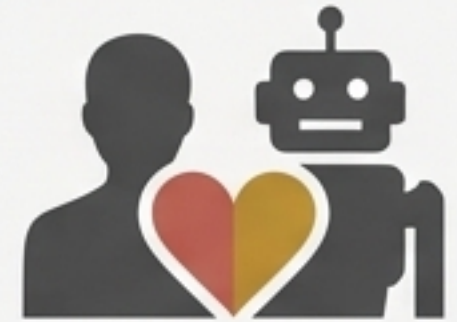Emotions as structured "judgments" arising from the evaluation of events relative to personal concerns.

## Psychology

Appraisal Theory, where emotional intensity scales with the gap between reality and expectation.

## Neuroscience

The brain as a prediction engine. The amygdala's role in responding to unexpected stimuli and prediction errors.

## Affective Computing

Emotion as integral to intelligence and rational decision-making, not a hindrance to it.

# A Step Towards Computational Proto-Consciousness.



**The Connection**

The ability to model its own internal state, compare it to learned expectations, and react to discrepancies is a foundational step towards self-monitoring systems.

**Definition**

ENNs implement a form of **computational proto-consciousness**: a pre-reflective capacity where the system is aware of mismatches between its internal model and the world. It is the "feeling of what happens" within the system.
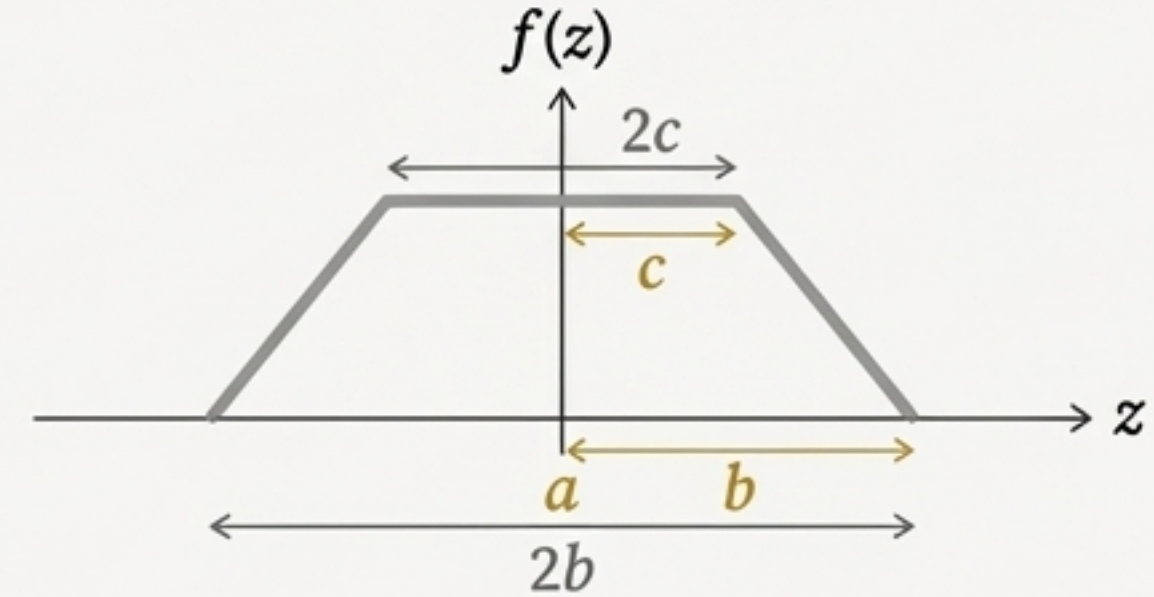
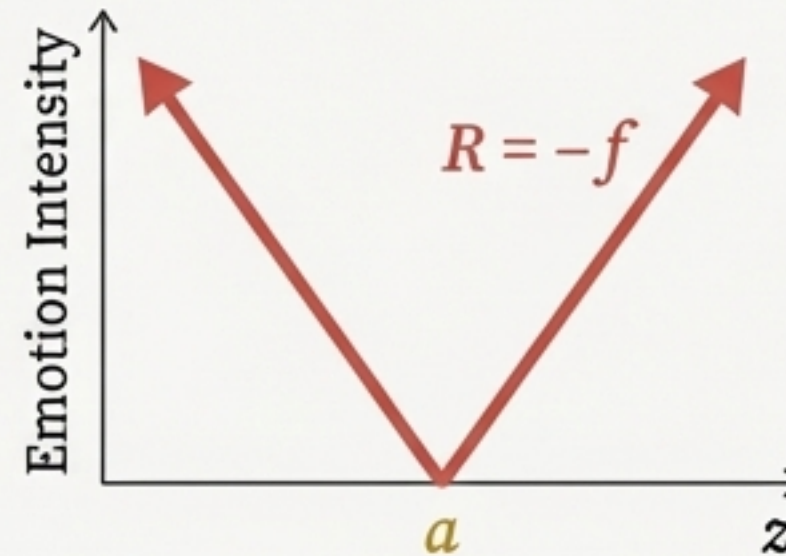# Beyond Anxiety: A Family of Emotional Activations

## General Framework & Logics

The anxiety function is a specific case of a generalized family of emotional activations based on a trapezoidal function $f(z; a, b, c)$.

- **Reflectors:** Amplify deviation from an anchor ($R = -f$). They model surprise, tension, and alertness. The $|z-a|$ anxiety function is an unbounded Linear Reflector.

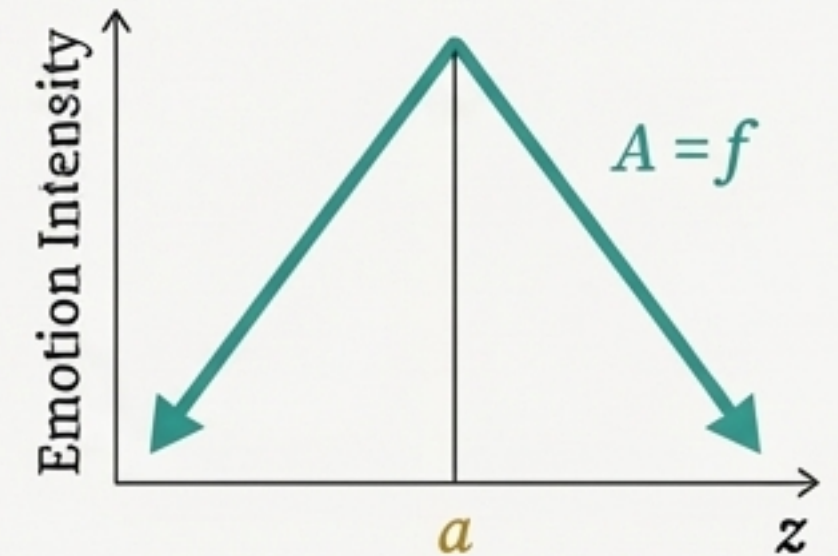- **Attractors:** Reward closeness to an anchor ($A = f$). They model satisfaction, comfort, and resonance.

## Generalized Trapezoidal Function $f(z; a, b, c)$



**Reflector** (Surprise)

$R = -f$

Emotion Intensity

$a$   $z$

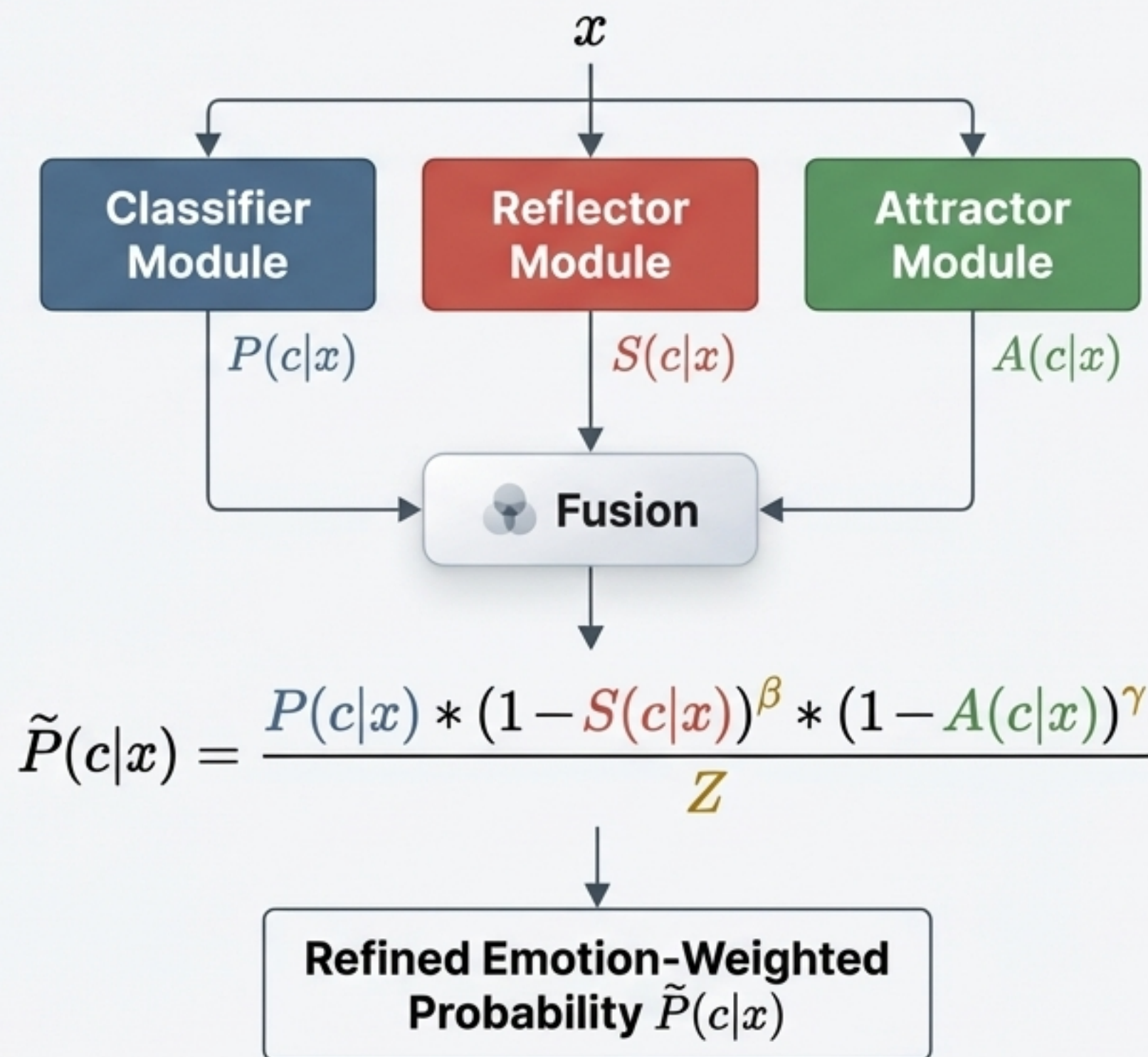**Attractor** (Satisfaction)

$A = f$

Emotion Intensity

$a$   $z$

# A Richer Inference: Fusing Probability, Surprise, and Satisfaction

A fully-realized ENN can produce three parallel distributions:

1. **Probability `$P(c|x)$**: From the rational stream.
2. **Surprise `$S(c|x)$**: A novelty score from Reflector neurons.
3. **Satisfaction Mismatch `$A(c|x)$**: A novelty score from Attractor neurons.

**Intuition:** A class is reinforced if it is probable AND emotionally typical (low surprise, low satisfaction mismatch). A class is penalized if it feels emotionally dissonant, even if it is statistically likely.

$x$

| Classifier Module | Reflector Module | Attractor Module |
|---|---|---|

$P(c|x)$ $\qquad$ $S(c|x)$ $\qquad$ $A(c|x)$

**Fusion**

$$\tilde{P}(c|x) = \frac{P(c|x) * (1 - S(c|x))^{\beta} * (1 - A(c|x))^{\gamma}}{Z}$$

**Refined Emotion-Weighted Probability $\tilde{P}(c|x)$**

# The Future of AI is Thinking *and* Feeling.

**Hybrid Architecture:** ENNs integrate rational prediction with an internal, learnable 'emotional' appraisal stream.

**Introspective Capabilities:** Emotional portfolios and dual-channel outputs enable self-monitoring, anomaly detection, and richer, more nuanced judgments.

**Powerful Core:** The 'Passionate Neuron' with its simple `|z-a|` anxiety activation is computationally efficient and improves model generalization.

**A New Paradigm:** This framework opens the door to more robust, interpretable, and affect-aware AI systems that can develop a sense of what is not just likely, but appropriate.

NotebookLM