# "Gold from Bronze": Stabilizing Synthetic Data in Neural Networks

Shivansh Sharma

*Undergraduate Student in Department of Emerging Technologies*
(Mahatma Gandhi Institute of Technology)
Hyderabad, India
shivanshsharma71231@gmail.com

*Abstract*—This paper presents a methodological framework for developing a robust, stateful "Emotional State Manager" using synthetic data generated by Large Language Models (LLMs). Addressing the scalability limitations of rule-based systems and the scarcity of physiological ground-truth data, we introduce a teacher-student pipeline using Gemini 2.5 Flash to generate a dataset of continuous hormone vectors. A compact Feed-Forward Network (FNN) is trained on this data, where we identify and resolve a critical "Dying ReLU" failure mode by substituting ELU activation. Further optimization via a custom Parametric Bounded Linear Activation (PBLA) and targeted edge-case augmentation yields a stable Test MAE of 5.25. The results demonstrate that small neural networks can learn to filter the volatility of synthetic teacher labels, achieving a "learned stability" suitable for real-time affective computing.

*Index Terms*—Affective Computing, Synthetic Data, Neural Networks, Large Language Models, PBLA.

## I. INTRODUCTION

Modeling human emotion in computational systems presents significant challenges, primarily due to the scarcity of labeled affective data and the inherently continuous and ambiguous nature of emotional states. Most existing production systems rely on rule-based emotional engines, which employ extensive sequences of conditional statements to encode hand-designed transitions between emotions. Such systems are inherently brittle, difficult to scale, and incapable of exhibiting emergent emotional behaviors.

The initial system developed was a rule-based engine heavily reliant on conditional branches. Attempts to expand its scope quickly highlighted the fundamental limitations of symbolic logic for continuous affective modeling. To overcome the constraints of brittle rule-based systems, the focus shifted to learning emotional state transitions directly from data.

However, real-world datasets containing physiological or hormone-level data are often scarce, proprietary, or subject to strict ethical constraints. This necessitated a novel approach: leveraging a large language model (LLM) as a synthetic teacher to generate plausible "hormone vectors" that represent the modeled internal emotional states.

This paper details the methodology for constructing this synthetic dataset, the subsequent training of a small neural network on the generated data, and the identification of several subtle failure modes, including dying Rectified Linear Units (ReLUs), noisy labels, activation saturation, and non-grounded volatility. The paper concludes by presenting the architectural and data-centric refinements employed to successfully resolve these issues.

## II. RELATED WORK

### A. Affective Computing & Dialogue Systems

Prior research in affective computing has predominantly focused on classification tasks using discrete emotion labels. Rashkin et al. (2019) introduced benchmarks for empathetic dialogue generation, relying on static labels like "happy" or "sad" assigned to individual text turns [7]. While effective for sentiment analysis, these models lack state persistence. Similarly, Gao et al. (2021) surveyed Dialogue State Tracking (DST) methods that track user goals (e.g., booking a flight) but noted a gap in tracking internal affective states over long horizons [2]. Our work addresses this limitation by modeling emotion not as a static label, but as a continuous, evolving hormone vector that persists across turns.

### B. Synthetic Data Distillation

The scarcity of physiological emotional data has led researchers to explore synthetic alternatives. Wang et al. (2024) demonstrated that Large Language Models (LLMs) can serve as effective data augmenters for low-resource NLP tasks, successfully distilling knowledge from large teacher models into smaller student networks [8]. We extend this methodology by employing Gemini 2.5 Flash [3] not just for augmentation, but as a primary "synthetic teacher" to simulate complex physiological dynamics that are otherwise unobservable.

### C. Neural Network Pathologies

Training regression models on synthetic data often introduces instability. Lu et al. (2019) analyzed the "Dying ReLU" problem, showing that standard ReLU activations often collapse in regression tasks where the output distribution has high variance or negative shifts [5]. Clevert et al. (2016) proposed the Exponential Linear Unit (ELU) to mitigate this by allowing negative values, pushing the mean activation closer to zero [1]. Our work validates these findings in the context of affective computing, demonstrating that ELU (and our proposed PBLA) prevents the zero-gradient collapse observed in baseline architectures.

## III. METHODOLOGY

### A. Dataset Generation

We utilized a two-stage data generation process. First, we employed a teacher-student framework where the Gemini 2.5 Flash API [3] was prompted to act as an expert emotional modeler. The model generated next-state hormone vectors based on randomly sampled current states and user intent logits. To address the lack of volatility in the initial synthetic data, we augmented the training set with 20 rule-based "edge case" samples (e.g., high stress combined with high happiness). The final dataset consisted of 1,000 samples.

### B. Model Architecture

The proposed architecture is a compact Feed-Forward Network (FNN) designed for real-time inference. The input layer accepts a 13-dimensional vector comprising the current 5 hormone values and 8 user intent logits. This is processed by a single hidden layer of 320 units. We selected the Exponential Linear Unit (ELU) [1] as the hidden activation function to eliminate the "Dying ReLU" problem observed in baseline experiments.

### C. Parametric Bounded Linear Activation (PBLA)

To manage output saturation, we introduced a custom activation function, PBLA. Unlike standard Sigmoid or Hard-Sigmoid functions which saturate at 0, causing "zero-bias collapse," PBLA learns a per-unit slope ($k$) and bias ($b$) constrained within fixed upper and lower bounds ($U$ and $L$ respectively). This allows the model to predict small-signal values (e.g., 1.7) without vanishing gradients, stabilizing the regression output against the volatility of the synthetic training labels.

$$PBLA(x) = clamp(k \cdot x + b, L, U) \tag{1}$$

### D. Training Procedure

The model was trained using the L1 (Mean Absolute Error) loss function, which offers robustness against outliers. Optimization was performed with the AdamW algorithm [4], employing a weight decay of 0.02. A batch size empirically determined to be optimal, ranging from 100 to 180, was used. The training was conducted over 4000 epochs, where convergence was visually confirmed and no overfitting was observed. The learning rate was decayed by a factor of 0.1 at 2000 epochs using a step decay schedule. Crucially, the train/test split was held static at $90/10$ across all experiments.

### E. Interpretation Layer (Physiological Mapping)

An essential component of the framework is a deterministic interpretation layer that maps the continuous 5-dimensional hormone vector to categorical emotion states. This layer, grounded in simplified physiological hypotheses, was designed to ensure that the model's numeric outputs translate into coherent and interpretable emotional behavior. The evaluation mapping rules are defined as follows:

- **Happy:** Modeled as reward-seeking (Serotonin, Dopamine) and gated by the absence of stress (Cortisol).
- **Sad:** Modeled as high stress (Cortisol) combined with low well-being (Serotonin).
- **Angry:** Modeled as a combination of threat (Adrenaline) and stress (Cortisol).
- **Caring:** Defined as a pro-social state driven by Oxytocin, suppressed by panic (Adrenaline).
- **Love:** Defined as a peak state requiring simultaneous bonding (Oxytocin) and well-being (Serotonin), making it highly fragile to any stressor (Cortisol or Adrenaline).

The constants embedded in these equations function as damping factors, empirically tuned to prevent easy saturation of the calculated emotion scores.

## IV. CHRONOLOGICAL PROGRESS & ABLATION FINDINGS

### A. Baseline Model Initialization

*Result: MAE (Train) $\approx 8.0$, MAE (Test) $\approx 9.0$*

The initial architecture utilized ReLU activations, which resulted in significant 'dying neuron' issues when compounded by the L1 loss function's sensitivity to outliers. Extensive hyperparameter optimization—including comprehensive sweeps of batch sizes (12-370), dropout rates, learning rates, and architectural configurations—failed to achieve convergence beyond this performance plateau.

### B. Activation Function Refinement (Phase 2)

The ReLU activation function was replaced with the Exponential Linear Unit (ELU) to mitigate the problem of gradient-dead regions. *Result: MAE (Train) $\approx 7.0$, MAE (Test) $\approx 7.0$* This modification marked the first successful stable convergence of the model.

### C. Data-Centric Refinement (Phase 3)

To address the instability observed in Phase 2, we implemented three targeted refinements. First, we streamlined the feature set by removing five engineered inputs (reducing dimensions from 18 to 13). Second, we augmented the training set with the 20 critical "edge-case" samples. Finally, we introduced a learning-rate decay schedule.

*Result: MAE (Train) $\approx 5.2$, MAE (Test) $\approx 6.0$*

This phase revealed the "Gold from Bronze" phenomenon, where the model began rejecting extreme, artificial spikes in the teacher labels, resulting in qualitatively smoother behavior despite higher numerical error against the volatile ground truth.

### D. Final Model Configuration (Phase 4)

The final configuration synthesized our architectural and data insights. We replaced the `relu6_scaled` output activation with our custom Parametric Bounded Linear Activation (PBLA). We switched the optimizer to AdamW and expanded the synthetic dataset to approximately 1,000 samples.

**Final Result:**
- MAE (Train): 5.60
- MAE (Test): 5.25

The observation that Test MAE is lower than Training MAE indicates highly effective regularization.

## V. Results

### A. Training Procedure

The progression of model performance across the four development phases is summarized in Table I. Approximately 80% of samples achieved $\leq 2$ MAE, while 5% of samples (high-volatility labels) inflated the total MAE by 0.7-1.0. The PBLA activation nearly eliminates zero-bias predictions.

TABLE I
TRAINING PROGRESS AND MAE REDUCTION

| Phase | Key Change | Train MAE | Test MAE |
|-------|-----------|-----------|----------|
| Baseline | ReLU | $\approx 8.0$ | $\approx 9.0$ |
| Phase 2 | ELU | $\approx 7.0$ | $\approx 7.0$ |
| Phase 3 | Data cleanup + LR decay | $\approx 5.2$ | $\approx 6.0$ |
| Phase 4 | PBLA + AdamW + Dataset | $\approx 5.60$ | $\approx 5.25$ |

### B. Qualitative Analysis of State Transitions

To validate the model's learned behavior, specific transition samples were analyzed to verify if the "Gold from Bronze" effect resulted in socially coherent interactions. Three distinct emergent behaviors were observed (Table II).

## VI. Discussion

### A. Learning in the Absence of Ground Truth

Real physiological signals are unobservable in this context. The model learns an approximation of emotional dynamics, not biology. Yet it captures consistent, human-interpretable transitions. The successful decoding of predicted vectors via the Interpretation Layer confirms that the learned manifold is structurally coherent.

### B. Rejecting Noise: The Gold from Bronze Principle

In several cases, teacher labels exhibit implausible hormone swings (e.g., oxytocin jumping from 90 to 15 in one step). The model correctly refuses to imitate these, predicting softer transitions instead. This increases MAE relative to the teacher but increases realism.

### C. Why This Model Works Better Than the Teacher

LLMs generate emotionally coherent but sometimes inconsistent hormone values. A small neural network, trained with strong regularization and successful architecture choices, becomes a stability filter.

### D. Lessons for Affective Computing

This work suggests a new workflow:

1) Use LLMs to generate approximate emotional data.
2) Train a small grounded model to refine that data.
3) Architect activations to prevent pathological failure modes.
4) Embrace smoothing of noisy supervision.

## VII. Conclusion

A compact and effective affective state-transition model has been demonstrated, trained entirely on synthetic hormone-inspired data created using a teacher LLM. Through systematic experimentation, architectural changes, data-centric iteration, and a novel output activation (PBLA), the Mean Absolute Error (MAE) was reduced from $\approx 9.0$ to $\approx 5.25$. Crucially, qualitative validation via a deterministic physiological mapping layer confirms that this numerical improvement corresponds to interpretable, socially intelligent behavior. The approach is suitable for real-time conversational agents, affective simulators, and reinforcement-learning agents requiring internal emotional states.

## VIII. Future Work

- **Multi-teacher prompting:** Using multiple teachers (Gemini, GPT-4, Claude) to reduce student model bias.
- **Recurrent/Attention architecture:** Incorporating GRUs or lightweight attention to capture higher-order temporal dynamics like emotional accumulation.
- **Long-horizon stability:** Testing multi-step rollouts with biological constraints to improve interpretability.

## References

[1] D. A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)," in *Proc. ICLR*, 2016.
[2] Q. Gao, L. He, and Z. Chen, "A Survey on Dialogue State Tracking," *arXiv:2105.02388*, 2021.
[3] Gemini Team, Google, "Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities," *arXiv:2507.06261*, 2025.
[4] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980*, 2014.
[5] L. Lu, Y. Shin, Y. Su, and G. E. Karniadakis, "Dying ReLU and Initialization: Theory and Numerical Examples," *arXiv:1903.06733*, 2019.
[6] A. Paszke et al., "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *NeurIPS*, 2019.
[7] H. Rashkin, E. M. Smith, M. Li, and Y. L. Boureau, "Towards Empathetic Open-domain Conversation Models: A New Benchmark and Dataset," in *Proc. ACL*, 2019.
[8] Y. Wang, J. Kaddour, and L. Winstanley, "Large Language Models as Data Augmenters," in *Proc. NAACL*, 2024.

TABLE II
QUALITATIVE FAILURE & SUCCESS MODES

| Scenario | Current State | User Input (Logits) | Pred. Next State | Interpretation |
|---|---|---|---|---|
| **De-escalation** | Angry (44.0) | Loving (94.1) | Happy (30.9) | **Success:** The model correctly identifies that high affection can resolve moderate anger, flipping the valence. |
| **Inertia (Resistance)** | Angry (94.5) | Thankful (74.9) | Angry (81.8) | **Success:** The model exhibits emotional inertia. A simple "thank you" is insufficient to resolve extreme rage immediately. |
| **Stability Filtering** | Happy (50.2) | Angry (96.0) | Angry (45.0) | **Success (Gold from Bronze):** While user input was extremely negative, the model predicted a moderate response rather than an extreme spike, rejecting teacher label volatility. |